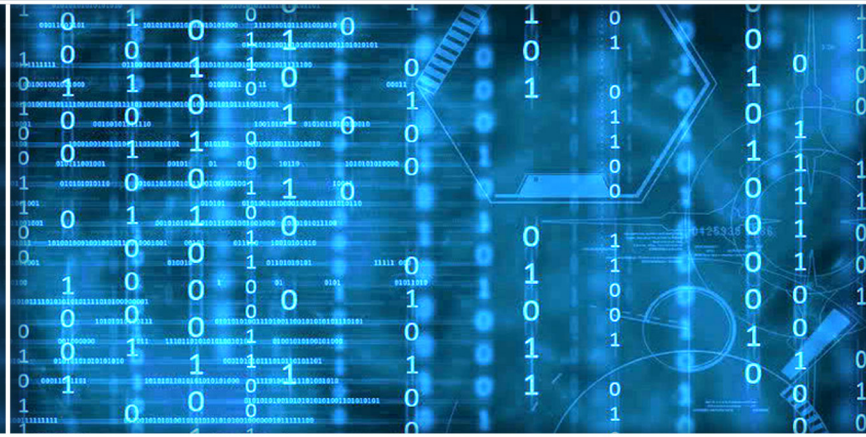


Volume 9 Issue 5

May 2018



ISSN 2156-5570(Online)

ISSN 2158-107X(Print)



Editorial Preface

From the Desk of Managing Editor...

It may be difficult to imagine that almost half a century ago we used computers far less sophisticated than current home desktop computers to put a man on the moon. In that 50 year span, the field of computer science has exploded.

Computer science has opened new avenues for thought and experimentation. What began as a way to simplify the calculation process has given birth to technology once only imagined by the human mind. The ability to communicate and share ideas even though collaborators are half a world away and exploration of not just the stars above but the internal workings of the human genome are some of the ways that this field has moved at an exponential pace.

At the International Journal of Advanced Computer Science and Applications it is our mission to provide an outlet for quality research. We want to promote universal access and opportunities for the international scientific community to share and disseminate scientific and technical information.

We believe in spreading knowledge of computer science and its applications to all classes of audiences. That is why we deliver up-to-date, authoritative coverage and offer open access of all our articles. Our archives have served as a place to provoke philosophical, theoretical, and empirical ideas from some of the finest minds in the field.

We utilize the talents and experience of editor and reviewers working at Universities and Institutions from around the world. We would like to express our gratitude to all authors, whose research results have been published in our journal, as well as our referees for their in-depth evaluations. Our high standards are maintained through a double blind review process.

We hope that this edition of IJACSA inspires and entices you to submit your own contributions in upcoming issues. Thank you for sharing wisdom.

Thank you for Sharing Wisdom!

Managing Editor
IJACSA
Volume 9 Issue 5 May 2018
ISSN 2156-5570 (Online)
ISSN 2158-107X (Print)
©2013 The Science and Information (SAI) Organization

Editorial Board

Editor-in-Chief

Dr. Kohei Arai - Saga University

Domains of Research: Technology Trends, Computer Vision, Decision Making, Information Retrieval, Networking, Simulation

Associate Editors

Chao-Tung Yang

Department of Computer Science, Tunghai University, Taiwan

Domain of Research: Software Engineering and Quality, High Performance Computing, Parallel and Distributed Computing, Parallel Computing

Elena SCUTELNICU

"Dunarea de Jos" University of Galati, Romania

Domain of Research: e-Learning, e-Learning Tools, Simulation

Krassen Stefanov

Professor at Sofia University St. Kliment Ohridski, Bulgaria

Domains of Research: e-Learning, Agents and Multi-agent Systems, Artificial Intelligence, Big Data, Cloud Computing, Data Retrieval and Data Mining, Distributed Systems, e-Learning Organisational Issues, e-Learning Tools, Educational Systems Design, Human Computer Interaction, Internet Security, Knowledge Engineering and Mining, Knowledge Representation, Ontology Engineering, Social Computing, Web-based Learning Communities, Wireless/ Mobile Applications

Maria-Angeles Grado-Caffaro

Scientific Consultant, Italy

Domain of Research: Electronics, Sensing and Sensor Networks

Mohd Helmy Abd Wahab

Universiti Tun Hussein Onn Malaysia

Domain of Research: Intelligent Systems, Data Mining, Databases

T. V. Prasad

Lingaya's University, India

Domain of Research: Intelligent Systems, Bioinformatics, Image Processing, Knowledge Representation, Natural Language Processing, Robotics

Reviewer Board Members

Aakash Ahmad Abbasi	Ali Asghar Pourhaji Kazem	Aris Skander Skander
Abbas Karimi	Ali Hamzeh	Arun D Kulkarni
Abbas M. Al-Ghaili	Ali Ismail Awad	Arun Kumar Singh
Abdelghni Lakehal	Ali Mustafa Qamar	Arvind K Sharma
Abdul Aziz Gill	Alicia Menchaca Valdez	Asadullah Shaikh
Abdul Hamid Mohamed Ragab	Altat Mukati	Asfa Praveen
Abdul Karim Assaf ABED	Aman Chadha	Ashok Matani
Abdul Razak	Amin Ahmad Shaqrah	Ashraf Hamdy Owis
Abdul Wahid Ansari	amine baina	ASIM TOKGOZ
Abdur Rashid Khan	Amir HAJJAM EL HASSANI	Asma Cherif
Abeer Mohamed ELkorany	Amirrudin Kamsin	Asoke Nath
ABRAHAM VARGHESE	Amitava Biswas	Athanasios Koutras
Adebayo Omotosho	Amjad Gawanmeh	Ayad Ghany Ismaeel
ADEMOLA ADESINA	Anand Nayyar	Ayman EL-SAYED
Aderemi A. Atayero	Anandhi Mohanraj Anu	Ayman Shehata Shehata
Adi A. Maaita	Andi Wahyu Rahardjo Emanuel	Ayoub BAHNASSE
Adnan Ahmad	Anews Samraj	Ayush Singhal
Adrian Nicolae Branga	Anirban Sarkar	Azam Moosavi
Ahmad A. Al-Tit	Anita Sofia V S	Babatunde Opeoluwa Akinkunmi
Ahmad A. Saifan	Anju Bhandari Gandhi	Bae Bossoufi
Ahmad Hoirul Basori	Anouar ABTOY	Balasubramanie Palanisamy
Ahmad Mousa Altamimi	Anshuman Sahu	BASANT KUMAR VERMA
Ahmed Boutejdar	Anthony Nosike Isizoh	Basem M. ElHalawany
Ahmed Nabih Zaki Rashed	Antonio Dourado	Basil Hamed
Ahmed S.A AL-Jumaily	Antonio Formisano	Basil M Hamed
Ahmed Z. Emam	ANUAR BIN MOHAMED KASSIM	Basim Almayahi
Ajantha Herath	Anuj Kumar Gupta	Bestoun S. Ahmed
Akram Belghith	Anuranjan misra	Bhanu Kaushik
Alaa F. Sheta	Appasami Govindasamy	Bhanu Prasad Pinnamaneni
Albert Alexander S	Arash Habibi Lashkari	Bharti Waman Gawali
Alci-nia Zita Sampaio	Aree Ali Mohammed	Bilian Song
Alexane Bouënard	Arfan Jaffar	Binod Kumar
ALI AMER ALWAN	ARINDAM SARKAR	Bogdan Belean

Bohumil Brtnik	Divya Kashyap	George D. Pecherle
Bouchaib CHERRADI	Djilali IDOUGHI	George Mastorakis
Brahim Raouyane	Dong-Han Ham	Georgios Galatas
Branko Karan	Dragana Becejski-Vujaklija	Gerard Dumancas
Bright Keswani	Duck Hee Lee	Ghalem Belalem Belalem
Brij Gupta	Duy-Huy NGUYEN	gherabi noredine
C Venkateswarlu Venkateswarlu Sonagiri	Ehsan Mohebi	Giacomo Veneri
Chanashekhhar Meshram	El Sayed A. Mahmoud	Giri Babu
Chao Wang	Elena Camossi	Goraksh Vithalrao Garje
Chao-Tung Yang	Elena SCUTELNICU	Govindarajulu Salendra
Charlie Obimbo	Elyes Maherzi	Grebenisan Gavril
Chee Hon Lew	Eric Tutu Tchao	Grigoras N. Gheorghe
CHERIF Med Adnen	Eui Chul Lee	Guandong Xu
Chien-Peng Ho	Evgeny Nikulchev	Gufran Ahmad Ansari
Chun-Kit (Ben) Ngan	Ezekiel Uzor OKIKE	Gunaseelan Devaraj
Ciprian Dobre	Fabio Mercorio	GYÖRÖDI ROBERT STEFAN
Constantin Filote	Fadi Safieddine	Hadj Hamma Tadjine
Constantin POPESCU	Fahim Akhter	Haewon Byeon
CORNELIA AURORA Gyorödi	Faizal Khan	Haibo Yu
Cosmina Ivan	FANGYONG HOU	Haiguang Chen
Cristina Turcu	Faris Al-Salem	Hamid Ali Abed AL-Asadi
Dai-Gyoung Kim	fazal wahab karam	Hamid Mukhtar
Daniel Filipe Albuquerque	Firkhan Ali Hamid Ali	Hamidullah Binol
Daniel Ioan Hunyadi	Fokrul Alom Mazarbhuiya	Hanan Elazhary
Daniela Elena Popescu	Fouad AYOUB	hanan habbi
Danijela Efnusheva	Francesco FP Perrotta	Hany Kamal Hassan
Dariusz Jakóbczak	Frank AYO Ibikunle	Harco Leslie Hendric SPITS WARNARS
Deepak Garg	Fu-Chien Kao	HARDEEP SINGH
Devena Prasad	G R Sinha	Hariharan Shanmugasundaram
DHAYA R	Gahangir Hossain	Harish Garg
Dheyaa Kadhim	Galya Nikolova Georgieva- Tsaneva	Hazem I. El Shekh Ahmed I. El Shekh Ahmed
Diaa Salama Dr	Gamil Abdel Azim	Heba Mahmoud Afify
Dimitris Chrysostomou	Ganesh Chandra Deka	Hela Mahersia
Dinesh Kumar Saini	Ganesh Chandra Sahoo	Hemalatha SenthilMahesh
Dipti Durgesh Patil	Gaurav Kumar	

Hesham G. Ibrahim	John P Sahlin	LATHA RAJAGOPAL
Hikmat Ullah Khan	JOHN S MANOHAR	Lazar Vojislav Stošić
Himanshu Aggarwal	JOSE LUIS PASTRANA	Le Li
Hongda Mao	José Santos Reyes	Leanos A Maglaras
Hossam Faris	Jui-Pin Yang	Leon Andretti Abdillah
Huda K. Kadhim AL-Jobori	Jungu J Choi	Lijian Sun
Hui Li	Jyoti Chaudhary	Liming Luke Chen
Hüseyin Oktay ERKOL	Jyoti Gautam	Ljubica B. Kazi
Ibrahim Adepoju Adeyanju	K V.L.N.Acharyulu	Ljubomir Jerinic
Ibrahim Missaoui	Ka-Chun Wong	Lokesh Kumar Sharma
Ikvinderpal Singh	Kamatchi R	Long Chen
Ilayaraja Muthalagu	Kamran Kowsari	M A Rabbani
Imad Zeroual	KANNADHASAN SURIYAN	M. Reza Mashinchii
Imed JABRI	KARTHIK MURUGESAN	M. Tariq Banday
Imran Ali Chaudhry	KASHIF MUNIR	Madiah Mohd Saudi
Imran Memon	Kashif Nisar	madjid khalilian
IRFAN AHMED	Kato Mivule	Mahdi H. Miraz
ISMAIL YUSUF	Kayhan Zrar Ghafoor	Mahmoud M Abd Ellatif
iss EL OUADGHIRI	Kennedy Chinedu Okafor	Mahtab Jahanbani Fard
Iwan Setyawan	KHAIRULLAH KHAN KHAN	Majharoddin Kazi Kazi
Jabar H Yousif	Khaled Loukhaoukha	majzoob kamal aldein omer
Jacek M. Czerniak	Khalid Mahmood	Malack Omae Oteri
Jafar Ahmad Alzubi	Khalid Nazim Sattar Abdul	Malik Muhammad Saad Missen
Jai Singh W	Khin Wee Lai	Mallikarjuna Reddy Doodipala
JAMAIAH HAJI YAHAYA	Khurram Khurshid	Man Fung LO
James Patrick Henry Coleman	KIRAN SREE POKKULURI	Manas deep
Jamil Abdulhamid Mohammed Saif	KITIMAPORN CHOOCHOTE	Manisha Gupta
Jatinderkumar Ramdass Saini	Kohei Arai	Manju Kaushik
Javed Anjum Sheikh	Kottakkaran Sooppy Nisar	Manmeet Mahinderjit Singh
Jayapandian N	kouki Mohamed	Manoharan P.S.
Jayaram M A	Krasimir Yankov Yordzhev	Manoj Manoj Wadhwa
Jerwinprabu A	Krassen Stefanov Stefanov	Manpreet Singh Manna
Ji Zhu	Krishna Kishore K V	Manuj Darbari
Jia Uddin Jia	Krishna Prasad Miyapuram	Marcellin Julius Antonio Nkenlifack
Jim Jing-Yan Wang	Labib Francis Gergis	Marek Reformat
	Lalit Garg	Maria-Angeles Grado-Caffaro

Marwan Alseid	Mohammed Shamim Kaiser	Naseer Ali Alquraishi
Mazin S. Al-Hakeem	Mohammed Tawfik Hussein	Nasrollah Pakniat
Md Ruhul Islam	Mohd Ashraf Ahmad	Natarajan Subramanyam
Md. Al-Amin Bhuiyan	Mohd Helmy Abd Wahab	Natheer Gharaibeh
Mehdi Bahrami	Mokhtar Beldjehem	Nayden V. Nenkov
Mehdi Neshat	Mona Elshinawy	Nazeeh Ghatasheh
Messaouda AZZOUZI	Monir Kaid	Nazeeruddin Mohammad
Milena Bogdanovic	Mostafa Mostafa Ezziyyani	Neeraj Kumar Tiwari
Miriampally Venkata Raghavendra	Mouhammad sharari sharari alkasassbeh	NEERAJ SHUKLA
Mirjana Popovic	Mounir Hemam	Nestor Velasco-Bermeo
Miroslav Baca	Mourad Amad	Nguyen Thanh Binh
Moamin Mahmoud	Mudasir Manzoor Kirmani	Nidhi Arora
Moeiz Miraoui	Mueen Uddin	NILAMADHAB MISHRA
Mohamed AbdelNasser	Muhammad Adnan Khan	Nilanjan Dey
Mohamed Mahmoud	Muhammad Abdul Rehman	Ning Cai
Mohamed Salah SALHI	Muhammad Asif Khan	Niraj Singhal
Mohamed A. El-Sayed	Muhammad Hafidz Fazli Bin Md Fauadi	Nithyanandam Subramanian
Mohamed Abdel Fatah Ashabrawy	Muhammad Naeem	Nizamud Din
Mohamed Ali Mahjoub	Muhammad Saeed	Noura Aknin
Mohamed Eldosoky	Muniba Memon	Obaida M. Al-Hazaimeh
Mohamed Hassan Saad Kaloup	MUNTASIR AL-ASFOOR	Olawande Justine Daramola
Mohamed Najeh LAKHOUA	Murphy Choy	Oliviu Matei
Mohamed SOLTANE Mohamed	Murthy Sree Rama Chandra Dasika	Om Prakash Sangwan
Mohammad Abdul Qayum	MUSLIHAH WOOK	Omaima Nazar Al-Allaf
Mohammad Ali Badamchizadeh	Mustapha OUJAOURA	Omar A. Alzubi
Mohammad Azzeh	MUTHUKUMAR S SUBRAMANYAM	Omar S. Gómez
Mohammad H. Alomari	N.Ch. Sriman Narayana Iyengar	Osama Ali Awad
Mohammad Haghighat	Nadeem Akhtar	Osama Omer
Mohammad Jannati	nafiul alam siddique	Ouchtati Salim
Mohammad Zarour	Nagy Ramadan Darwish	Ousmane THIARE
Mohammed Abdulhameed Al- shabi	Najeed Ahmed Khan	P.V. Praveen Sundar
Mohammed A. Akour	Najib A. Kofahi	Paresh V Virparia
Mohammed Ali Hussain	Namrata Dhanda	Parminder Singh Kang
Mohammed Sadgal	Nan Wang	PAUL CELICOURT
		Peng Xia
		Ping Zhang

Piyush Kumar Pareek	Reza Fazel-Rezai	Senol Piskin
Poonam Garg	Reza Ghasemy Yaghin Dr Reza Ghasemy Yaghin	SENTHIL P Prof
Prabhat K Mahanti	Riaz Ul-Amin	Sérgio André Ferreira
PRASUN CHAKRABARTI	Ricardo Ângelo Rosa Vardasca	Seyed Hamidreza Mohades Kasaei
Praveen Kumar	Ritaban Dutta	Shadi Mahmoud Atalla
PRISCILLA RAJADURAI	Rodica Doina Zmaranda	Shafiqul Abidin
PROF DURGA PRASAD SHARMA (PHD)	Rohini Ravi	Shahab Shamshirband
Purwanto Purwanto	Rohit Raja	Shahanawaj Ahamad
Qaisar Abbas	Roopali Garg	Shaidah Jusoh
Qifeng Qiao	roslina ibrahim	Shaiful Bakri Ismail
Rachid Saadane	Ruchika Malhotra	Shailesh Kumar
Radwan R. Tahboub	Rutvij H. Jhaveri	Shakir Gayour Khan
raed Kanaan	SAADI Slami	Shashi Dahiya
Raghuraj Singh	Sachin Kumar Agrawal	Shawki A. Al-Dubae
Rahul Malik	Sagarmay Deb	Sheeraz Ahmed Dr.
Raja Ramachandran	Sahar Abd EL_RAhman Ismail	Sheikh Ziauddin
raja sarath kumar boddu	Said Ghoniemy	Sherif E. Hussein
Rajesh Kumar	Said Jadid Abdulkadir	Shishir Kumar
Rakesh Chandra Balabantaray	Sajal Bhatia	SHOBA MOHAN
Rakesh Kumar Dr.	Saman Hina	Shriniwas Vasantrao Chavan
Ramadan Elaiess	SAMSON OLUWASEUN FADIYA	Shriram K Vasudevan
Ramani Kannan	Sanam Shahla Rizvi	Siddeeq Ameen
RAMESH MUTHUSAMY	Sandeep R Reddivari	Siddhartha Jonnalagadda
RAMESH VAMANAN	Sangeetha SKB	Sim-Hui Tee
Rana Khudhair Abbas Ahmed	Sanskruti V Patel	Simon L. R. Vrhovec
Rashad Abdullah Al-Jawfi	Santosh Kumar	Simon Uzezi Ewedafe
Rashid Sheikh	Sasan Adibi	Siniša Opic
Ratnesh Litoriya	Sattar Bader Sadkhan	Sivakumar Poruran
Ravi Kiran Varma P	Satyena Prasad Singh	sivaranjani reddy
Ravi Prakash	Sebastian Marius Rosu	Slim BEN SAOUD
RAVINA CHANGALA	Secui Dinu Calin	Sobhan Roshani
Ravisankar Hari	Seema Shah	Sofien Mhatli
Rawya Y. Rizk	Seifedine Nimer Kadry	sofyan Mohammad Hayajneh
Rayed AlGhamdi	Selem Charfi	Sohail Jabbar
Reshmy Krishnan	SENGOTTUVELAN P	Sri Devi Ravana

Sudarson Jena	Taskeed Jabid	Wenzhao Zhang
Sudipta Roy	Tasneem Bano Rehman	Wichian Sittiprapaporn
Suhail Sami Owais Sami Owais Owais	thabet Mohamed slimani	Xi Zhang
Suhas J Manangi	Totok R. Biyanto	Xiao Zhang
SUKUMAR SENTHILKUMAR	Touati Youcef	Xiaojing Xiang
Süleyman Eken	Tran Xuan Sang	Xiaolong Wang
Sumazly Sulaiman	TSUNG-CHUAN MA	Xunchao Hu
Sumit Goyal	Tsvetanka Georgieva-Trifonova	Y Srinivas
Sunil Phulre	Uchekukwu Awada	Yanping Huang
Suparerk Janjarasjitt	Udai Pratap Rao	Yao-Chin Wang
Suresh Sankaranarayanan	Urmila N Shrawankar	Yasser M. Alginahi
Surya Narayan Panda	V Baby Deepa	Yaxin Bi
Susarla Venkata Ananta Rama Sastry	Vaidas Giedrimas	Yi Fei Wang
Suseendran G	Vaka MOHAN	YI GU
Suxing Liu	Venkata Raghavendran Chaluvadi	Yihong Yuan
Syed Asif Ali	VENKATESH JAGANATHAN	Yilun Shang
T C.Manjunath	Vijay Bhaskar Semwal	Yu Qi
T V Narayana rao Rao	Vijayarani Mohan S	Zacchaeus Oni Omogbadegun
T. V. Prasad	Vijendra Singh	Zaffar Ahmed Shaikh
Taghi Javdani Gandomani	Vinayak K Bairagi	Zairi Ismael Rizman
Taiwo Ayodele	VINCE PAUL A	Zarul Fitri Zaaba
Talal Bonny	Visara Urovi	Zeki Yetgin
Tamara Zhukabayeva	Vishnu Narayan Mishra	Zenzo Polite Ncube
Taner Tuncer	Vitus S.W. Lam	ZHENGYU YANG
Tanvi Banerjee	VNR SAIKRISHNA K	Zhigang Yin
Tanweer Alam	Voon Ching Khoo	Zhihan Lv
Tanzila Saba	VUDA SREENIVASARAO	Zhixin Chen
TAOUFIK SALEM SAIDANI	Wali Khan Mashwani	Zia Ur Rahman Zia
Tarek Fouad Gharib	Wei Wei	Ziyue Xu
tarig ahmed	Wei Zhong	Zlatko Stapic
	Wenbin Chen	Zne-Jung Lee
		Zuraini Ismail

CONTENTS

Paper 1: Cardiocotographic Diagnosis of Fetal Health based on Multiclass Morphologic Pattern Predictions using Deep Learning Classification

Authors: Julia H. Miao, Kathleen H. Miao

PAGE 1 – 11

Paper 2: Flow-Length Aware Cache Replacement Policy for Packet Processing Cache

Authors: Hayato Yamaki

PAGE 12 – 20

Paper 3: Fuzzy Logic-Controlled 6-DOF Robotic Arm Color-based Sorter with Machine Vision Feedback

Authors: Alexander C. Abad, Dino Dominic Ligutan, Elmer P. Dadios, Levin Jaeron S. Cruz, Michael Carlo D.P. Del Rosario, Jho Nathan Singh Kudhal

PAGE 21 – 31

Paper 4: Mixed Profile Method of Speed and Location for Robotic Arms Motion used for Precise Positioning

Authors: Liliانا Marilena Matica, Cornelia Győrödi, Helga Silaghi, Andrei Silaghi

PAGE 32 – 36

Paper 5: A Novel E-Mail Network Evolution Model based on user Information

Authors: Lejun ZHANG, Tongxin ZHOU, Chunhui ZHAO, Zilong JIN

PAGE 37 – 45

Paper 6: Design of Traffic Flow Simulation System to Minimize Intersection Waiting Time

Authors: Jang, Seung-Ju

PAGE 46 – 50

Paper 7: Application of the Hierarchy Analysis Method at the Foodstuffs Quality Evaluation

Authors: Nikitina Marina Aleksandrovna, Nikitin Igor Alekseevich, Semenkinа Natalya Gennadievna, Zavalishin Igor Vladimirovich, Goncharov Andrey Vitalievich

PAGE 51 – 59

Paper 8: 3D Visualization of Sentiment Measures and Sentiment Classification using Combined Classifier for Customer Product Reviews

Authors: Siddhaling Urologin, Sunil Thomas

PAGE 60 – 68

Paper 9: A Novel Energy Efficient Mobility Aware MAC Protocol for Wireless Sensor Networks

Authors: Zain ul Abidin Jaffri, Asif Kabir, Gohar Rehman Chughtai, S. Sabahat H. Bukhari, Muhammad Arshad Shehzad Hassan

PAGE 69 – 74

Paper 10: Routing Optimization in WBAN using Bees Algorithm for Overcrowded Hajj Environment

Authors: Ghassan Ahmed Ali, Shah Murtaza Rashid Al Masud

PAGE 75 – 79

Paper 11: An Opportunistic Dissemination Protocol for VANETs

Authors: Amina SEDJELMACI, Fedoua DIDI, Ahmed ABDUL RAHUMAN

PAGE 80 – 87

Paper 12: Student Facial Authentication Model based on OpenCV's Object Detection Method and QR Code for Zambian Higher Institutions of Learning

Authors: Lubasi Kakwete Musambo, Jackson Phiri

PAGE 88 – 94

Paper 13: BLOT: A Novel Phase Privacy Preserving Framework for Location-Based Services

Authors: Abdullah Albelaihy, Jonathan Cazalas, Vijey Thayananthan

PAGE 95 – 104

Paper 14: Development of Mobile-Interfaced Machine Learning-Based Predictive Models for Improving Students' Performance in Programming Courses

Authors: Fagbola Temitayo Matthew, Adeyanju Ibrahim Adepoju, Oloyede Ayodele, Obe Olumide, Olaniyan Olatayo, Esan Adebimpe, Omodunbi Bolaji, Egbetola Funmilola

PAGE 105 – 115

Paper 15: New Techniques to Enhance Data Deduplication using Content based-TTDD Chunking Algorithm

Authors: Hala AbdulSalam Jasim, Assmaa A. Fahad

PAGE 116 – 121

Paper 16: Security Improvement in Elliptic Curve Cryptography

Authors: Kawther Esaa Abdullah, Nada Hussein M. Ali

PAGE 122 – 131

Paper 17: Classification of Affective States via EEG and Deep Learning

Authors: Jason Teo, Lin Hou Chew, Jia Tian Chia, James Mountstephens

PAGE 132 – 142

Paper 18: Rainfall Prediction using Data Mining Techniques: A Systematic Literature Review

Authors: Shabib Aftab, Munir Ahmad, Noureen Hameed, Muhammad Salman Bashir, Ifikhar Ali, Zahid Nawaz

PAGE 143 – 150

Paper 19: An Intelligent Bio-Inspired Algorithm for the Faculty Scheduling Problem

Authors: Sarah Al-Negheimish, Fai Alnuhait, Hawazen Albrahim, Sarah Al-Mogherah, Maha Alrajhi, Manar Hosny

PAGE 151 – 159

Paper 20: A Lightweight Multi-Message and Multi-Receiver Heterogeneous Hybrid Signcryption Scheme based on Hyper Elliptic Curve

Authors: Abid ur Rahman, Insaf Ullah, Muhammad Naeem, Rehan Anwar, Noor-ul-Amin, Hizbullah Khattak, Sultan Ullah

PAGE 160 – 167

Paper 21: A Chatbot for Automatic Processing of Learner Concerns in an Online Learning Platform

Authors: Mamadou BAKOUAN, Beman Hamidja KAMAGATE, Tiemoman KONE, Souleymane OUMTANAGA, Michel BABRI

PAGE 168 – 176

Paper 22: Formalization of Behavior Change Theories to Accomplish a Health Behavior

Authors: Adnan Manzoor, Imtiaz Ali Halepoto, Sohail Khokhar, Nazar Hussain Phulpoto, Engr. Muhammad Sulleman Memon

PAGE 177 – 182

Paper 23: Implementation of Winnowing Algorithm with Dictionary English-Indonesia Technique to Detect Plagiarism

Authors: Anton Yudhana, Sunardi, Iif Alfiatul Mukaromah

PAGE 183 – 189

Paper 24: Multi-Stage Algorithms for Solving a Generalized Capacitated P-median Location Problem

Authors: Mohammed EL AMRANI, Youssef BENADADA

PAGE 190 – 196

Paper 25: Communicator for Hearing-Impaired Persons using Pakistan Sign Language (PSL)

Authors: Muhammad Wasim, Adnan Ahmed Siddiqui, Abdulbasit Shaikh, Lubaid Ahmed, Syed Faisal Ali, Fauzan Saeed

PAGE 197 – 202

Paper 26: TPACK Adaptation among Faculty Members of Education and ICT Departments in University of Sindh, Pakistan

Authors: Saira Soomro, Arjumand Bano Soomro, Najma Imtiaz Ali, Tariq Bhatti, Nazish Basir, Nazia Parveen Gill

PAGE 203 – 209

Paper 27: Quality Aspects of Continuous Delivery in Practice

Authors: Maryam Shahzeydi, Taghi Javdani Gandomani, Rasool Sadeghi

PAGE 210 – 212

Paper 28: Integration of Heterogeneous Requirements using Ontologies

Authors: Ahmad Mustafa, Wan M.N. Wan-Kadir, Noraini Ibrahim, Muhammad Arif Shah, Muhammad Younas

PAGE 213 – 218

Paper 29: Effect of Service Broker Policies and Load Balancing Algorithms on the Performance of Large Scale Internet Applications in Cloud Datacenters

Authors: Ali Mefteh, Ahmed E. Youssef, Mohammad Zakariah

PAGE 219 – 227

Paper 30: Multiple Trips Pattern Mining

Authors: Riaz Ahmed Shaikh, Kamelsh Kumar, Razaqat Hussain Arain, Hidayatullah Shaikh, Imran Memon, Safdar Ali Shah

PAGE 228 – 232

Paper 31: E-shape Multiband Patch Antenna for 4G, C-band and S-band Applications

Authors: Mehr-e-Munir, Khalid Mahmood, Saad Hassan Kiani

PAGE 233 – 237

Paper 32: An Optimized Inset Feed Circular Cross Strip Antenna Design for C-Band Satellite Links

Authors: Faisal Ahmed Dahri, Riaz A. Soomro, Sajjad Ali Memon, Zeeshan Memon, Majid Hussain Memon

PAGE 238 – 242

Paper 33: Hybrid Ensemble Framework for Heart Disease Detection and Prediction

Authors: Elham Nikookar, Ebrahim Naderi

PAGE 243 – 248

Paper 34: M/M/1/n+Flush/n Model to Enhance the QoS for Cluster Heads in MANETs

Authors: Aleem Ali, Neeta Singh, Poonam Verma

PAGE 249 – 254

Paper 35: Binary PSO GSA for Load Balancing Task Scheduling in Cloud Environment

Authors: Thanaa S. Alnusairi, Ashraf A. Shahin, Yassine Daadaa

PAGE 255 – 264

Paper 36: The P System Design Method based on the P Module

Authors: Ping Guo, Xixi Peng, Lian Ye

PAGE 265 – 274

Paper 37: Cascades Neural Network based Segmentation of Fluorescence Microscopy Cell Nuclei

Authors: Sofyan M. A. Hayajneh, Mohammad H. Alomari, Bassam Al-Shargabi

PAGE 275 – 285

Paper 38: An Automatic Segmentation Algorithm for Solar Filaments in H-Alpha Images using a Context-based Sliding Window

Authors: Ibrahim A. Atoum

PAGE 286 – 291

Paper 39: Relative Humidity Profile Estimation Method with AIRS (Atmospheric Infrared Sounder) Data by Means of SDM (Steepest Descend Method) with the Initial Value Derived from Linear Estimation

Authors: Kohei.Arai

PAGE 292 – 299

Paper 40: Tuning of Customer Relationship Management (CRM) via Customer Experience Management (CEM) using Sentiment Analysis on Aspects Level

Authors: Hamed AL-Rubaiee, Khalid Alomar, Renxi Qiu, Dayou Li

PAGE 300 – 312

Paper 41: An Accurate Multi-Biometric Personal Identification Model using Histogram of Oriented Gradients (HOG)

Authors: Mostafa A. Ahmad, Ahmed H. Ismail, Nadir Omer

PAGE 313 – 319

Paper 42: Detection of Mass Panic using Internet of Things and Machine Learning

Authors: Gehan Yahya Alsalat, Mohammad El-Ramly, Aly Aly Fahmy, Karim Said

PAGE 320 – 329

Paper 43: Motif Detection in Cellular Tumor p53 Antigen Protein Sequences by using Bioinformatics Big Data Analytical Techniques

Authors: Tariq Ali, Sana Yasin, Umar Draz, M. Ayaz Arshad, Tayyaba Tariq, Sarah Javaid

PAGE 330 – 338

Paper 44: Investigating Methods of Resource Provisioning Mechanisms in Cloud: A Review

Authors: Babur Hayat Malik, Talia Anwar, Sadaf Ilyas, Farheen Jafar, Munazza Iffikhar, Maryam Malik, Noreen Islam Deen

PAGE 339 – 348

Paper 45: University Notification Subscription System using Amazon Web Service

Authors: Babur Hayat Malik, Zaheer Mehmood Dar, Sabah Mubarik Kayani, Mahnoor Dar, Muhammad Hassan Shafiq, Imran Kabir, Fatima Masood, Hamna Zakriya, Asad Ali

PAGE 349 – 354

Paper 46: Performance Measurement Model of Mobile User Connectivity in Femtocell/Macrocell Networks using Fractional Frequency Re-use Scheme

Authors: Mehrin Anannya, Riad Mashrub Shourov

PAGE 355 – 362

Paper 47: Heart Failure Prediction Models using Big Data Techniques

Authors: Heba F. Rammal, Ahmed Z. Emam

PAGE 363 – 371

Paper 48: Towards Privacy Preserving Commutative Encryption-Based Matchmaking in Mobile Social Network

Authors: Fizza Abbas, Ubaidullah Rajput, Adnan Manzoor, Imtiaz Ali Halepoto, Ayaz Hussain

PAGE 372 – 375

Paper 49: Towards Security as a Service to Protect the Critical Resources of Mobile Computing Devices

Authors: Abdulrahman Alreshidi

PAGE 376 – 383

Paper 50: Comparison of Task Scheduling Algorithms in Cloud Environment

Authors: Babur Hayat Malik, Mehwashma Amir, Bilal Mazhar, Shehzad Ali, Rabiya Jalil, Javaria Khalid

PAGE 384 – 390

Paper 51: Technical and Perceived Usability Issues in Arabic Educational Websites

Authors: Mohamed Benaida, Abdallah Namoun

PAGE 391 – 400

Paper 52: Automatic Sign Language Recognition: Performance Comparison of Word based Approach with Spelling based Approach

Authors: Shazia Saqib, Syed Asad Raza Kazmi, Khalid Masood, Saleh Alrashed

PAGE 401 – 405

Paper 53: Geographical Distance and Communication Challenges in Global Software Development: A Review

Authors: Babur Hayat Malik, Saeed Faroom, Muhammad Nauman Ali, Nasir Shehzad, Sheraz Yousaf, Hammad Saleem

PAGE 406 – 414

Paper 54: Gaze Direction based Mobile Application for Quadriplegia Wheelchair Control System

Authors: Muayad Sadik Croock, Salih Al-Qaraawi, Rawan Ali Taban

PAGE 415 – 426

Paper 55: A Study on Usability Awareness in Local IT Industry

Authors: Mahmood Ashraf, Lal Khan, Muhammad Tahir, Ahmed Alghamdi, Mohammed Alqarni, Thabit Sabbah, Muzafar Khan

PAGE 427 – 432

Paper 56: Monitoring Vaccine Cold Chain Model with Coloured Petri Net

Authors: Fatima Ouzayd, Hajar Mansouri, Manal Tamir, Raddouane Chiheb, Zied Benhouma

PAGE 433 – 438

Paper 57: Framework for Rumors Detection in Social Media

Authors: Rehana Moin, Zahoor-ur-Rehman, Khalid Mahmood, Mohammad Eid Alzahrani, Muhammad Qaiser Saleem

PAGE 439 – 444

Paper 58: Modeling of Arduino-based Prepaid Energy Meter using GSM Technology

Authors: Uzair Ahmed Rajput, Khalid Rafique, Abdul Sattar Saand, Mujtaba Shaikh, Muhammad Tarique

PAGE 445 – 449

Paper 59: Koch Island Fractal Patch Antenna (KIFPA) for Wideband Applications

Authors: Meryem HADJI, Sidi Mohammed MERIAH, Djamila ZIANI

PAGE 450 – 455

Paper 60: Investigating Saudi Parents' Intention to Adopt Technical Mediation Tools to Regulate Children's Internet Usage

Authors: Ala'a Bassam Al-Naim, Md Maruf Hasan

PAGE 456 – 464

Paper 61: Control of Industrial Systems to Avoid Failures: Application to Electrical System

Authors: Yamen EL TOUATI, Saleh ALTOWAIJRI, Mohamed AYARI

PAGE 465 – 470

Paper 62: Mobility Management Using the IP Protocol

Authors: Imtiaz A. Halepoto, Adnan Manzoor, Nazar H. Phulpoto, Sohail A. Memon, Muzamil Hussain

PAGE 471 – 475

Paper 63: Experimental Results on Agent-Based Indoor Localization using WiFi Signaling

Authors: Stefania Monica, Federico Bergenti

PAGE 476 – 488

Paper 64: Divide and Conquer Approach for Solving Security and Usability Conflict in User Authentication

Authors: Shah Zaman Nizamani, Waqas Ali Sahito, Shafique Awan

PAGE 489 – 495

Paper 65: A High-Performing Similarity Measure for Categorical Dataset with SF-Tree Clustering Algorithm

Authors: Mahmoud A. Mahdi, Samir E. Abdelrahman, Reem Bahgat

PAGE 496 – 509

Paper 66: Efficient Community Detection Algorithm with Label Propagation using Node Importance and Link Weight

Authors: Mohsen Arab, Mahdieh Hasheminezhad

PAGE 510 – 518

Paper 67: An Indefinite Cycle Traffic Light Timing Strategy

Authors: Ping Guo, Daiwen Lei, Lian Ye

PAGE 519 – 524

Paper 68: Search Manager: A Framework for Hybridizing Different Search Strategies

Authors: Yousef Abdi, Yousef Seyfari

PAGE 525 – 540

Paper 69: A Study of Feature Selection Algorithms for Predicting Students Academic Performance

Authors: Maryam Zaffar, Manzoor Ahmed Hashmani, K.S. Savita, Syed Sajjad Hussain Rizvi

PAGE 541 – 549

Cardiotocographic Diagnosis of Fetal Health based on Multiclass Morphologic Pattern Predictions using Deep Learning Classification

Julia H. Miao¹, Kathleen H. Miao^{1,2}

¹Cornell University, Ithaca, NY 14850, USA

²New York University School of Medicine, New York, NY 10016, USA

Abstract—Medical complications of pregnancy and pregnancy-related deaths continue to remain a major global challenge today. Internationally, about 830 maternal deaths occur every day due to pregnancy-related or childbirth-related complications. In fact, almost 99% of all maternal deaths occur in developing countries. In this research, an alternative and enhanced artificial intelligence approach is proposed for cardiotocographic diagnosis of fetal assessment based on multiclass morphologic pattern predictions, including 10 target classes with imbalanced samples, using deep learning classification models. The developed model is used to distinguish and classify the presence or absence of multiclass morphologic patterns for outcome predictions of complications during pregnancy. The testing results showed that the developed deep neural network model achieved an accuracy of 88.02%, a recall of 84.30%, a precision of 85.01%, and an *F*-score of 0.8508 in average. Thus, the developed model can provide highly accurate and consistent diagnoses for fetal assessment regarding complications during pregnancy, thereby preventing and/or reducing fetal mortality rate as well as maternal mortality rate during and following pregnancy and childbirth, especially in low-resource settings and developing countries.

Keywords—Activation function; deep learning; deep neural network; dropout; ensemble learning; multiclass; regularization; cardiotocography; complications during pregnancy; fetal heart rate

I. INTRODUCTION

In 2012, approximately 213 million pregnancies occurred worldwide [1]. Of those pregnancies, 190 million (89%) occurred in developing countries and 23 million (11%) were in developed countries. In 2013, complications of pregnancy resulted in 293,336 deaths due to maternal bleeding, complications of abortion, high blood pressure, maternal sepsis, and obstructed labor [2]. According to the World Health Organization [3], roughly 303,000 women died during and following pregnancy and childbirth in 2015, and in 2016, about 830 women died every day from pregnancy-related or childbirth-related complications around the world.

Medical complications of pregnancy and pregnancy-related deaths, impacting mothers and/or their babies, continue to remain a major global challenge today. Maternal death rate is especially concentrated in several specific areas of the world. In fact, almost 99% of all maternal deaths occur in developing countries [3]. This high and uneven mortality distribution reflects global inequities of access to medical services and

medical treatment. There are large mortality differences not only between countries but also within countries. These disparities in mortality rates persist even between high-income and low-income women, as well as between women living in rural areas and urban areas. Complications during pregnancy and childbirth are thus among the leading causes of death in developing countries [2], [3]. Most of these complications develop during pregnancy, while other complications may happen before pregnancy but are further worsened over the course of pregnancy. However, almost all of these maternal deaths during pregnancy occurred in low-resource settings, and most of them could have been prevented or treated.

Complications of pregnancy may include disorders of high blood pressure, gestational diabetes, infections, preeclampsia, pregnancy loss and miscarriage, preterm labor, and stillbirth. Other complications include severe nausea, vomiting, and iron-deficiency anemia [4], [5]. Thus, pregnancy may be affected due to these conditions, which require additional ways of assessing and evaluating fetal well-being. These conditions may involve medical problems in the mother that could impact on the fetus, pregnancy-specific problems, and diseases of the fetus [6]. In association with increased risk to the fetus, medical problems in the mother include essential hypertension, pre-eclampsia, renal and autoimmune disease, maternal diabetes, and thyroid disease [7]-[10]. Other medical problems in pregnancy, which pose increased risk to fetal health, are prolonged pregnancy, vaginal bleeding, reduced fetal movements, and prolonged ruptured membranes [11]. Additionally, intrauterine growth restriction, fetal infection, and multiple pregnancies increase the risks to the fetuses [11], [12]. As a result, these conditions could lead to neurodevelopmental problems in infancy, such as non-ambulant cerebral palsy, developmental delay, auditory and visual impairment, and fetal compromise, which might lead to morbidity or mortality in the newborn.

In order to assess fetal well-being and monitor for increased risk of complications of pregnancy, cardiotocography (CTG) is widely used as a technical method of continuously measuring and recording the fetal heart rate (FHR) and uterine contractions during pregnancy. This provides the possibility of monitoring the development of fetal hypoxia and intervening appropriately before severe asphyxia or death occurs [13]. In association with uterine contractions, the FHR along with its variability, reactivity, and possible decelerations are important measurements for assessment of fetal well-being [14]. The

FHR can be obtained via an ultrasound transducer placed on the mother's abdomen. The CTG, which depends on FHR, uterine contractions, and fetal movement activity, is utilized to detect and identify dangerous situations for the fetus. During the antepartum and intrapartum periods in pregnancy and childbirth, the CTG is often used for assessment and evaluation of fetal conditions by obstetricians.

Recently, advanced technologies in modern medical practices have successfully allowed robust and reliable machine learning and artificial intelligence techniques to be utilized in providing automated prognosis based on early detection outcomes in many medical applications [15]-[18]. The implementation and feasibility of machine learning tools can significantly aid medical practitioners in making informed medical decisions and diagnoses, potentially reducing maternal and fetal mortality and complications during pregnancy and childbirth and significantly aiding populations in both developing and developed countries. Diagnosing the FHR multiclass morphologic pattern is a challenging task, but computer-aided detection (CAD) based on machine learning technologies have been developed to provide automated classifications for fetal state during pregnancy [19]. Previous research reports used CAD approaches to diagnose the fetal state in pregnancy based on a method of support vector machines (SVM) with a Gaussian kernel function [20], [21]. Other research reports included classification of cardiocograms using Neural Network and Random Forest classifiers [22], [23]. However, these above mentioned machine learning methods and approaches were designed and developed to classify and predict only binary outcomes of normal versus abnormal cases during medical diagnosis in patients or as normal versus pathological cases in pregnancy using clinical diagnostic datasets of patients and CTG data in pregnancy, respectively.

In this research, an alternative and enhanced artificial intelligence approach is proposed for CTG diagnosis of fetal assessment based on multiclass morphologic pattern predictions, including 10 target classes, using deep learning classification models. The designed and developed deep learning classification and prediction models include two systems: a deep learning-based training classification model and a deep learning-based prediction model (also known as a diagnosis model). The training classification model is based on a multilayer perceptron with a multiclass softmax classification using deep learning technologies. The diagnosis model is used to distinguish and classify the presence or absence of multiclass morphologic patterns for outcome predictions of complications during pregnancy. By uniquely integrating multiclass morphologic patterns and predictions instead of binary predictions of normal versus pathological cases, the models provide a more reliable and specific diagnosis based on fetal health assessment with CTG. The performances of the deep learning-based classification and prediction model for diagnosing multiclass morphologic patterns in pregnancy were evaluated using multiclass measures based on averages of recalls (also known as sensitivities), precisions, *F*-scores, misclassification errors, and diagnostic accuracies.

II. CARDIOTOCOGRAPHY DATA DESCRIPTION

In this section, descriptions and characteristics of the CTG data sets regarding complications in pregnancy are introduced. The CTG data sets, which have been used in this research paper, were obtained from the CTG databases available in the UCI Machine Learning Repository [24]. These databases consist of data information on measurements of FHR and uterine contraction features during pregnancy based on Cardiocograms, which were contributed by the Biomedical Engineering Institute, Porto, Portugal, and the Faculty of Medicine, University of Porto, Portugal in September 2010. These data sets were collected based on clinical instances in pregnancy in 1980 and periodically from 1995 to 1998, resulting in a constantly increasing dataset size.

There are a total of 2,126 clinical instances, representing different complications of pregnancy on fetal cardiocograms in the CTG dataset. The clinical instances on the fetal cardiocograms were automatically processed, and their respective diagnostic features were measured. These clinical instances were also classified with respect to a morphologic pattern by three expert obstetricians and had consensus classification labels assigned to each of them. Each clinical instance in the CTG dataset contains 21 input attributes and one multiclass attribute as well as one fetal state. The multiclass attribute represents the multiclass morphologic patterns, which includes the 10 target classes. Additionally, this multiclass attribute is represented by an integer valued from "1" to "10", where each of integers represents one of the morphologic patterns in pregnancy. The fetal state is assigned one of the 3 classes, including normal, suspect, or pathologic cases. Thus, the CTG dataset can be used for building classification and predictive models based on the 10-class, 3-class, or even the 2-class classification experiments by eliminating the suspect class category in fetal state. In previous reports [20]-[23], several machine learning-based classification models eliminated all suspect cases in fetal state and were established based on only a binary classification of the fetal state in terms of normal and pathologic cases.

In this research paper, the multiclass morphologic patterns, including all 10 of the target classes, have been utilized for developing the deep neural network classification and prediction models. Pattern recognition and prediction of multiple target outcomes is a challenging task in the field of machine learning and artificial intelligence; since multiclass morphologic patterns with imbalanced sample sizes of the 10 target classes in the CTG data will be used, this task thus proves advanced. Ultimately, however, the integration of all 10 target classes and multiclass morphologic patterns, compared to previous research model's use of binary classification, allows a more reliable and realistic diagnosis and prediction of multiclass outcomes, thus aiding patients with an accurate and more specific fetal health assessment.

III. DEEP NEURAL NETWORK ARCHITECTURES AND CLASSIFIERS

In this section, we present a CTG classification and diagnosis model for the multiclass morphologic pattern prediction, representing the 10 target classes for fetal outcome forecasting in pregnancy using the deep neural network

classification and prediction models along with corresponding algorithms, methods, approaches, and architecture. Furthermore, we discuss some of the special techniques that can be used to prevent overfitting and enhance the deep neural network classification and prediction model performances for multiclass morphologic pattern prediction and CTG diagnosis.

Deep learning consists of neural networks that teach themselves and make decisions autonomously. Deep learning methods and architectures have gained significant attention in the area of artificial intelligence in recent years. It has recently expanded exponentially in both academic communities and global high-tech industries since 2011 [25]. One of the most important deep learning architectures is a Deep Belief Network, which is built by stacking a set of restricted Boltzmann machines (RBM) [26]-[28]. The RBM is a generative stochastic artificial neural network that can learn from a probability distribution over its input data. Depending on an objective task, the RBM can also be trained either for supervised learning or for unsupervised learning. Another important deep learning architecture is called deep auto-encoder [29]. The deep auto-encoder is also an artificial neural network, usually used for unsupervised learning [29]-[31]. The deep auto-encoder is capable of learning an encoding representation based on a set of input data and has become more widely used for learning generative models of data.

A traditional multilayer perceptron neural network model can be considered as a processor that acquires and stores experiential knowledge through a machine learning process during a training process [15], [17]. In order to retain the knowledge, synaptic weights that resemble interneuron connections are used. The training process of a learning algorithm involves the modification of the synaptic weights of the model in order to obtain a desired objective. The multilayer perceptron neural network model uses a back-propagation approach for training the neural network classification unit during the training process. The back-propagation approach based on the Widrow-Hoff learning rule [15], [17], [32] can be used to minimize the objective function for the neural network model. The input data and the corresponding output data are used to train the neural network classification model until the model appropriately approximated a function within a prior defined error value. During the training process, a learning algorithm is used to adjust weights and biases by utilizing the derivative vectors of errors back-propagated through the neural network classification unit.

In theory, deep neural networks and the multilayer perceptron neural networks, which have the exact same network structure and computations, perform similarly if the deep neural networks and multilayer perceptron neural networks have been given the same conditions. Both deep neural networks and the multilayer perceptron neural networks consist of perceptrons in terms of linear and nonlinear transformation functions. The nonlinear transformation functions between layers of perceptrons enable neural networks to be used for modeling nonlinear behaviors.

Deep neural networks differ from multilayer perceptron neural networks in terms of the network depth, which is determined according to the number of hidden layers in the network. In general, a neural network with three or more hidden layers is considered as a deep neural network. In that case, the higher layers are building new abstractions on top of previous layers, usually leading to learning a better solution with the deep neural network. On the other hand, the number of hidden layers in the network also entails difficulties to train the network in practice. This is because increasing the number of hidden layers in the networks leads to two major issues:

1) Vanishing gradients: The back-propagation approach [15], [17] becomes less helpful in passing information from the higher layers to the lower layers. The gradients become small relative to the weights of the networks and begin to almost vanish when information is passed back.

2) Overfitting: The deep neural network classification model performs very well with a training dataset, but the model shows poorer performances on a real testing dataset. Overfitting is the central problem in the field of machine learning and artificial intelligence.

Fig. 1 shows the deep neural network architectures based on the multilayer perceptron neural networks in detail. Architecturally, the simplest form of the deep neural networks is a feedforward and non-recurrent neural network very similar to the multilayer perceptron neural network, which has an input layer (green color), an output layer (green color), and one or more hidden layers connecting them, but the number of hidden layers consist of a set of active nodes (blue color) and in-active nodes (red color). In this research paper, we have used this type of deep neural network architecture as a fundamental network system to develop the deep learning-based neural network classification and prediction models. The designed deep neural network architecture allows us to classify the multiclass morphologic patterns with imbalanced sample sizes of the 10 target classes in the CTG data for fetal assessment during pregnancy.

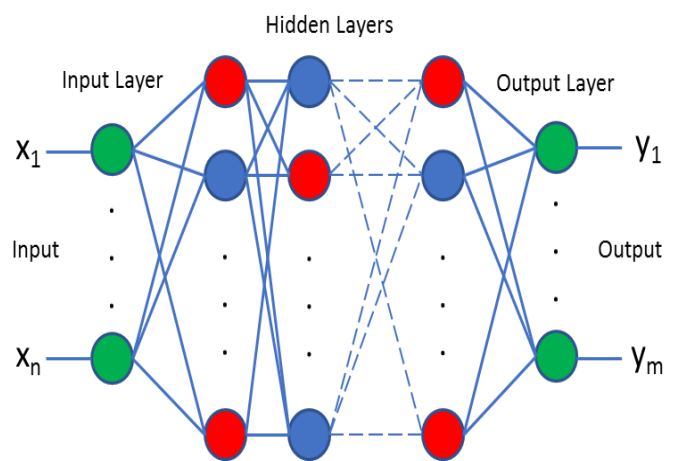


Fig. 1. A deep neural network architecture with an input layer, the number of hidden layers, and an output layer, where the blue cycles are active nodes, the red cycles are in-active nodes in the number of the hidden layers, and the green cycles present an input layer and an output layer.

The deep neural network architecture is composed of the multiple perceptrons, which are stacked one after the other in a layer-wise fashion. The input matrix data X is fed into the input layer, which is a multidimensional perceptron with a weight matrix W_1 , bias vector B_1 , and a transfer function Φ_1 . The output of the input layer is then fed into the first hidden layer, which is a perceptron with another weight matrix W_2 , bias vector B_2 , and a transfer function Φ_2 . This process continues for every one of the L hidden layers, which is again a perceptron with another weight matrix W_L , bias vector B_L , and a transfer function Φ_L until we reach the output layer. As can be seen, according to Fig. 1, we refer to the first layer as the input layer, the last layer as the output layer, and every other layer as a hidden layer in the network architecture.

The deep neural network architecture with one hidden layer has a mathematical representation:

$$Y = \Phi_2(\Phi_1(XW_1 + B_1)W_2 + B_2), \quad (1)$$

where Y is an output matrix data. The deep neural network architecture with two hidden layers computes a function in the following:

$$Y = \Phi_3(\Phi_2(\Phi_1(XW_1 + B_1)W_2 + B_2)W_3 + B_3), \quad (2)$$

and, in general, the deep neural network architecture with the number of $(L-1)$ hidden layers calculates an output function given by:

$$Y = \Phi_L(\dots \Phi_3(\Phi_2(\Phi_1(XW_1 + B_1)W_2 + B_2)W_3 + B_3) \dots)W_L + B_L, \quad (3)$$

where the transfer function $\Phi_n, n = 1, 2, \dots, L$, can be either a linear or a nonlinear transfer function.

A. Activation Function

An activation function of a neural node in the neural network defines an output of that neural node given a set of inputs. In an artificial neural network or deep neural network architecture, this activation function is also called a transfer function, which can be a linear or non-linear transfer function. The most common transfer functions that are used in deep learning or deep neural network architectures are as follows:

1) Scaled exponential linear unit

$$F(x) = \lambda \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases} \quad (4)$$

where α and λ are hyper-parameters to be adjusted, $\alpha > 0$ and $\lambda > 0$. When $\alpha = 0$ and $\lambda = 1$, (4) is called the *Rectified linear unit* (ReLU); where $\lambda = 1$, (4) is known as the *exponential linear unit* (ELU).

2) Sigmoid function unit

$$F(x) = \frac{1}{1+e^{-x}} \quad (5)$$

Equation (5) is a *sigmoid function*, which is real-valued and differentiable, having a non-negative or non-positive first derivative, one local minimum, and one local maximum. In general, the sigmoid function exists as a range from 0 to 1 and is used for binary classification. It is especially used for classification models where we want to predict a probability as an output.

3) Softmax function

$$F(x_j) = \frac{e^{x_j}}{\sum_k e^{x_k}} \quad (6)$$

where x is a vector of the inputs to the output layer, and $j = 1, 2, \dots, K$, indexes the output units. The *softmax function* is often used for any number of multiclass classifications.

In this research paper, the ReLU function has been used in the input and hidden layers, and the softmax function is used in the output layer for the deep neural network architecture.

B. Dropout

One of the primary pitfalls of machine learning and artificial intelligence is overfitting when the model catastrophically sacrifices generalization performances for the purposes of minimizing training loss. In other words, a deep neural network model performs really well based on a training dataset. However, in practice, the deep neural network model performs much more poorly on testing data or real unseen testing data. Indeed, overfitting is one of the key critical problems in the field of machine learning and artificial intelligence.

Deep neural networks, which consist of multiple linear and non-linear hidden layers, have a self-learning capability of capturing very complicated relationships between their inputs and outputs. However, with a limited size of a training dataset, many of these complicated relationships could be the result of sampling noise, that is, they may exist in the training set but not in the real testing data. This leads to overfitting. In addition, large deep neural networks are slow to train for use in applications, thereby making it difficult to handle overfitting by combining the predictions of many different large neural nets at testing time.

Dropout is one of several effective and powerful regularization techniques to prevent deep neural network architectures from overfitting [33], [34]. The key idea of the dropout, based on a theory of the role of sex in evolution [35], is to randomly eliminate units along with their network connections from the deep neural networks during a training process. In other words, the central idea of the dropout is to take a deep neural network classification model, which overfits easily, and to train only smaller subsets of the classification models from the deep neural network architectures. As a result, the dropout technique can prevent the deep neural network units from co-adapting too much, thereby avoiding a single node specializing to a task.

Additionally, a dropout technique for the deep neural networks can be viewed as an alternative form of ensemble learning, in which each member of the ensemble learning is trained based on a different subsample of the input data, thereby resulting in learning only a subset of the entire input feature space. At each training step within a batch size, the dropout technique creates a different deep neural network by randomly removing some of the neural units from the hidden layer and/or even input and output layers. Conceptually, the dropout technique actually achieves a similar outcome such that an ensemble learning system uses many different deep neural networks at each of the steps (or batch size) with a subset of input data during the training process. During the

testing process, the deep neural network is only used with the scaled down weights (or partial weights in the network) instead of using entire neural units. Thus, from a point of mathematical view, the dropout technique approximates ensemble averaging using the geometric mean as average [36].

The dropout technique for the deep neural networks has been especially successful in applications because of its simplicity and remarkable effectiveness as a regularization function as well as its interpretation as an exponentially large ensemble learning for the deep neural networks. As a result, this dropout technique implemented in the deep neural network model significantly reduces overfitting issues and provides major improvements over other regularization methods.

C. Regularization

Regularization is one of the key elements of machine learning and artificial intelligence, especially in deep learning. The regularization allows deep neural networks to generalize well to testing data even when the networks are trained based on a finite training set or an imperfect optimization procedure [37], [38]. In other words, regularization can be considered as any modification to a learning algorithm in which is intended to reduce its network test error but not its network training error. In essence, regularization is a supplementary technique that can be used to make the model performance better in general and to produce better results on the testing data [38]. In conjunction with the dropout technique, regularization is another mathematical method for combating overfitting for the deep neural networks.

One of the most popular regularizations is L_2 regularization also known as weight decay, which takes a more direct approach than the dropout technique for regularizing. Generally, a common underlying cause for overfitting is that the deep neural network classification model is too complex in terms of large parameters for the problem based on a training data set. In other words, the regularization can be used to decrease complexity of the deep neural network classification model while maintaining the same number of large parameters. Thus, in order to minimize a L_2 norm, the regularization does so by penalizing weights with large magnitudes using a hyper-parameter λ to specify the relative importance of the L_2 norm for minimizing the loss on the training data set.

Formally, training a deep neural network f_θ is to find a weight function $\theta(\mathbf{w}, \mathbf{b})$, where \mathbf{w} and \mathbf{b} denote weights and bias, respectively, such that the expected regularized loss can be minimized:

$$E(\theta, D) = \arg \min_{\theta} \left\{ \frac{1}{D} \sum_{(x_i, t_i) \in D} E(f_\theta(x_i), t_i) \right\} + \lambda \|\theta\|_p, \quad (7)$$

where D is a training data and (x_i, t_i) are samples in the training data D ; the x_i are inputs and t_i are targets. The hyper-parameter λ can be used to control the relative importance of the regularization function. The first item and second item in (7) are referred to as an error function and a regularization error, where

$$\|\theta\|_p = \left(\sum_{j=0}^N |\theta_j|^p \right)^{\frac{1}{p}}, \quad (8)$$

which is the L_p norm of θ . If $p = 1$, (8) is L_1 regularization. If $p = 2$, (8) is L_2 regularization. Note that the error function in

(7), which is dependent on the targets, assigns a penalty to model predictions according to whether or not the model predictions are consistent with the targets. The regularization error assigns a penalty to the model depending on anything except the targets.

D. Initialization

Training deep neural networks is difficult because of vanishing or exploding activations and gradients. The central challenge in training deep neural networks is about how to deal with the strong dependencies that exist during training between the parameters across a large number of hidden layers. This is because a solution to a non-convex optimization algorithm, such as the method of stochastic gradient descent, heavily depends on initialization weights in the deep neural networks. In other words, if the initialization weights in the deep neural networks start too small, then the initialization weights shrink as they pass through each of the hidden layers until they are too tiny to be useful. On the other hand, if the initialization weights in the deep neural networks begin too large, then the initialization weights quickly rise as they pass through each hidden layer until they are too large to be useful. These behaviors are referred to as saturation in training deep neural networks because of nonlinear activation functions embedded in the hidden layers.

Note that deep neural networks with linear and/or nonlinear activation functions initialized from unsupervised pre-training methods, such as deep RBM and deep auto-encoder [39]-[41], do not suffer from these saturation behaviors. Consequently, another important note is that even in the presence of very large amounts of training data in a supervised learning, stochastic gradient descent (SGD) is still subject to a degree of overfitting to the training data. In that sense, unsupervised pre-training method based on the deep RBM and deep auto-encoder interacts intimately with the optimization process. The positive effect of the unsupervised pre-training method is seen not only in generalization error but also in training error when the amount of training data becomes large.

Training deep RBM and deep auto-encoder as an unsupervised pre-training method for the deep neural networks can be considered a breakthrough in effective training strategies [26], [42]-[44]. The unsupervised pre-training method is generally based on greedy layer-wise unsupervised pre-training followed by supervised fine-tuning [41]. Each layer is pre-trained with an unsupervised learning algorithm by learning nonlinear activation functions of their inputs from the previous layers, which capture the major variations in their inputs. Lastly, the unsupervised pre-training method establishes the stage for a final training phase in which the deep neural networks is fine-tuned with respect to a supervised learning criterion of the gradient-based optimization.

Another initialization method for the deep neural networks is known as *Xavier* initialization [39], which is used to make sure that the weights are in a reasonable range of values throughout many hidden layers. Assume that there is an input X with N components and a linear neuron (or combination) with random weights W :

$$Y = \sum_{i=1}^N w_i X_i. \quad (9)$$

The variance of this liner combination Y is given by [45]:

$$Var(\sum_{i=1}^N w_i X_i) = \sum_{i=1}^N w_i^2 Var(X_i) + \sum_{i \neq j} w_i w_j Cov(X_i, X_j). \quad (10)$$

If the random variables X_1, \dots, X_N are independent and identically distributed, this always leads to uncorrelated random variables such that

$$Cov(X_i, X_j) = 0, \text{ for } i \neq j. \quad (11)$$

Thus, (11) is rewritten to

$$Var(\sum_{i=1}^N w_i X_i) = \sum_{i=1}^N w_i^2 Var(X_i). \quad (12)$$

In addition, we further assume that the deep neural network weights w_i and inputs X_i are uncorrelated and both have zero-mean:

$$\sum_{i=1}^N w_i^2 Var(X_i) = \sum_{i=1}^N Var(W) Var(X) = NVar(W) Var(X). \quad (13)$$

Comparing (12) to (13), we obtain a result as follows:

$$Var(\sum_{i=1}^N w_i X_i) = NVar(W) Var(X). \quad (14)$$

Equation (14) implies that the variance of the output is the variance of the input with a scaled function by $NVar(W)$. If we further want to make the variance of the input to be the same as the variance of the output, it must hold $Var(W) = \frac{1}{N}$ for the inputs so that we are able to preserve variance of the inputs after passing through a number of the hidden layers. For the backpropagation update, we also need to ensure that $Var(W) = \frac{1}{M}$ for the outputs. Thus, in general, for implementation of the initialization on the deep neural networks, the variance of the weights for the deep neural networks can be set to their average based on the inputs and outputs, that is,

$$Var(W) = \frac{1}{N+M}. \quad (15)$$

IV. MULTICLASS EVALUATION METHODS ON DEEP NEURAL NETWORK

In this section, multiclass evaluation methods for the performances of the deep neural network classification model in multiclass morphologic pattern prediction based on the CTG data are discussed in detail.

The evaluation of the model performances for the deep neural networks is typically based on testing data sets, rather than analytically in the field of machine learning and artificial intelligence. The classification effectiveness of machine learning models, deep neural networks, and/or any other type of models can usually be measured in terms of model sensitivity (also known as recall), specificity, precision, F -score, accuracy, and misclassification error [45], [16]. In this section, we extend the evaluation methods for the effectiveness measurements of the deep neural network classification model from a binary classification to a multiclass classification problem.

Let C_1, \dots, C_K be multiclass labels, in which we want to predict K labels using deep neural network classification models. For correct decisions, let TP be a decision to assign similar multiclass to the same cluster, and let TN be a decision to assign dissimilar multiclass to different clusters. On the

other hand, for incorrect decisions, let FP be a decision to assign dissimilar multiclass to the same cluster, and let FN be a decision to assign similar multiclass to different clusters.

For the effectiveness measurement of the deep neural network classification models, a global calculation of the TP , TN , FP , and FN can be computed in the following:

$$TP = \sum_{i=1}^N TP_i, \quad (16)$$

$$TN = \sum_{i=1}^N TN_i, \quad (17)$$

$$FP = \sum_{i=1}^N FP_i, \quad (18)$$

and

$$FN = \sum_{i=1}^N FN_i, \quad (19)$$

where TP_i , FP_i , TN_i , and FN_i denote the local measures, representing the number of *true positive*, *false positive*, *true negative*, and *false negative* test examples with respective to the i -th class label. This evaluation method is referred to as *micro average* [46], [47].

Then, *sensitivity* for the multiclass classification is defined as the probability of correctly identifying those with the true positive rate [16]:

$$Sensitivity = \frac{TP}{TP+FN}, \quad (20)$$

where sensitivity is known as *recall* in the field of machine learning and deep learning. The *specificity* is defined as the probability of correctly identifying true negative rate:

$$Specificity = \frac{TN}{FP+TN}, \quad (21)$$

The *precision* is defined as

$$Precision = \frac{TP}{TP+FP}, \quad (22)$$

In performance measures of the deep neural networks, the recall in (20) is a measure of quantity while the precision in (22) is a measure of quality. Both the recall and precision are in a mutual relationship based on the understanding and measure of relevance.

Another method, which is similar to the one-vs-all classification technique, is to calculate all the local recalls R_i and precisions P_i for each C_i of the multiclass. This method is referred to as *macro average*. Then, the average of recalls is as follows:

$$\bar{R} = \frac{1}{K} \sum_{i=1}^K R_i, \quad (23)$$

and average of precisions is given by,

$$\bar{P} = \frac{1}{K} \sum_{i=1}^K P_i. \quad (24)$$

Note that micro and macro average methods represent different calculation behaviors, thereby leading to different results in the multiclass evaluation of the classification model effectiveness for the deep neural networks.

V. RESULTS

In this research paper, the deep neural network classification model is proposed for accurate diagnosis of fetal state based on the CTG data and the multiclass morphologic pattern outcome predictions. The proposed deep neural network classification model has a deep neural network

architecture, including 21 input units, first and second hidden layers, and 11 binary output units. The 11 binary output units, which allow us to form 10 unique sequences, can be used to represent the 10 target classes in the morphologic pattern outcomes on fetal assessment for multiclass classifications and predictions. The first hidden layer contains 105 units with each of the ReLU activation functions and 25% dropout rate of the network. The second hidden layer has 42 units, also connected with the ReLU activation functions, and 20% dropout rate of the network. Each of the 11 output units in the last stage of the deep neural network architecture is connected to a softmax activation function. For each batch process during the training of the deep neural network, the dropout rates in the first and second hidden layer are randomly applied to the deep neural network, thereby resulting in random connections within the deep neural network architecture. Doing so allows us to generate an alternative form of ensemble learning as well as reduce and/or prevent overfitting issues for the deep neural network classification model.

In the CTG data, each clinical instance consists of 40 raw attributes. Among all of the raw attributes, only 23 of them can be used for developing the deep neural network classification and prediction models. The other 13 attributes are not recommended to be used according to the attribute restriction in the CTG data. Table I lists the detailed 23 raw attributes, which had been used for the development of the deep neural network classification model. The variable of the “Class” in Table I is referred to as a target variable, which includes 10 integers from 1 to 10, representing different morphologic pattern behaviors of complications in pregnancy.

The proposed deep neural network classification and prediction models were applied to all clinical instances, which represent complications of pregnancy based on fetal assessments in the CTG data, and were used to predict the multiclass morphologic patterns with the 10 target class outcomes. Table II shows the details of the clinical instances in terms of the number of cases and percentages of the presence or absence of complications during pregnancy based on fetal assessments in each of the 10 morphologic pattern outcome data sets.

As can be seen in Table II, there are large differences in terms of percentages of the number of cases within each of the multiclass morphologic pattern outcomes in the CTG data. Class C3 has the lowest percentage of number of cases at 2.49%, while class C2 has the highest percentage of number of cases at 27.23%. As can be seen, the number of the sample distributions in the multiclass morphologic pattern outcomes would lead to a challenge in multiclass classification with imbalanced sample sizes for the CTG data in the field of machine learning and deep learning.

In order to evaluate the effectiveness of the performances of the developed deep neural network classification and prediction models, the model accuracy and misclassification error as well as the recall, precision, and *F*-score were estimated using a nonparametric approach based on a holdout method [45]. The holdout method applied by partitioning the CTG data into two mutually exclusive data sets, training data and testing data, respectively. The deep neural network

classification model was first trained using the training data, and then it was tested using the testing data.

In this research, the entire CTG data was randomly separated into 70% training and 30% testing data sets using the holdout method. The deep neural network classification and prediction models were trained and tested by using the 70% training and 30% testing data sets, respectively. The training and testing processes for the deep neural network classification and prediction models were repeated 24 times based on the different 70% training and 30% testing data sets. Doing so would determine an average of the testing performance results for the deep neural network classification and prediction models.

TABLE I. THE RAW ATTRIBUTES OF THE VARIABLE NAMES AND DESCRIPTIONS IN THE CTG DATA SET

Variable Name	Descriptions	Variable Name	Descriptions
LB	FHR baseline (beats per minute)	Min	Minimum of FHR histogram
AC (Second)	Number of accelerations	Max	Maximum of FHR histogram
FM (Second)	Number of fetal movements	Nmax	Number of histogram peaks
UC (Second)	Number of uterine contractions	Nzeros	Number of histogram zeros
DL (Second)	Number of light decelerations	Mode	Histogram mode
DS (Second)	Number of severe decelerations	Mean	Histogram mean
DP (Second)	Number of prolonged decelerations	Median	Histogram median
ASTV	Percentage of time with abnormal short term variability	Variance	Histogram variance
MSTV	Mean value of short term variability	Tendency	Histogram tendency
ALTV	Percentage of time with abnormal long term variability	Class	FHR pattern class code (1 to 10): 1-calm sleep 2-REM sleep 3-calm vigilance 4-active vigilance 5-shift pattern 6-accelerative or decelerative pattern (stress situation) 7-decelerative pattern (vagal stimulation) 8-largely decelerative pattern; 9-flat-sinusoidal pattern (pathological state) 10-suspect pattern
MLTV	Mean value of long term variability	NSP	fetal state class code (N=normal; S=suspect; P=pathologic)
Width	Width of FHR histogram		

TABLE II. CLINICAL INSTANCES OF COMPLICATIONS DURING PREGNANCY BASED ON FETAL ASSESSMENTS OF THE MULTICLASS MORPHOLOGIC PATTERNS IN THE 10 TARGET CLASS OUTCOMES IN THE CTG DATA

MP	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
NC	384	579	53	81	72	332	252	107	69	197
%	18.0	27.2	2.4	3.8	3.3	15.6	11.8	5.0	3.2	9.2
	6	3	9	1	9	2	5	3	5	7

Note: MP means Morphologic patterns; NC means Number of cases.

Fig. 2 shows a graph plot of the designed and developed deep neural network classification model performances using the training dataset at each of the epochs for 80,000 iterations during the training process. The accuracy of the deep neural network classification model is 97.32% based on the training dataset. Furthermore, Fig. 3 illustrates a graph plot of the deep neural network classification model loss function error throughout the 80,000 iterations during the training process. The loss function error of the deep neural network classification model is 0.0941 in an optimal sense of minimum mean square error (MMSE). In general, the higher the accuracy that the deep neural network classification model can achieve the lower the loss function error is obtained during the training process.

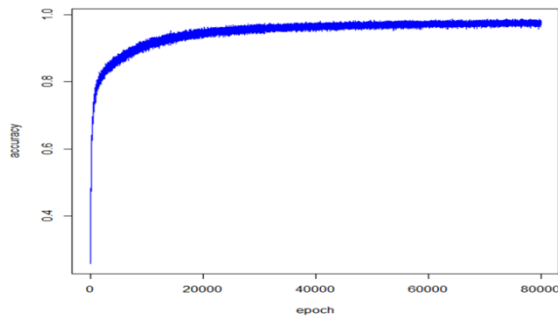


Fig. 2. A graph plot of the deep neural network classification model performance in terms of accuracy at each of the epochs for 80,000 iterations during the training process, where x-axis denotes each of the epochs and y-axis represents the model accuracy (the value of 1.0 at the y-axis represents a trained model with 100% accuracy).

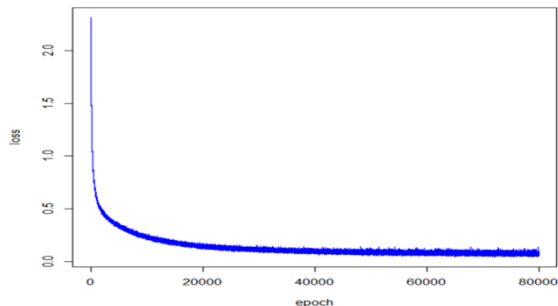


Fig. 3. A graph plot of the deep neural network classification model loss function error at each of the epochs for 80,000 iterations during the training process, where x-axis denotes each of the epochs and y-axis represents the model loss function error in MMSE.

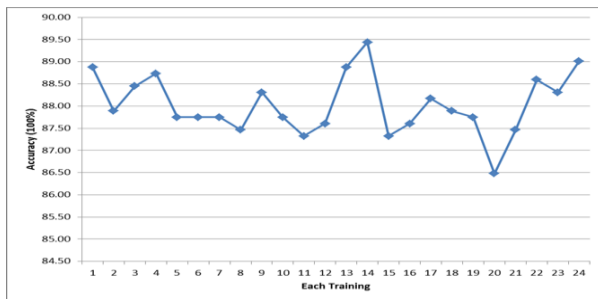


Fig. 4. A graph plot of the deep neural network classification model performances for 24 accuracy measures based on the 24 different testing data sets, where x-axis represents each of the 24 tests and y-axis is the tested model accuracy in percentage.

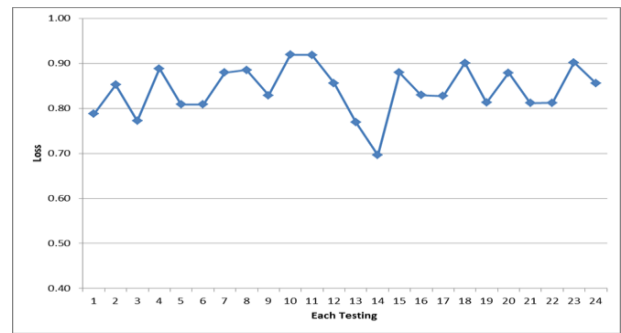


Fig. 5. A graph plot of the deep neural network classification model performances with 24 loss function error measures based on the 24 different testing data sets, where x-axis represents each of the 24 tests and y-axis is the tested model loss function error.

In this research paper, we repeated the same training and testing processes 24 times for the deep neural network classification model. For each of the 24 times, the entire CTG data was randomly divided into 70% training and 30% testing data sets. Then the deep neural network classification and prediction models were trained using the training data sets and tested using the testing data sets, respectively. The deep neural network classification model performances were recorded based on the 24 testing data sets. Displaying the results, Fig. 4 shows a graph plot of the deep neural network classification model performances in terms of 24 accuracy measures based on the 24 different testing data sets. As can be seen, the highest testing accuracy of the deep neural network classification model is 89.44%; the lowest testing accuracy is 86.48%. This leads to an average model accuracy of 88.02% with a standard deviation of 0.67%. The average misclassification error of the model is 11.98%. Correspondingly, Fig. 5 is a graph plot of the deep neural network classification model performances regarding the 24 loss function error measures based on the 24 different testing data sets as well. The best MMSE of the loss function error is 0.70 while the worst MMSE is 0.92. The average MMSE of the loss function errors is 0.84 with a standard deviation of 0.05.

In general, according to the training and testing results, whether the deep neural network classification model falls into a global or local minimum in an optimal sense is inconsequential. If the deep neural network overfitting in a minimum sense can be controlled, this deep neural network classification model would be determined to have realistically accurate diagnoses for fetal assessment during pregnancy based on the multiclass morphologic patterns of the 10 target class predictions.

Thus, in order to optimize the deep neural network classification model, the “rmsprop” method was used during the training process in this research. The “rmsprop” method is one of the mini-batch learning methods, which divides the learning rate for a weight by a running average of the magnitudes of recent gradients for the weight [48] and keeps a moving average of the squared gradient for each weight. In this research, a mini-batch size of 80 was used along with a learning rate of 0.00005 during the training processes.

Table III shows a combined confusion matrix (also known as an error matrix), which is a special table, using a summation

of the 24 individual confusion matrices based on 24 independently individual testing results. This combined confusion matrix allows us to visualize the deep neural network classification model performances in detail. Each row of the combined confusion matrix represents the number of clinical instances in a predicted class, while each of the columns represents the number of clinical instances in an actual multiclass. In statistics, this combined confusion matrix can also be called a contingency table along with two dimensions in terms of “actual” and “predicted” as well as identical sets of “classes” in both dimensions.

Based on results of the combined confusion matrix in Table III, corresponding values of recall, precision, and *F*-score for each of the multiclass morphologic patterns based on the 10 target classes using the deep neural network classification model were estimated as shown in Table IV. As can be seen, the averages of recall and precision are 84.30% and 84.91% along with standard deviations 8.39% and 6.89%, respectively. This leads to an average of the *F*-scores that is equal to 0.8453 with a standard deviation of 0.0737.

TABLE III. A COMBINED CONFUSION MATRIX BASED ON THE 24 INDEPENDENTLY INDIVIDUAL CONFUSION MATRICES USING THE 24 DIFFERENT TESTING DATA SETS

		Actual Multiclass Cases									
		C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
P C	2584	97	12 7	0	12 3	22	27	0	3	154	
	101	4736	0	60	14	140	0	0	10	0	
	126	3	37 1	0	0	2	4	0	0	0	
	0	110	0	54 0	1	0	0	0	0	0	
	31	27	0	0	45 7	3	9	0	9	28	
	24	67	0	0	0	2278	55	1	4	0	
	15	40	6	0	0	28	1798	75	4	0	
	0	0	0	0	0	47	0	61 7	32	0	
	0	0	0	0	0	0	0	1	55 6	100	
	143	8	0	0	29	0	3	2	12 6	1062	

Note: PC means predicted cases.

TABLE IV. RECALL, PRECISION, AND *F*-SCORE FOR EACH OF THE MULTICLASS MORPHOLOGIC PATTERNS

Multiclass	Recall	Precision	<i>F</i> -Score
C1	0.8545	0.8237	0.8388
C2	0.9308	0.9358	0.9333
C3	0.7361	0.7332	0.7347
C4	0.9000	0.8295	0.8633
C5	0.7324	0.8103	0.7694
C6	0.9040	0.9378	0.9206
C7	0.9483	0.9145	0.9311
C8	0.8865	0.8865	0.8865
C9	0.7473	0.8463	0.7937
C10	0.7902	0.7735	0.7817
Average	0.8430	0.8491	0.8453
Standard Deviations	0.0839	0.0689	0.0727

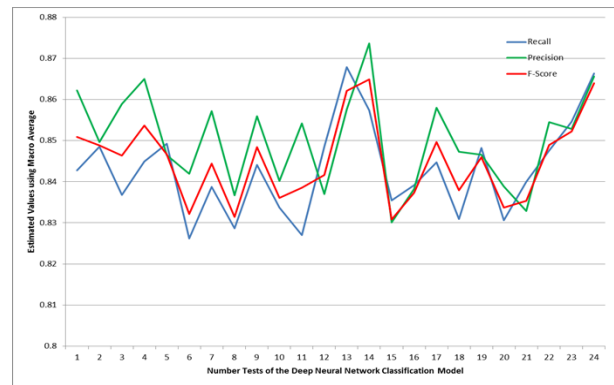


Fig. 6. A graphic plot of the macro average of recall, precision, and *F*-score based on the 24 individual confusion matrices for the deep neural network classification model using the 24 different testing data sets.

In this research, another method utilizing the macro average, which is the one-vs-all classification technique, was also used to estimate recall, precision, and *F*-score. Fig. 6 shows a graphic plot regarding the estimated values of recall, precision, and *F*-score based on the 24 independently individual confusion matrices. The macro average was used to compute all of the local recall and precision values for each of the multiclass morphologic patterns based on each of the 24 confusion matrices. The averages of recall and precision were then calculated according to (23) and (24) based on each of the confusion matrices, thereby leading to the 24 estimated values of recall, precision, and *F*-score as shown in Fig. 6. Based on the macro average, the averages of the recall, precision, and *F*-score are 84.30%, 85.01%, and 0.8508, respectively. Correspondingly, the standard deviations for the averaged recall, precision, and *F*-score are 1.13%, 1.14%, and 0.0100, respectively. As can be seen, in general, the *F*-score curve is displayed between the recall and precision curves.

Furthermore, it is noted that there are some differences in terms of averages and standard deviations of the recall, precision, and *F*-score for the deep neural network classification model using the combined confusion matrix as shown in Tables III and IV as well as using the macro average as shown in Fig. 6. However, generally speaking, there are no significant differences for the averages and standard deviations of the recall, precision, and *F*-score for the deep neural network classification model using the methods of micro and macro averages. Utilizing the 24 independently individual confusion matrices based on the 24 different testing data sets for estimating the deep neural network classification model performances allows us to establish and demonstrate a new and alternative way of representing the model recall, precision, and *F*-score in a dynamic representation.

VI. CONCLUSION AND FUTURE WORK

In this research paper, the deep neural network classification and prediction models were designed and developed for CTG diagnosis and prediction based on fetal assessment in pregnancy with the multiclass morphologic patterns of the 10 target classes with imbalanced samples. In conjunction with the dropout technique, regularization was applied to combat overfitting for the deep neural networks during the training process. As a result, the developed deep

neural network architecture allowed us to not only show a strong, alternative form of largely exponential ensemble learning but also reduce overfitting issues for the deep neural network classification and prediction models. Therefore, the developed deep neural network classification and prediction models can provide highly accurate and consistent diagnoses for fetal assessment regarding complications during pregnancy based on the multiclass morphologic patterns, thereby preventing and/or reducing fetal morbidity or mortality rate as well as maternal mortality rate during and following pregnancy and childbirth, especially in developing countries or in low-resource settings.

The dropout technique was used to enable us to randomly drop neural units with their connections in the deep neural network architecture. It can be treated as a large exponential ensemble learning for the deep neural networks. This significantly reduces overfitting and provides major improvements over traditional regularization methods. However, one of the problems with the dropout technique is that the training period is typically longer than that of a standard deep neural network architecture. This is because the parameter updates in the networks are very noisy. Moreover, the dropout technique can be considered as an alternative way of adding noise to the hidden units in the networks. This becomes a trade-off requirement between overfitting and training time. In other words, by increasing training time, one can already use a high dropout rate and encounter fewer overfitting problems for the deep neural network architectures. Thus, for future work, an interesting direction to take is to speed up the dropout technique during the training processes despite the large deep neural network architecture. Furthermore, another future direction is to use the dropout technique as an adaptive regularization for adaptive ensemble learning to further prevent overfitting, thereby enhancing the model performances of the deep neural network architectures and diagnoses of fetal health assessment with cardiotocography in clinical cases.

ACKNOWLEDGMENT

The authors would like to thank Dr. Joaquim P. Marques de Sa from the Biomedical Engineering Institute, Porto, Portugal; and Dr. Joao Bernardes and Dr. Diogo Ayres-de-Campos from the Faculty of Medicine, University of Porto, Portugal, whose Cardiotocography datasets of clinical instances were contributed to and made available in the Cardiotocography Databases of the UCI Machine Learning Repository.

REFERENCES

- [1] G. Sedgh, S. Singh, and R. Hussain, "Intended and unintended pregnancies worldwide in 2012 and recent trends," *Studies in Family Planning*, Vol. 45, Issue 3, pp. 301-314, September 2014.
- [2] C. J. Murray, "Global, regional, and national age-sex specific all-cause and cause-specific mortality for 240 causes of death, 1990-2013: a systematic analysis for the global burden of disease study 2013." *Lancet*, Vol 385, pp. 117-171, January 2015.
- [3] World Health Organization, Maternal mortality: fact sheet, Updated November 2016. Available <http://www.who.int/mediacentre/factsheets/fs348/en/>.
- [4] National Institutes of Health, What are some common complications of pregnancy? US Department of Health and Human Services: Available <https://www.nichd.nih.gov/health/topics/pregnancy/conditioninfo/Pages/complications.aspx>.
- [5] Office on Women's Health, Pregnancy Complications, US Department of Health and Human Services: Available <https://www.womenshealth.gov/pregnancy/youre-pregnant-now-what/pregnancy-complications>.
- [6] R. M. Grivell, Z. A. Gillian, M. L. Gyte, and D. Devane, "Antenatal cardiotocography for fetal assessment." *Cochrane Database of Systematic Reviews*, Issue 9. No. CD007863, pp. 1-57, 2015, John Wiley & Sons, Ltd.
- [7] C. Nelson-Piercy and C. Williamson, Medical Disorders in pregnancy, In: Chamberlain G, Steer P editor(s), *Turnbull's Obstetrics*, 3rd Edition, Edinburgh: Churchill Livingstone, pp. 275-97, 2001.
- [8] C. Lloyd, Hypertensive disorders of pregnancy, In: Fraser DM, Cooper MA editor(s), *Myles Textbook for Midwives*. 14th Edition, Edinburgh: Churchill Livingstone, pp. 357-71, 2003.
- [9] C. Lloyd, Common medical disorders associated with pregnancy, In: Fraser DM, Cooper MA editor(s), *Myles Textbook for Midwives*, 14th Edition, Edinburgh: Churchill Livingstone, pp. 321-55, 2003.
- [10] National Institute for Health and Clinical Excellence, Diabetes in pregnancy: management of diabetes and its complications from pre-conception to the postnatal period, London, pp. 1-39, March 2008.
- [11] C. Gribbin and J. Thornton, Critical evaluation of fetal assessment methods, In: James DK, Steer PJ, Weiner CP editor(s), *High Risk Pregnancy Management Options*, Elsevier, 2006.
- [12] N. M. Fisk and R. P. Smit, Fetal growth restriction; small for gestational age, In: Chamberlain G, Steer P editor(s), *Turnbull's Obstetrics*, 3rd Edition, Edinburgh: Churchill Livingstone, pp. 197-209, 2001.
- [13] I. Ingemarsson, "Fetal monitoring during labor," *Neonatology*, Vol. 95, No. 4, June 2009.
- [14] FIGO News, "Report of the FIGO study group on the assessment of new technology: evaluation and standardization of fetal monitoring," Organized by G. Rooth, A. Huch, and R. Huch, *International Journal of Gynecology & Obstetrics*, Vol. 25, pp. 159-167, 1987.
- [15] G. J. Miao, K. H. Miao, and J. H. Miao, "Neural pattern recognition model for breast cancer diagnosis," *Multidisciplinary Journals in Science and Technology, Journal of Selected Areas in Bioinformatics*, August Edition, pp. 1-8, September 2012.
- [16] K. H. Miao, J. H. Miao, and G. J. Miao, "Diagnosing coronary heart disease using ensemble machine learning," *International Journal of Advanced Computer Science and Applications*, Vol. 7, No. 10, pp. 30-39, 2016.
- [17] K. H. Miao and G. J. Miao, "Mammographic diagnosis for breast cancer biopsy predictions using neural network classification model and receiver operating characteristic (ROC) curve evaluation," *Multidisciplinary Journals in Science and Technology, Journal of Selected Areas in Bioinformatics*, September Edition, Vol. 3, Issue 9, pp. 1-10, October 2013.
- [18] J. H. Miao, K. H. Miao, and G. J. Miao, "Breast cancer biopsy predictions based on mammographic diagnosis using Support Vector Machine learning," *Multidisciplinary Journals in Science and Technology, Journal of Selected Areas in Bioinformatics*, Vol. 5, No. 4, pp. 1-9, 2015.
- [19] D. Ayres-de-Campos, J. Bernardes, A. Garrido, J. Marques-de-Sá, and L. Pereira-Leite, "Sisporto 2.0: A program for automated analysis of cardiotocograms," *The Journal of Maternal-Fetal Medicine*, Vol. 9, Issue 5, pp. 311-318, October 2000.
- [20] P. A. Warrick, E. F. Hamilton, R. E. Kearney, and D. Precup, "A machine learning approach to the detection of fetal hypoxia during labor and delivery," *Proceedings of the Twenty-Second Innovative Applications of Artificial Intelligence Conference*, pp. 1865-1870, 2010.
- [21] Z. Comert and A. F. Kocamaz, "Comparison of machine learning techniques for fetal heart rate classification," *Special issue of the 3rd International Conference on Computational and Experimental Science and Engineering*, Vol. 132, pp. 451-454, 2017.
- [22] C. Sundar, M. Chitradevi, and G. Geetharamani, "Classification of cardiotocogram data using neural network based machine learning technique," *International Journal of Computer Applications*, Vol. 47, No. 14, pp. 19-25, June 2012.
- [23] M. Arif, "Classification of cardiotocograms using Random Forest classifier and selection of important features from cardiotocogram

- signal," *Biomaterials and Biomedical Engineering*, Vol. 2, No. 3, pp. 173-183, 2015.
- [24] The UCI Machine Learning Repository, Cardiotocography Data Set: Available <http://archive.ics.uci.edu/ml/datasets/Cardiotocography>.
- [25] J. Wang, "Deep learning: an artificial intelligence revolution," *ARK Invest*, pp. 1-41, New York, June 2017.
- [26] G. E. Hinton, S. Osindero, and Y. W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, Vol. 18, pp. 1527-1554, 2006.
- [27] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural Computation*, Vol. 14, No. 8, pp. 1711-1800, 2002.
- [28] G. E. Hinton, *A Practical Guide to Training Restricted Boltzmann Machines*, Department of Computer Science, University of Toronto, Canada, UTML TR 2010-003, Version 1, August 2010.
- [29] Y. Bengio, "Learning deep architectures for AI," *Journal Foundations and Trends in Machine Learning*, Vol. 2, pp. 1-127, January 2009.
- [30] C. Y. Liou, J. C. Huang, and W. C. Yang, "Modeling word perception using the Elman network," *Neurocomputing*, pp. 3150-3157, 2008.
- [31] P. Baldi, "Autoencoders, unsupervised learning, and deep architectures," *JMLR: Workshop and Conference Proceedings*, pp. 37-50, 2002.
- [32] S. Haykin, *Neural Network: A Comprehensive Foundation*, Macmillan College Publishing Company, 1994.
- [33] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *Neural and Evolutionary Computing*, pp. 1-18, July 2012.
- [34] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, Vol. 15, pp. 1929-1958, 2014.
- [35] A. Livnat, C. Papadimitriou, N. Pippenger, and M. W. Feldman, "Sex, mixability, and modularity," *Proceedings of the National Academy of Sciences*, Vol. 107, No. 4, pp. 1452-1457, 2010.
- [36] D. Warde-Farley, I. J. Goodfellow, A. Courville, and Y. Bengio, "An empirical analysis of dropout in piecewise linear networks," *Machine Learning*, pp. 1-10, January 2014.
- [37] I. J. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, MIT Press, 2016.
- [38] J. Kukačka, V. Golkov, and D. Cremers, "Regularization for deep learning: a taxonomy," *Artificial Intelligence*, pp 1-27, October 2017.
- [39] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, Vol. 9, Chia Laguna Resort, Sardinia, Italy, 2010.
- [40] C. R. Yali, G. Nallamala, W. Fedus, and Y. Prabhuzantye, "Efficient encoding using deep neural networks," pp. 1-8, 2015.
- [41] D. Erhan, Y. Bengio, A. Courville, P. Manzagol, P. Vincent, and S. Bengio, "Why does unsupervised pre-training help deep learning?" *Journal of Machine Learning Research*, Vol. 11, pp. 625-660, 2010.
- [42] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," In *Advances in Neural Information Processing Systems*, Vol. 12, pp. 153-160, 2007.
- [43] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: a review and new perspectives," *Department of computer science and operations research, University of Montreal, Canadian Institute for Advanced Research*, pp. 1-30, April 2014.
- [44] M. Ranzato, C. Poultney, S. Chopra, and Y. LeCun, "Efficient learning of sparse representations with an energy-based model," *Proceedings of the 19th International Conference on Neural Information Processing Systems*, pp. 1137-1144, 2006.
- [45] G. J. Miao and M. A. Clements, *Digital Signal Processing and Statistical Classification*, Artech House, Inc., 2002.
- [46] M. L. Zhang and Z. H. Zhou, "A review on multi-label learning algorithms," *IEEE Transactions on Knowledge and Data Engineering*, Vol. 26, Issue 8, pp. 1819-1837, March 2013.
- [47] E. Gibaja and S. Ventura, "A tutorial on multi-label learning," *ACM Computing Surveys*, Vol. 47, Issue 3, April 2015.
- [48] G. Hinton, with N. Srivastava, and K. Swersky, "Neural network for machine learning: Lecture 6a overview of mini-batch gradient descent," *Computer Science Department, University of Toronto*, Winter 2014.

Flow-Length Aware Cache Replacement Policy for Packet Processing Cache

Hayato Yamaki

Dept. of Computer and Network Engineering
The University of Electro-Communications
Chofu, Japan

Abstract—Recent core routers are required to process packets not only at high throughput but also with low power consumption due to the increase in the network traffic amount. Packet processing cache (PPC) is one of the effective approaches to meet the requirements. PPC enables to process a packet without accessing to a ternary content addressable memory (TCAM) by storing the TCAM lookup results of a flow in a cache. Because the cache miss rate of PPC directly impacts on the packet processing throughput and the power consumption of core routers, it is important for PPC to reduce the number of cache misses. In this study, we focus on characteristics of flows and propose an effective cache replacement policy for PPC. The proposed policy, named Hit Dominance Cache (HDC), divides the cache into two areas and assigns flows to the appropriate area to evict mice flows rapidly and to retain elephant flows preferentially. Simulation results with 15 real network traces show that HDC can reduce the number of cache misses in PPC by up to 29.1% and 12.5% on average when compared to 4-way LRU, conventionally used in PPC. Furthermore, the hardware implementation using Verilog-HDL shows that the hardware costs of HDC is comparable to those of 4-way LRU though HDC performs as if the cache was composed of 8-way set associativity. Finally, we show that HDC can achieve 503 Gbps with 88.8% energy of conventional PPC (20.5% energy of TCAM only architecture).

Keywords—Router; packet processing; cache replacement

I. INTRODUCTION

Internet traffic has increased year by year due to the popularization of internet applications which generate a large number of packets, such as file sharing, cloud services, and video streaming. Because the traffic concentrates on routers, the processing load of routers has increased and becomes a serious problem. According to The Ministry in Japan [1], [2], it is reported that the total amount of internet traffic and the power consumption of network devices will increase approximately 190 times and 10 times, respectively, in 2025 compared to 2006. The power consumption of routers is no longer negligible because it will account for several percentages of total power consumption in the world [3], [4]. Not only high throughput but also low power consumption is required for routers and especially for core routers, which handle the huge traffic close to the center of the internet.

In a core router, table lookups for packet processing is known as the main cause of both degrading the throughput and

consuming the power [5]-[7]. To determine how to process a packet (i.e., where to transmit or how to filter the packet, etc.), routers are required to retrieve tables such as the routing table, the address resolution protocol (ARP) table, the access control list (ACL), and the quality of service (QoS) table. In recent core routers, these tables are stored in a ternary content addressable memory (TCAM), which is a memory specialized in high-speed data search. While the TCAM can obtain a table lookup result with one cycle, it consumes approximately 16 times as large power as a same sized static random access memory (SRAM) [8]. Due to this, it is indicated that the TCAM accounts for 40% of all power consumed in a router [9], [10]. Make matter worse, the lookup performance of the TCAM cannot reach the throughput required for future internet (i.e., more than 400 Gbps) because of the low operation frequency. Thus, improvement of the TCAM lookups is important for future core routers to achieve both high throughput and low power consumption.

As one of the solutions, optimizing the TCAM use is the most popular approach. Nawa et al. proposed a novel searching scheme for the TCAM [9]. They enabled to reduce the dynamic energy of the TCAM lookups by dividing the all TCAM entries into several groups and searching only appropriate group. Gamage et al. proposed a method for high-throughput table lookups with parallelized TCAMs [5]. The proposed method enabled to accelerate the packet processing by assigning packets to the suitable TCAMs. However, it is still difficult to achieve more than 400-Gbps throughput and SRAM-like power. Other approaches are required for further improvement.

Packet processing cache (PPC) is another approach and recently reevaluated because it can improve the table lookup performance without impeding the TCAM-based approaches and thus adopt concurrently with them. PPC includes a cache which retains the TCAM lookup results of packets and reuses them to process following packets. If the TCAM lookup result of a packet is in the cache, PPC can process the packet without accessing the TCAM using the cache. The more PPC can process packets with the cache, the larger PPC can acquire the throughput and the power reduction because of a small number of TCAM accesses. In other words, the performance of PPC depends on the cache hit/miss rate. Thus, to reduce the number of cache misses is an important issue for PPC. In this study, we investigate causes of the cache misses in PPC and propose an effective cache replacement policy for PPC to reduce the number of cache misses.

This work was supported by KAKENHI 18K18022, a research grant from The Mazda Foundation, and a research grant from Sumitomo Foundation.

The contribution of this paper is summarized as follows:

- Major Causes of the cache misses in PPC are suggested. We show that two types of flows make a large number of cache misses.
- Our simulation shows that Hit Dominance Cache (HDC), proposed in this paper, reduces the number of cache misses in PPC by up to 29.1% (12.5% on average) with comparable hardware cost to 4-way Least Recently Used (LRU), typically used in PPC.
- The performance of cache replacement policies to various types of traffic patterns is evaluated. Although the cache access patterns in PPC (i.e., behavior of packets) differ depending on the network structure, previous studies simulated with only a few network traces. This study uses 15 network traces for evaluation.
- This paper is an expansion of the paper [11], which published in the proceedings of Future of Information and Communication Conference (FICC) 2018. It newly reveals more detailed relation between flows and the PPC cache misses. In addition, the performance difference of variable HDC designs and the more detailed hardware costs are also evaluated.

The rest of this paper is organized as follows. We first show the more details of PPC in Section 2 and introduce the related works of reducing the number of cache misses in PPC in Section 3. After that, we investigate major causes of the cache misses in PPC and propose our technique HDC in Section 4. Section 5 evaluates HDC performance from the aspects of the cache miss reduction, the implementation costs, the throughput, and the energy consumption. Finally, we conclude this paper in Section 6.

II. PACKET PROCESSING CACHE

In the routers, the TCAM lookup results depend on several fields in the packet header (e.g., IP addresses) because they are used as a key for retrieval. In particular, the five-tuple (i.e., the source and destination IP addresses, the source and destination port numbers, and the protocol number) are used in most tables in a router. Based on this fact, PPC defines packets which have the same five-tuple as a flow and stores the TCAM lookup results of the first packets of flows to a cache memory. If the TCAM lookup result of a flow is in the cache, PPC can process the subsequent packets of the flow without accessing the TCAM using the cached result and reduce the number of TCAM accesses. Because the packets of the same flow arrive in a router at short intervals [12], [13], routers can benefit from the cache. The outline of the packet processing by PPC is shown in Fig. 1. The latency and the energy consumption of the cache are considerably lower than those of the TCAM, and therefore PPC can process packets at high speed with low energy consumption if the cache hits occur.

The cache is composed of 13-byte flow information as a cache tag and 15+-byte TCAM lookup results as cache data. The cache data include 1 byte as a result of the routing table lookup (output interface number), 12 bytes as a result of the ARP table lookup (MAC address), 1 byte as a result of the ACL (filtering decision), 1 byte as a result of the QoS table (priority value). Moreover, the cache can also store other processing results such as filtering results of NIDS (Network Intrusion Detection System), encapsulation results, and encryption results by expanding the cache data field. PPC can perform many packet processing with one cache access.

In PPC, TCAM accesses are required only when the TCAM lookup results of the corresponding flows are not in the cache. For this reason, the performance of PPC depends on the cache miss rate. We define the throughput and the energy consumption of the table lookups per packet with PPC as T_{PPC} and E_{PPC} , respectively. These variables are calculated as:

$$T_{PPC} = \begin{cases} \frac{1}{l_{Cache}} \cdot 64 \text{ byte} & (l_{Cache} > l_{TCAM} \cdot m) \\ \frac{1}{l_{TCAM} \cdot m} \cdot 64 \text{ byte} & (l_{Cache} < l_{TCAM} \cdot m) \end{cases} \quad (1)$$

$$E_{PPC} = (DE_{Cache} + DE_{TCAM} \cdot n \cdot m) + (SE_{Cache} + SE_{TCAM}). \quad (2)$$

Here, l_{Cache} and l_{TCAM} represent the latencies of the cache and the TCAM, respectively. Likewise, DE_{Cache} and DE_{TCAM} represent the dynamic energy of the cache and the TCAM per access, respectively, while SE_{Cache} and SE_{TCAM} represent the static power of them. Conventionally, the impact of the static power in the memories can be ignored because the dynamic energy dominates the total energy consumed by a memory in a router [14].

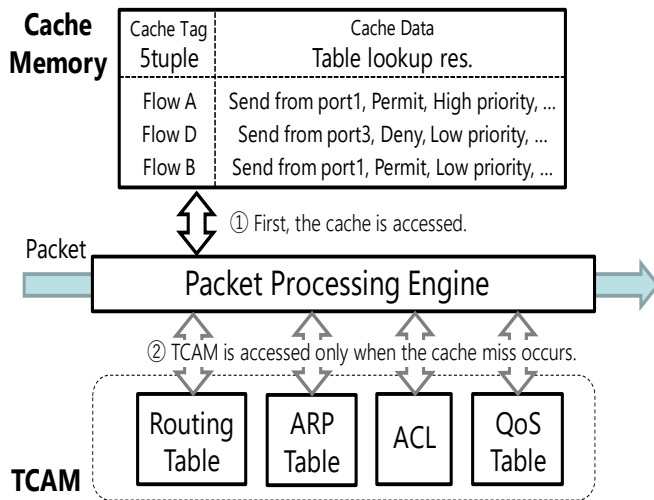


Fig. 1. Outline of packet processing by PPC.

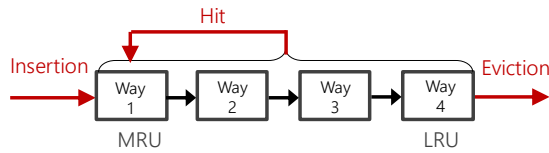


Fig. 2. Outline of LRU entry replacement.

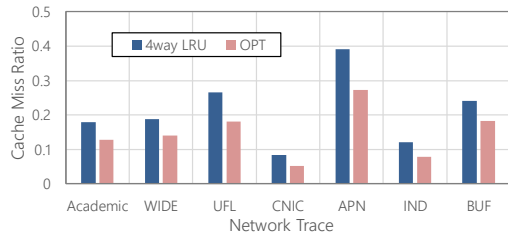


Fig. 3. Comparison of cache miss rates between LRU and OPT.

However, we consider both of energies in this paper because the static power is not negligible in case of a small cache memory. m represents the cache miss rate. n represents the number of TCAM accesses needed to process a packet. In this paper, we suppose $n = 4$ because popular routers have four types of tables as shown in Fig. 1. 64 bytes referred in (1) represents the shortest packet length. Equation (1) means that the table lookup throughput with PPC is limited by the smaller throughput of the cache or the TCAM. These equations indicate that the throughput and the power consumption with PPC are mainly determined by the cache miss rate due to the large latency and high energy consumption of the TCAM. Thus, to achieve low cache miss rate is important to improve both the throughput and the power consumption with PPC.

To increase the total number of cache entries is the most simple and effective solution to reduce the number of cache misses for many cache systems. However, for PPC, this approach leads to the increase in the cache capacity easily due to the large entry size (28 bytes per entry). As a result, it increases in both the latency and the power consumption of the cache. Unlike processor caches, the gap of the latencies between the cache and the TCAM is small. Thus, it is desirable for PPC to use a small cache memory such as level-1 (L1) caches in processors. In this paper, we aim to reduce the number of cache misses in L1-sized PPC without increasing the cache capacity by improving the use of the cache space.

Conventionally, PPC is designed by a 4-way set associative cache from the aspects of the cache miss rate and the hardware complexity [15]-[17]. It means that each cache line has four entries, and PPC can keep useful entries in the cache by applying a suitable cache replacement policy. The cache replacement policy is used to determine which entry should be replaced when the cache line is full. It is known that the cache replacement policy impacts on the cache miss rate significantly. In PPC, Least Recently Used (LRU), whose concept is to evict the entry hit oldest, is empirically used as the cache replacement policy. As shown in Fig. 2, LRU inserts a new entry into the Most Recently Used (MRU) position and evicts it from the Least Recently Used (LRU) position. In addition, when an entry is referenced, LRU shifts it to the MRU position. LRU performs good in many cache systems because it can

utilize the temporal locality of data; however, it is not certain that 4-way LRU is suitable for PPC. Fig. 3 shows the difference of the cache miss rates measured by a simulation with real network traffics. The details of the simulation are described in Section 5. In this simulation, we designed 4-way LRU and optimal page replacement algorithm (OPT) [18] as the cache replacement policy in PPC. OPT is an ideal replacement policy which uses the information of all future arrived data. Thus, the performance of OPT is the best of all cache replacement policies; however, it cannot be implemented for practical use in most cases. Fig. 3 indicates that approximately 30% of all the cache misses in PPC have the opportunities to be reduced by replacing entries more effectively than LRU. In this study, in order to close the gap of the cache miss rates between LRU and OPT, we consider the effective replacement policy.

III. RELATED WORK

Various cache replacement policies have been proposed in previous studies for many cache systems. However, they are not always effective for PPC because of the difference of the cache access patterns. In this section, we introduce a few studies focusing on the reduction in the number of cache misses in PPC and reveal the problems of the proposed methods.

Chang et al. pointed out that PPC cannot prepare a large number of cache entries due to the large tag size and proposed a method to compress the cache tag [19]. They used a 32-bit hash value calculated from the five-tuple as the cache tag instead of the 104-bit flow information. However, adding extra hardware is needed to avoid the conflicts of the hash values. Similarly, Ata et al. compressed the cache tag of PPC by using only three fields of the five-tuple: the source and destination IP addresses and the smaller number of ports [20]. However, it cannot meet demands of recent routers. For example, the QoS table requires the five-tuple to determine the QoS value. Compressing the cache tag sacrifices the information stored in the cache or requires to add extra hardware.

Li et al. discussed the appropriate cache design for PPC [17]. In order to use the cache space efficiently, they focused on three viewpoints: the cache associativity, the cache replacement policy, and the hash function. In [17], the authors concluded that a 4-way set associativity with LRU is the best design from the balance between the implementation costs and the cache hit rate. Additionally, they show that the difference of the hash functions does not impact on the cache hit rate largely. While the authors evaluated three policies (LRU, least frequently used (LFU), and round robin), other policies were not considered.

Kim et al. proposed an effective cache replacement policy for PPC [21]. They indicated that LRU was not suitable for PPC because LRU focused on only temporal locality and cannot utilize activities of networks. The proposed policy classifies the cache entries into two types: a switching entry and a non-switching entry. The switching entry is the entry hit at least once; non-switching entry is the entry never hit. The entry is replaced from the non-switching entries. Furthermore, they proposed two types of cache replacement policies called Weighted Priority LRU Scheme and L2A Cache Scheme. In

Weighted Priority LRU Scheme, the non-switching entries cannot be replaced until a threshold time passes because it is expected that the non-switching entries are referenced again in a short time. In L2A Cache Scheme, the replaced entry is decided by the amount of the timestamp values in last two packets. L2A Cache Scheme can reduce the number of cache misses compared to LRU in case that the cache size is small. However, concrete hardware requirements, such as the number of bits for storing the timestamp to the cache and the way for getting the time, were not referred. The increase in the memory costs becomes a critical problem especially in PPC.

Yamaki et al. considered the methods of denying the cache registration of one-packet flows because they have no opportunities to hit in the cache [22], [23]. They proposed several methods to specify applications which create flows composed of only one packet, such as flows created by domain name system (DNS) or by several types of network attacks, and not to store these flows in the cache. While the proposed methods can reduce the number of cache misses in PPC by approximately 8%, the misses caused by various small factors cannot be improved because the methods handle the misses of only specific applications.

IV. CACHE REPLACEMENT POLICY

In this paper, we focus on the cache replacement policy because it has a potential to reduce a large number of cache misses as shown in Fig. 2 with small hardware modification. In this section, we first analyze the flow behavior in PPC to reveal the cause of the cache misses in PPC. After that, we propose a novel cache replacement policy based on the above analysis.

A. Analysis of Flow Behavior in PPC

Flows composed of a few packets (referred to as mice flows) are one of the main causes of increasing the cache misses in PPC [23]. The mice flows cause cache pollution due to the occupation of entries in spite of a few cache hits. In particular, flows composed of one packet are not needed for the cache because they never hit in the cache. In [23], it was indicated that 99% of all flows are the mice flows composed of less than 10 packets. Moreover, the flows composed of one packet accounts for about half of all flows in networks. Therefore, mice flows occupy most entries in the cache and impact on the cache performance significantly. It is important for PPC to evict the mice flows from the cache rapidly.

Unlike the mice flows, there are flows composed of a large number of packets (referred to as elephant flows), such as video flows. These elephant flows are composed of more than 1,000 packets. Although the number of elephant flows are few, they have great impact on the cache miss rate because of a large number of cache references. These trends that a large number of mice flows and a few number of elephant flows account for most packets in networks is known as the elephant and mice phenomenon [24]. We analyze the behavior of the elephant flows in the cache and show the examples in Fig. 4. This graph shows logs of the cache hits and misses in elephant flows. Contrary to expectation, many packets in an elephant flow make the cache misses. It means that entries of the elephant flows are replaced repeatedly though they are referenced many times.

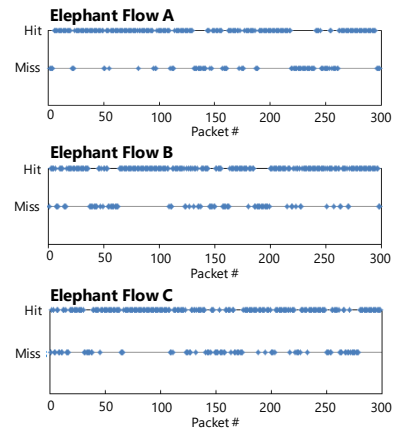


Fig. 4. Logs of the cache hits and misses in elephant flows.

We also investigate the impact of the mice flows and elephant flows on the number of cache misses. Fig. 5 shows the total amount of packets, compulsory misses, and other misses in the flows composed of the specific number of packets. Compulsory misses shown in the figure are misses caused by first packets of flows, and thus it is difficult for PPC to prevent them. In contrast, other misses shown in the figure can be reduced by retaining entries of the corresponding flows appropriately. This figure indicates that a large number of misses, especially compulsory misses, are caused by mice flows, and it pollutes the cache. In addition, we show the ratios of each cache miss to all packets in Fig. 6. As shown in Fig. 4, it also indicates that packets in elephant flows cause the cache misses at a higher rate than we expect. Thus, it has better for PPC to prevent the cache pollution caused by mice flows and to prioritize elephant flows.

B. Hit Dominance Cache

Based on above considerations, we propose Hit Dominance Cache (HDC) to prevent the cache pollution caused by mice flows and retain elephant flows preferentially. The main concept of HDC is to provide difference cache priorities to elephant flows and mice flows by assigning difference cache areas to them. Because it is difficult to identify whether a packet is belonging to an elephant flow or a mice flow accurately when the packet comes, HDC judges it from the number of cache hits. HDC prioritizes entries which hit many times as elephantish flows.

Fig. 7 shows the outline of HDC entry replacement. In HDC, the cache area is divided into two areas: the hit area and the primary area. Entries are inserted in and evicted from the LRU position of the primary area. Thereby, new entries can be evicted by one replacement at the shortest. It enables to evict the mice flows rapidly. When an entry in each area is referenced, the entry is shifted to the MRU position in each area. In addition, one notable behavior of HDC is that the entry in the primary area is swapped for the LRU position entry in the hit area if the entry is referenced threshold times (we define it as the HDC threshold). It enables to retain elephant flows preferentially because mice flows disturb entries in only the primary area. In Fig. 7, HDC is depicted as the combination of 4-way hit area and 4-way primary area; however, each area can be designed variably.

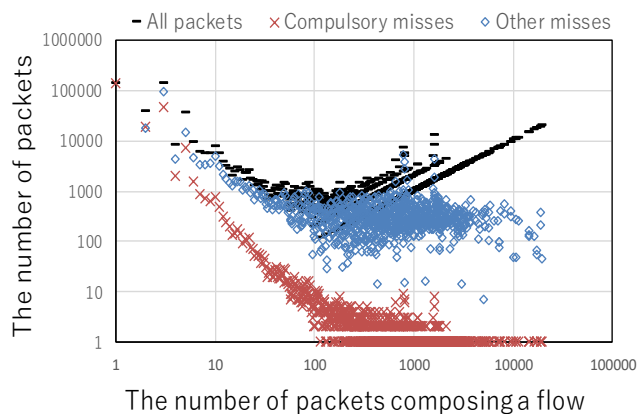


Fig. 5. Total amount of all packets, compulsory misses, and other misses in the flows composed of the specific number of packets.

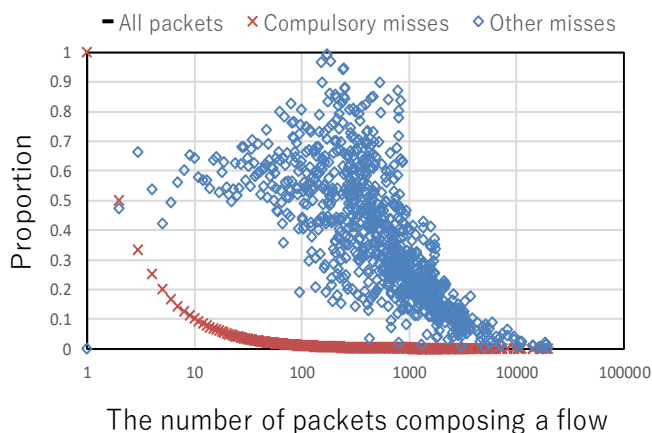


Fig. 6. Average ratios of two types of cache misses to all packets in the corresponding flow. For example, this graph shows all packets in the flows composed of one packet cause the compulsory misses.

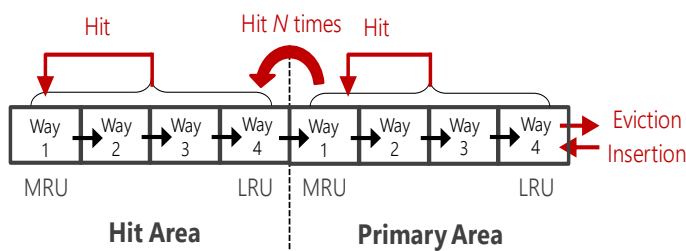


Fig. 7. Outline of HDC entry replacement.

V. EVALUATION

This section evaluates the usefulness of HDC. In this paper, we simulated the packet processing with PPC including HDC using a software PPC simulator and implemented HDC using a hardware description language to evaluate HDC from following three aspects:

- Cache miss reduction
- Implementation costs
- Overall throughput and energy consumption.

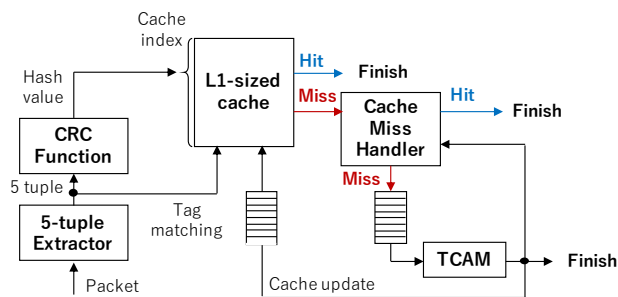


Fig. 8. Block diagram of PPC simulator.

A. Cache Miss Reduction

To measure the cache miss rate of PPC, we prepared a cycle-based PPC simulator written in C++. The block diagram of the PPC simulator is shown in Fig. 8. The simulator models the PPC behavior such as reading packets, extracting the flow information, calculating the cache index, referring to and updating the cache, and referring to the TCAM. First, pcap-format file is read, and a packet is extracted based on the timestamp value. After that, the cache index is calculated from the five-tuple of the packet using CRC hash function. Next, the cache is referenced using the cache index. If the packet hits in the cache, the PPC simulator finishes processing the packet and extracts the next packet. On the other hand, if the packet misses in the cache, the packet is forwarded to Cache Miss Handler (CMH). CMH is a module to prevent the cache misses caused by the time-lag between the cache reference and the cache update. Because it takes a little time to insert a new entry in the cache, cache misses may be occurred if subsequent packets of the same flow come before updating the cache. CMH manages the flows just processed in the TCAM and stores the subsequent packets of the same flows until the corresponding entry is prepared in the cache. More details of CMH is described in [25]. If the packet misses in CMH, the packet is forwarded to the TCAM module. After the TCAM access latency, the PPC simulator finishes processing the packet and updates the cache.

Table I shows the parameters of the PPC simulator. The cache was estimated as an L1-sized cache. The latencies of the cache and the TCAM were set to 0.5 ns and 5 ns, respectively. Note that we assumed the sizes of CMH and queues shown in Fig. 8 were enough large, and the simulator can process packets without any packet losses. Besides, we used 15 types of network traces shown in Table II as workloads to reveal the performance of the cache replacement policies without depending on the network traffic patterns. The network traces were acquired from RIPE Network Coordination Centre [26] and Widely Integrated Distributed Environment (WIDE) [27]. Furthermore, an academic trace acquired from a core network in Japan was used as a high-bandwidth workload.

First, the suitable design for each area of HDC was evaluated. As described in Section 4, HDC can variably design the number of associativity sets in each area. We implemented 2-way hit area 4-way primary area HDC (2-4 HDC), 4-way hit area 2-way primary area HDC (4-2 HDC), and 4-way hit area 4-way primary area HDC (4-4 HDC) and compared the cache miss rates of them. In this simulation, each HDC design had

the same total number of entries and set the HDC threshold = 8. Note that 8-way design was not evaluated in this simulation because it cannot be implemented for practical use due to the high implementation costs as described later. Fig. 9 shows the cache miss rates of 4-way LRU and each HDC design. In this figure, we showed the results of only three networks (i.e., TXG, IPLS, and UFL) because the trends of all the results are mostly the same. Fig. 9 indicates that 4-4 HDC is the best design to achieve the low cache miss rate. 4-4 HDC can reduce the number of cache misses by up to 20.8% and 16.7%, when compared to 2-4 HDC and 4-2 HDC, respectively. We consider it is because the total amounts of the mice flow packets and the elephant flow packets are almost the same as shown in Fig. 5, and one to one design is fitting to split these flows symmetrically. From this result, we adopt 4-4 HDC design hereafter.

Next, the impact of the HDC threshold was evaluated. Figure 10 shows the difference of the cache miss rates among 4-way LRU and HDCs with various HDC thresholds. We set the HDC threshold to 1, 2, 4, 8, and 16 and represented them from HDC 1 to HDC 16 in the figure. As well as Fig. 9, we showed the results of only three networks because of the similar trend. This result indicates that it is the best to set the HDC threshold to eight.

TABLE I. PARAMETERS OF PPC SIMULATOR

Parameter	Value
Total cache entries	1,024 entries
Latency of the cache	0.5 ns
Latency of the TCAM	5 ns

TABLE II. DETAILS OF NETWORK TRACES

Trace	Captured date	Average # of packets [pps (packets per second)]
IND [26]	2003/1/6	15,540
BUF [26]	2003/1/18	8,380
TXG [26]	2004/3/26	12,475
APN [26]	2004/3/26	20,330
IPLS3 [26]	2004/6/1	116,778
BWY [26]	2004/10/7	17,922
COS [26]	2005/1/8	8,051
CNIC [26]	2005/3/17	28,440
MRA [26]	2005/3/21	41,372
UFL [26]	2005/3/21	50,769
FRG [26]	2006/1/10	32,955
PSC [26]	2006/2/20	26,912
PUR [26]	2006/2/20	42,515
Academic	2010/6/17	371,013
WIDE [27]	2017/4/12	58,776

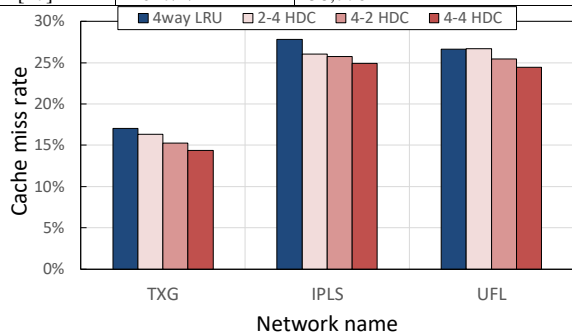


Fig. 9. Cache miss rates of various HDC designs.

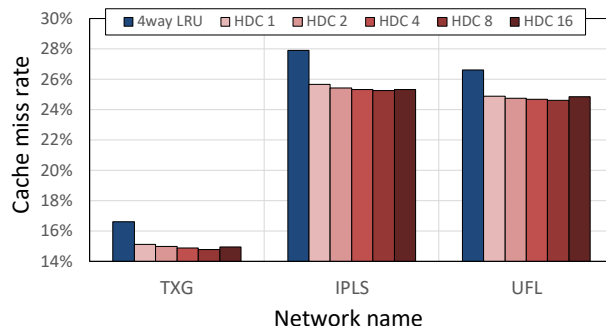


Fig. 10. Cache miss rates of various HDC thresholds.

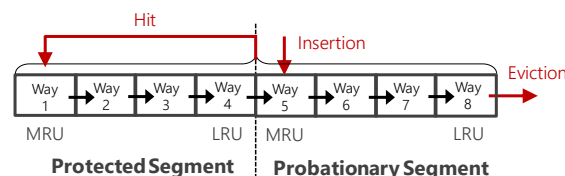


Fig. 11. Outline of SLRU entry replacement.

Finally, the cache miss reduction of HDC was evaluated. In this evaluation, 4-way LRU, 8-way LRU, OPT, and segmented LRU (SLRU) [28] were used as the cache replacement policies for comparison. SLRU is a cache replacement policy which resembles HDC and divides the cache area into two areas: the probationary segment and the protected segment. Fig. 11 shows the outline of SLRU entry replacement. In SLRU, a new entry is inserted into the MRU position of the probationary segment and evicted from the LRU position of the probationary segment. When an entry is referenced, the entry is set on the MRU position of the protected segment. Although SLRU divides the cache into two areas, it is the same as 8-way LRU whose inserted position of a new entry is changed to the middle of the 8-way entries.

Fig. 12 and 13 show the cache miss rates of various LRUs, HDC, and OPT with 15 types of network traces and the improvement ratios of them to 4-way LRU. The results showed that HDC performed the best in nine networks. HDC can reduce the cache misses by up to 29.1% (12.5% on average) compared to 4-way LRU, while SLRU can reduce the cache misses by up to 20.2% (11.1% on average) compared to 4-way LRU. SLRU performed better than HDC in six network traces; however, SLRU is not suitable for practical use because of the high implementation cost like 8-way LRU, as mentioned later. By contrast, 8-way LRU cannot achieve major improvement to 4-way LRU though 8-way LRU can achieve considerably lower cache miss rate than 4-way LRU in many other cache systems. It is because 8-way LRU cannot evict mice flows rapidly, and thus the cache entries are polluted by them. In this simulation, HDC provided up to 84.6% of the OPT performance; however, the average performance of HDC is 44.5% of the OPT performance. It means that there is still room for improvement.

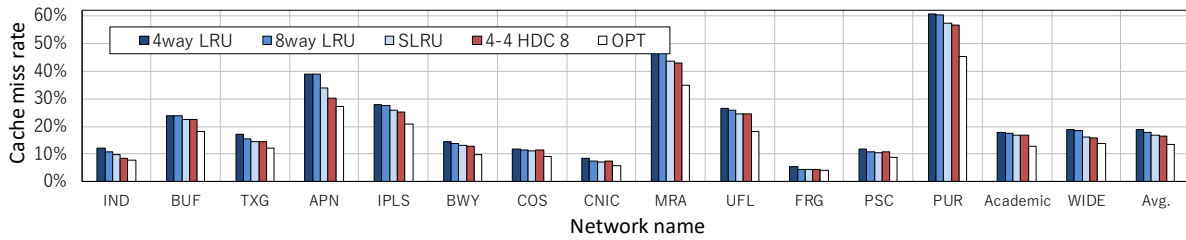


Fig. 12. Cache miss rates of various LRUs, HDC, and OPT in 15 types of network traces.

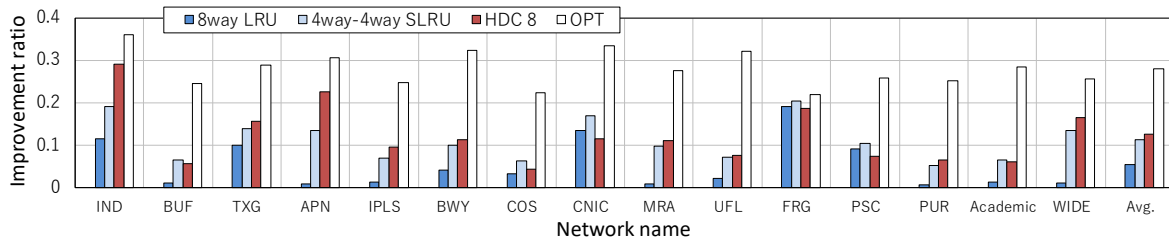


Fig. 13. Improvement ratios of the cache misses of various LRUs, HDC, and OPT to 4-way LRU.

B. Implementation Costs

When we consider a cache replacement policy, not only the cache miss reduction but also the implementation costs are an important issue because the cache replacement policy requires a large circuit area which cannot be ignored in some cases. In general, two modules are required to implement a cache replacement policy. One is a memory to store the replacement priorities of entries per cache line, and the other is a module to update these priorities when a cache hit or a new entry insertion occurs.

In the case of 4-way LRU, 5 bits are required per cache line to handle 24 possible combinations of the replacement priorities (because each entry is ranked from 1 to 4). On the other hand, 8-way LRU needs 16 bits per cache line to handle 40,320 possible combinations of those (because each entry is ranked from 1 to 8). Furthermore, the module to update the replacement priorities in the 8-way LRU is significant larger than those in the 4-way LRU because there are $40,320 * 7$ possible patterns of the replacement priority transitions in 8-way LRU though 4-way LRU needs $24 * 3$ possible patterns of the replacement priority transitions. As a result, it is not realistic to implement 8-way LRU on hardware. Similarly, SLRU cannot be implemented with realistic hardware costs.

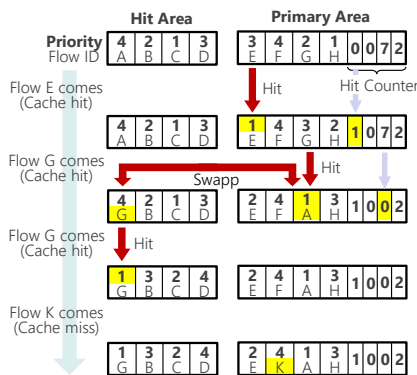


Fig. 14. An example of the method to manage the replacement priorities in 4-4 HDC 8.

Against above consideration, although HDC performs as if the cache was composed of 8-way set associativity, it enables to manage the replacement priorities in the same way as the 4-way LRU. Fig. 14 shows an example of how to manage the replacement priorities in 4-4 HDC. In order to swap the entries between the primary area and the hit area, hit counters are needed for each entry in the primary area to count the number of references. In the case of HDC 8, $3 \text{ bit} * (1,024 \text{ entries} / 2)$, namely 1.5K bits are needed as the hit counter. When a hit counter overflows, the corresponding entry is swapped for the LRU position entry in the hit area of the same cache line. At this time, the replacement priorities are updated in only the primary area because the replacement priority of the swapped entry in the hit area is not changed and keeps LRU position. Furthermore, unlike LRU, HDC does not need to update the replacement priorities when a new entry is inserted into the primary area.

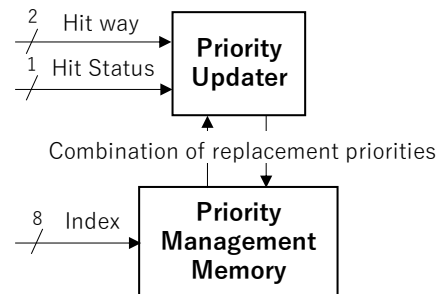


Fig. 15. Hardware architecture of general cache replacement policies.

TABLE III. SIMULATION ENVIRONMENT

Item	Tool Name
Hardware description language	Verilog-HDL
Logic simulation	Cadence NC-Verilog LDV5.7
ASIC synthesis	Synopsys Design Compiler X-2005.09
Libraries for ASIC synthesis	Free PDK OSU Library (45nm) [29]

TABLE IV. SYNTHESIS RESULTS OF 4-WAY LRU AND HDC.

	4-way LRU	HDC
Combinational circuit area	110.66 μm^2	80.332 μm^2
Memory requirement	1,280 bit	2,816 bit

TABLE V. LATENCIES AND ENERGIES OF CACHE AND TCAM

	Cache	TCAM
Latency	0.598 ns	5 ns
Dynamic energy	0.0539 nJ / access	30 nJ / access
Static power	0.0159 J/s	0.85 J/s

TABLE VI. THROUGHPUT AND ENERGY OF TABLE LOOKUPS

	Only TCAM	PPC with 4-way LRU	PPC with HDC
Throughput	102 Gbps	445 Gbps	503 Gbps
Energy per packet	124 nJ	28.6 nJ	25.4 nJ

In order to evaluate the implementation costs of HDC, the hardware implementation of HDC and LRU was simulated using a hardware description language. Fig. 15 and Table III show the hardware architecture of the cache replacement policies and the tools used for the evaluation, respectively. The priority updater shown in Fig. 14 receives the cache-hit status and the hit way number from the cache and the combination of the replacement priorities from the priority management memory when a cache hit occurs. After updating the replacement priorities, the priority updater writes them back to the priority management memory. The index of the priority management memory is the same as that of the cache in PPC.

Table IV shows the ASIC synthesis results of 4-way LRU and HDC. Note that we did not implement SLRU and 8-way LRU in this simulation because the hardware costs of them were obviously oversized for practical use. Table IV indicates that HDC can be implemented with 72.6% of the circuit area of 4-way LRU. It is because HDC does not need to update the replacement priorities when a new entry is inserted. On the other hand, the priority management memory size of HDC is 2.2 times as large as that of 4-way LRU. However, this increase is negligible because the priority management memory is small when compared to the cache. As a result, the implementation costs of HDC are comparable to 4-way LRU.

C. Overall Throughput and Energy Consumption

We finally estimated the overall throughput and the energy consumption of the table lookups. The throughput and the energy consumption with PPC can be calculated from (1) and (2), introduced in Section 2. Here, the latency, the dynamic energy, and the static power of the cache were estimated using a cache model CACTI 6.5 [30]; those of the TCAM were estimated using a TCAM power and timing model [8] (1 Mbit TCAM with 70 nm process was assumed). We show each estimated value in Table V. It indicates that both energies of the cache are remarkably smaller than those of the TCAM, and thus the cache miss of PPC significantly impacts on the table lookup performance.

Taking these estimations, we calculated the throughput and the energy of the table lookups and showed the result in Table VI. It was shown that PPC with HDC can achieve 503 Gbps throughput with 25.4 nJ energy per packet. It is 4.93

times high throughput and 20.4% energy, when compared to the conventional TCAM only architecture. In addition, when compared to PPC with 4-way LRU, PPC with HDC can improve the throughput and the energy by 13.0% and 11.2%, respectively.

VI. CONCLUSION

In this paper, an efficient cache replacement policy named Hit Dominance Cache (HDC) was proposed to reduce the number of cache misses in Packet Processing Cache (PPC) without increasing the cache size. Conventionally, Least Recently Used (LRU) is used as the cache replacement policy of PPC because it is known that LRU performs good in many cache systems. However, from the difference of the cache access patterns, LRU is not suitable for PPC. LRU cannot evict mice flows, which account for most flows in networks and make few hits in the cache, rapidly. As a result, the cache entries are polluted by the mice flows.

HDC divides the cache into two areas and assigns flows to the appropriate area depending on the number of references in the cache. HDC can evict the mice flows rapidly by inserting a new entry into the least recently used position and retain the elephant flows preferentially by shifting the entry hit many times to another area. The simulation result with 15 real network traces showed that HDC (4-way hit area, 4-way primary area, and HDC threshold = 8) can reduce the cache misses by up to 29.1% (12.5% on average) compared to 4-way LRU. Furthermore, the hardware implementation using Verilog-HDL showed that the hardware costs of HDC are comparable to those of 4-way LRU. Finally, we showed that PPC with HDC can achieve approximately 500 Gbps with 88.8% energy of conventional PPC (20.5% energy of TCAM only architecture).

REFERENCES

- [1] The Ministry, "Tabulation and Estimation of Internet Traffic in Japan," 2016, Available: http://www.soumu.go.jp/main_content/000462459.pdf. [Accessed May. 6, 2018]
- [2] METI, "Green IT Initiative in Japan," Available: <http://www.meti.go.jp/english/policy/GreenITInitiativeInJapan.pdf>. [Accessed May. 6, 2018]
- [3] J. Fan, C. Hu, K. He, J. Jiang, and B. Liuy, "Reducing power of traffic manager in routers via dynamic on/off-chip scheduling," 2012 Proc. IEEE INFOCOM, Orlando, FL, 2012, pp. 1925-1933.
- [4] X. Zheng, X. Wang, "Comparative study of power consumption of a NetFPGA-based forwarding node in publish-subscribe Internet routing," Computer Communications, vol. 44, 2014, pp. 36-43.
- [5] S. Gamage and A. Pasqual, "High performance parallel packet Classification architecture with Popular Rule Caching," 2012 18th IEEE Int'l. Conf. on Networks (ICON), Singapore, 2012, pp. 52-57.
- [6] B. Talbot, T. Sherwood, and B. Lin, "IP caching for terabit speed routers," Global Telecommunications Conference (GLOBECOM '99), Brazil, vol.2, 1999, pp.1565-1569.
- [7] N. B. Guinde, R. Rojas-Cessa, and S. G. Zivarras, "Packet classification using rule caching," 2013 Fourth International Conference on Information, Intelligence, Systems and Applications (IISA 2013), Piraeus, 2013, pp.1-6.
- [8] B. Agrawal and T. Sherwood, "Ternary CAM Power and Delay Model: Extensions and Uses," in IEEE Trans. on Very Large Scale Integration (VLSI) Systems, vol. 16, no. 5, 2008, pp. 554-564.
- [9] M. Nawa et al., "Energy-efficient high-speed search engine using a multi-dimensional TCAM architecture with parallel pipelined subdivided structure," 2016 13th IEEE Annual Consumer

- Communications & Networking Conference (CCNC), Las Vegas, NV, 2016, pp. 309-314.
- [10] Hewlett-Packard Development Company, "Energy Efficient Networking - Business white paper," 2011, Available: <http://h17007.www1.hp.com/docs/mark/4AA3-3866ENW.pdf>. [Accessed May. 6, 2018]
- [11] H. Yamaki, "Flow Characteristic-Aware Cache Replacement Policy for Packet Processing Cache," In Proc. of Future of Information and Communication Conference (FICC 2018), Singapore, 2018, pp.1-8.
- [12] G. S. Shenoy, J. Tubella, A. Gonzalez, "Exploiting temporal locality in network traffic using commodity multi-cores," 2012 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), 2012, pp.110-111.
- [13] C. Girish and R. Govindarajan, "Improving performance of digest caches in network processors," In Proc. of the 15th Int'l. Conf. on High performance computing (HiPC'08), India, 2008, pp.6-17.
- [14] B. Agrawal and T. Sherwood, "Modeling TCAM power for next generation network devices," 2006 IEEE International Symposium on Performance Analysis of Systems and Software, 2006, pp. 120-129.
- [15] C. Kim, M. Caesar, A. Gerber, and J. Rexford, "Revisiting Route Caching: The World Should Be Flat," In Proc. of the 10th International Conference on Passive and Active Network Measurement (PAM '09), Berlin, 2009, pp.3-12.
- [16] K. Y. Ho and Y. C. Chen, "Performance evaluation of ipv6 packet classification with caching," 2008 Third Int'l Conf. on Communications and Networking in China, Hangzhou, 2008, pp. 669-673.
- [17] K. Li, F. Chang, D. Berger, F. Wu-chang, "Architectures for packet classification caching," The 11th IEEE International Conference on Networks (ICON2003), Sydney, 2003, pp. 111-117.
- [18] L. A. Belady, "A study of replacement algorithms for a virtual-storage computer." IBM Syst. J. vol. 5, no. 2, 1966, pp. 78-101.
- [19] F. Chang, W. C. Feng, and K. Li, "Efficient Packet Classification with Digest Caches," Proc. Third Workshop Network Processors and Applications (NP-3), 2005.
- [20] S. Ata, M. Murata, and H. Miyahara, "Efficient cache structures of IP routers to provide policy-based services," IEEE Int'l. Conf. on Communications (ICC 2001), Helsinki, vol.5, 2001, pp. 1561-1565.
- [21] N. Kim, S. Jean, J. Kim, and H. Yoon, "Cache replacement schemes for data-driven label switching networks," 2001 IEEE Workshop on High Performance Switching and Routing, Dallas, TX, 2001, pp. 223-227.
- [22] H. Yamaki and H. Nishi, "An Improved Cache Mechanism for a Cache-based Network Processor," In Proc. of the Int'l. Conf. on Parallel and Distributed Processing Techniques and Applications (PDPTA '12), Las Vegas, NV, 2012, pp. 1-7.
- [23] H. Yamaki and H. Nishi, "Line Replacement Algorithm for L1-scale Packet Processing Cache," In Adjunct Proc. of the 13th Int'l. Conf. on Mobile and Ubiquitous Systems: Computing Networking and Services (MOBIQUITOUS 2016), Hiroshima, Japan, 2016, pp. 12-17.
- [24] T. Mori, M. Uchida, R. Kawahara, J. Pan, and S. Goto, "Identifying elephant flows through periodically sampled packets," In Proceedings of the 4th ACM SIGCOMM conference on Internet measurement (IMC '04). ACM, New York, USA, 2004, pp.115-120.
- [25] M. Okuno and H. Nishi, "Network-Processor Acceleration-Architecture Using Header-Learning Cache and Cache-Miss Handler," The 8th World Multi-Conference on Systemics, Cybernetics and Informatics (SCI2004), 2004, pp. 108-113.
- [26] RIPE Network Coordination Centre, "Réseaux IP Européens Network Coordination Centre RIPE NCC," Available: <http://www.ripe.net/>. [Accessed May. 6, 2018]
- [27] WIDE MAWI WorkingGroup, "MAWI Working Group Traffic Archive" Available: <http://mawi.wide.ad.jp/mawi/>. [Accessed May. 6, 2018]
- [28] R. Karedla, J. S. Love, B. G. Wherry, "Caching strategies to improve disk system performance," in Computer, vol. 27, no. 3, 1994, pp. 38-46.
- [29] North Carolina State University, "FreePDK45:Contents," Available: <http://www.eda.ncsu.edu/wiki/FreePDK45:Contents>. [Accessed May. 6, 2018]
- [30] N. Muralimanohar et al., "Optimizing NUCA organizations and wiring alternatives for large caches with CACTI 6.0," In Proc. of the 40th Annual IEEE/ACM Int'l. Symposium on Microarchitecture (MICRO 40), Chicago, USA, 2007, pp.3-14.

Fuzzy Logic-Controlled 6-DOF Robotic Arm Color-based Sorter with Machine Vision Feedback

Alexander C. Abad¹, Dino Dominic Ligutan¹,
Elmer P. Dadios²

¹Electronics and Communications Engineering Department
²Manufacturing Engineering and Management Department
Gokongwei College of Engineering
De La Salle University
Manila, Philippines

Levin Jaeron S. Cruz, Michael Carlo D.P. Del Rosario,
Jho Nathan Singh Kudhal

Electronics and Communications Engineering Department
Gokongwei College of Engineering
De La Salle University – Laguna Campus
Biñan Laguna, Philippines

Abstract—A demonstration of the application of fuzzy logic-based joint controller (FLJC) to a 6-DOF robotic arm as a color-based sorter system is presented in this study. The robotic arm with FLJC is integrated with a machine vision system that can discriminate different colors. Additionally, the machine vision system composed of Kinect camera and computer were used to extract the coordinates of the gripper and the objects within the image of the workspace. A graphical user interface with an underlying sorting algorithm allows the user to control the sorting process. Once the system is configured, the computed joint angles by FLJC are transmitted serially to the microcontroller. The results show that the absolute error of the gripper coordinates is less than 2 cm and that the machine vision is capable of achieving at least 95% accuracy in proper color discrimination both for first and second level stacked color objects.

Keywords—Color-based sorter; degrees of freedom; fuzzy logic; joint controller; machine vision; robotic arm

I. INTRODUCTION

The development of machines has been a valuable tool ever since the dawn of civilization. Machines had been the humanity's innovative creations whose sole purpose was to achieve efficiency and effectiveness to different tasks that are either routine or almost impossible for humans to do by hand. Machines were meant to be driven by a human operator, until the last century [1] where automation began to be favored by industry, specially deployed in car manufacturing process. This greatly reduced the manpower needed and at the same time was efficient in terms of resources and time. From thereon, autonomous machines came into existence and diverse forms of such machines were developed for specific purposes. One such machine is the autonomous robotic arm whose design was primarily inspired by the human arm. Due to the flexibility that the human arm can do varied tasks, the development of an autonomous robotic arm has been a subject of research [2] since its development in 1960s.

Autonomous robotic arms had numerous advantages as compared to human arm. Robotic arm machines are immune to fatigue and can be made to be invulnerable in wide environment settings. Additionally, it is the most viable alternative when deployed to environments that are too harmful for humans [3] and can be programmed to perform routine

tasks efficiently. Amidst these benefits, a robotic arm is also a complex mechanical machine that exhibits time-varying inertia and friction and as such is more challenging to control by means of classical linear-based controllers. To achieve autonomous operation, the machine must have a controller that is able to sense its current state and decide its course action in much the same way humans decide. Non-classical or intelligent controllers had been developed throughout the years, such as fuzzy logic based controllers [5], [6] that mimics the way humans think, artificial neural network based controllers [7], [8] that emulates the biological human brain, genetic algorithm based controllers [9], [10] inspired by evolutionary processes or hybrid types [11]. One such controller developed in this study is the fuzzy logic-based joint controller (FLJC) [4] that is capable of dealing with system nonlinearities by moving the joints of the robotic arm at proper rate and interval according to the task at hand. Fuzzy logic controllers has been shown as an effective controller in a number of robot systems like the micro soccer robots [12]-[15], micro-golf robot [16], ball-beam balancing robot [17] and simulated and actual robotic arms [4], [6], [18]-[21].

Aside from the controller developed in [4], this study will give emphasis on the integration of the controller with a machine vision system to demonstrate the use of the fuzzy logic controlled autonomous robotic arm system into a color-based sorter system. The machine vision system will be thoroughly discussed as well as the algorithm deployed to perform the sorting process to realize a fully functional color-based sorter. Test results of accuracy of the gripper to move towards the target coordinates as well as the reliability of the machine vision system are laid out and discussed. Lastly, several points are enumerated with regards to the possible improvements that could be made for the system.

II. SYSTEM CONFIGURATION

The color-based sorter system is similar to the configuration in [4] but with the following modifications: 1) the robotic arm's end-effector are embedded with limit switches to improve tactile sensing of the object, 2) the machine vision system is now capable of discriminating at most four different object colors, and 3) the sorter is capable of sorting out stacked objects up to second level. The robotic arm

itself is composed of a 4-DOF M100RAK [22] modular arm attached with 2-DOF gripper [23]. The shoulder, elbow and wrist joints are each mounted with MPU6050 Six-Axis Gyroscope and Accelerometer [24] Inertial Measurement Units (IMU) as sensors to acquire the robotic arm's pose in real time. Attached to the robotic arm's end-effector are the force sensing resistor (FSR) and miniature limit switches as its haptic feedback sensors. The robotic arm's servo motors are controlled directly by the Arduino [25] microcontroller that communicates with the computer. Set atop on the workspace is the Kinect sensor [26] that serves as the main sensory input for machine vision system in the computer. Fig. 1. shows how the components are connected to form the fuzzy logic-controlled color-based sorter.

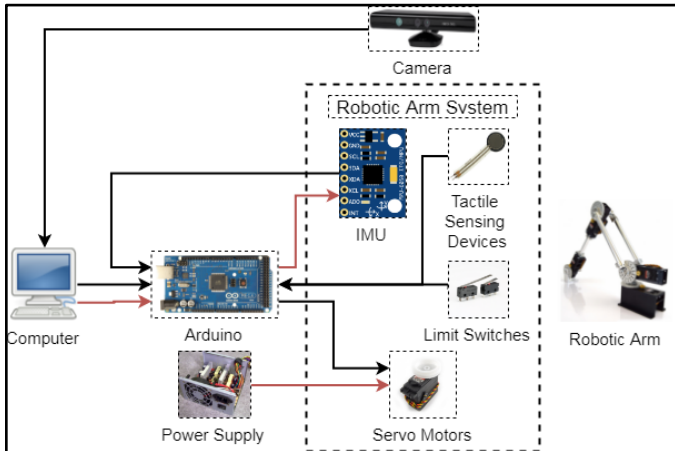


Fig. 1. Architecture of the color-based sorter.

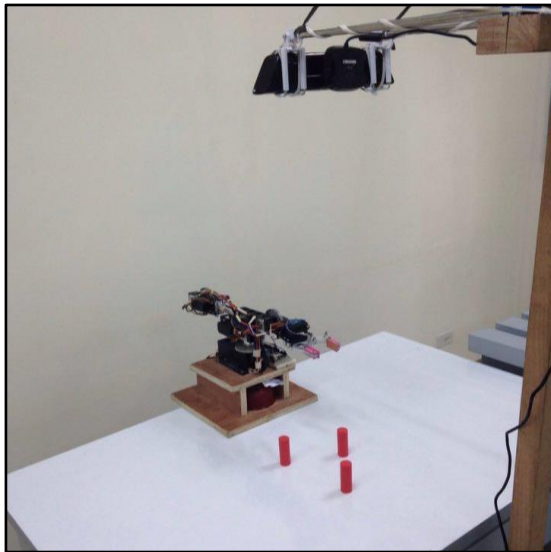


Fig. 2. Configuration of the color-based object sorter.

Shown in Fig. 2 is the hardware configuration used in this study with the robotic arm on the center in front of the cylinder objects and the camera on top. The pertinent dimensions of the workspace are shown in Fig. 3. The study focuses in the application of fuzzy logic-based controller of the robotic arm as well as the algorithm devised to properly sort the cylinder objects in place.

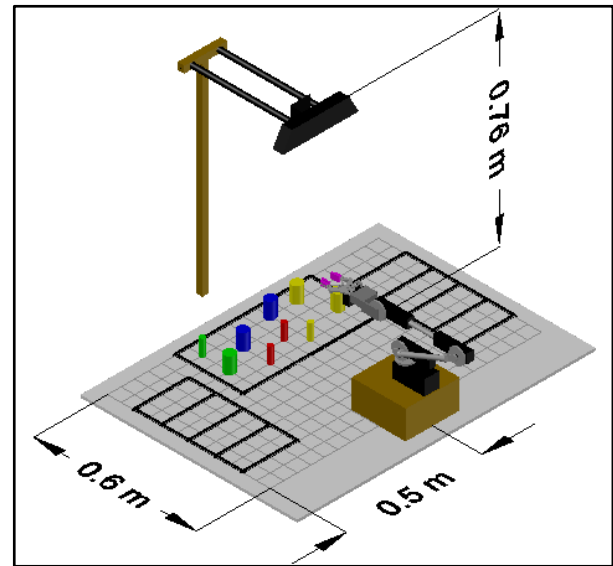


Fig. 3. Workspace dimensions of the color-based sorter.

III. FUZZY LOGIC-BASED JOINT CONTROLLER

The theory of fuzzy sets was first described by Lotfi Zadeh [27] and found its applications as a controller such as for plant processes [28]. Fuzzy sets are an extension of the bi-valued logic in that it can be used to describe half-truth statements to varying degrees. The concept of a fuzzy set can be exploited to emulate the way humans think when in control of a process by employing a human-like language describing how a complex system should be controlled. To achieve a descriptive language for control, a fuzzy logic controller consists of: 1) a fuzzifier block that converts real-world crisp values into fuzzy sets through membership functions, 2) an inference engine that interprets the input fuzzy set based on a set of human-defined language for control known as fuzzy rules to decide the output fuzzy sets, and 3) a defuzzifier block that converts the output fuzzy set back into real-world crisp values [29]. These crisp values are now used to directly control any process variables [5], [11], [17]. Shown in Fig. 4. is the conceptual block diagram of a fuzzy logic controller. The goal of the fuzzy logic controller is to move the end-effector to the desired target as close as possible. The controller is part of a closed-loop system composed of the sensors mounted on the robotic arm, the controller itself and the mechanically actuated robotic arm. The fuzzy logic controller dictates the microcontroller the amount and direction at which the servo motors are to be turned and the microcontroller in turn, through pulse width modulation signals controls the servo motors.

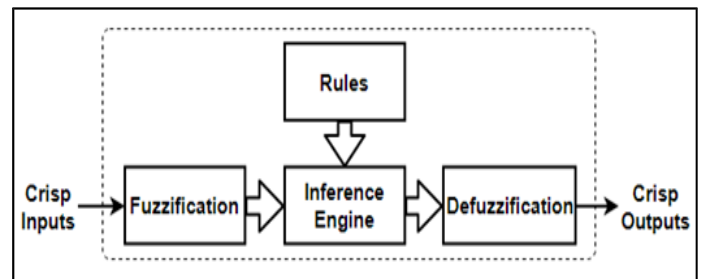


Fig. 4. Fuzzy logic system [5].

A. Input and Output Parameters

The top and side views of the robotic arm with pertinent dimensions are shown in Fig. 5. and 6. Excluding the 2-DOF gripper, there are four (4) joint angles that can be controlled to change the end-effector's position: the base angle (θ_b), shoulder angle (θ_s), elbow angle (θ_e) and wrist angle (θ_w). The fuzzy logic controller must control these joints so that the input errors in x-coordinates (e_x), y-coordinates (e_y) and z-coordinates (e_z) are close to zero as possible.

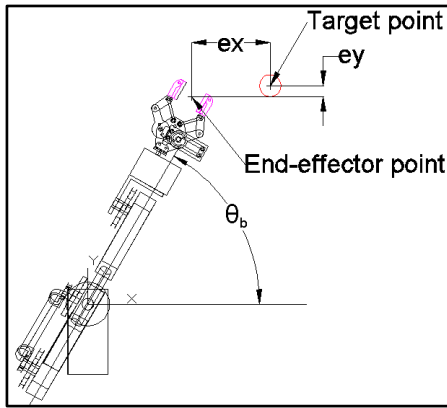


Fig. 5. Top view of the robotic arm relative to target.

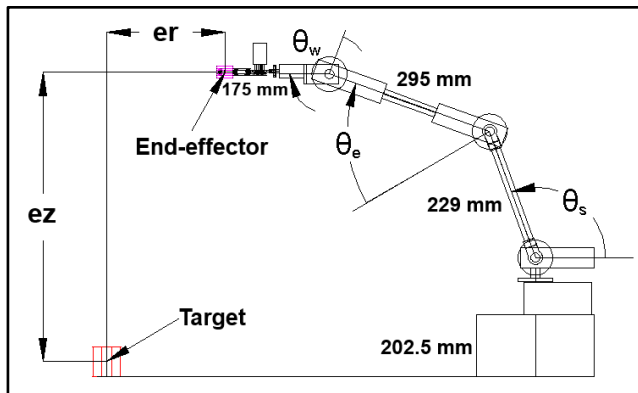


Fig. 6. Side view of the robotic arm relative to target.

The inputs to the controller were chosen according the following criteria: 1) the controller must know how close the end-effector is to the target, and 2) the controller must determine the current pose of the robotic arm to properly move the actuators in the desired direction. With these in mind, listed below are the input parameters for the fuzzy logic controller:

- 1) e_x is error in x-coordinate, defined to be difference between the x-coordinate of the end-effector and the x-coordinate of the target.
- 2) e_y is error in y-coordinate, defined to be difference between the y-coordinate of the end-effector and the y-coordinate of the target.
- 3) e_z is error in z-coordinate, defined to be difference between the z-coordinate of the end-effector and the z-coordinate of the target.
- 4) θ_b is the base angle, defined as the angle between the robotic arm and the x-axis.

- 5) θ_e is the elbow joint angle.
- 6) θ_η is defined as the gripper angle with respect to horizontal.
- 7) θ'_η is defined as the rate of change of gripper angle with respect to horizontal.

The outputs of the fuzzy logic controller are as follows:

- 1) $\Delta\theta_b$ is the change in base joint angle.
- 2) $\Delta\theta_s$ is the change in shoulder joint angle.
- 3) $\Delta\theta_e$ is the change in elbow joint angle.
- 4) $\Delta\theta_w$ is the change in wrist joint angle.

B. Membership Functions

Once the input and output parameters are defined, the appropriate membership functions for each parameter are defined according to the limitations of the robotic arm itself as well as the magnitude of the change produced by each parameter. These membership functions are then tuned and finalized through a series of tests and experimentations [5]. The input membership functions are tuned by considering the sensitivity of the controller to these inputs. In this study, the unit of measurement for the range of values sampled to discrete grades of membership functions for input errors is in millimeters while those for angular displacement are in radians. Trapezoidal membership functions were used at the extreme values of input joint angles to avoid self-collision. Shown in Fig. 7. through Fig. 13. are the membership functions of the seven input parameters. For the sake of brevity, the membership functions are labeled accordingly as follows:

Fuzzy Membership Acronyms:

- | | | |
|------------|-----------------------|---------------------|
| L – left | NL – negative large | P – positive |
| M – middle | N – negative | PL – positive large |
| R – right | Z – zero (negligible) | |

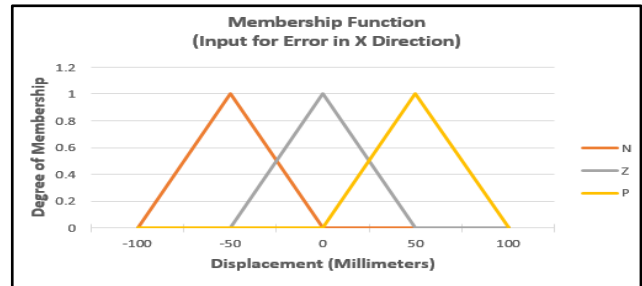


Fig. 7. Membership function for error in x-coordinate.

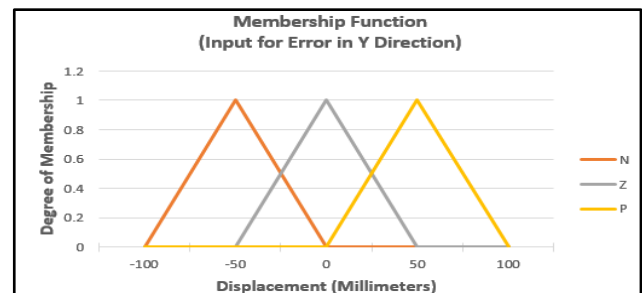


Fig. 8. Membership function for error in y-coordinate.

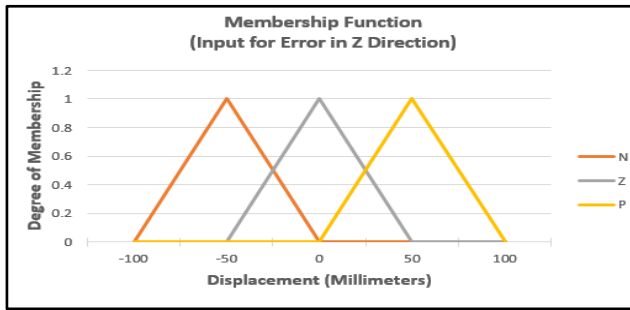


Fig. 9. Membership function for error in z-coordinate.

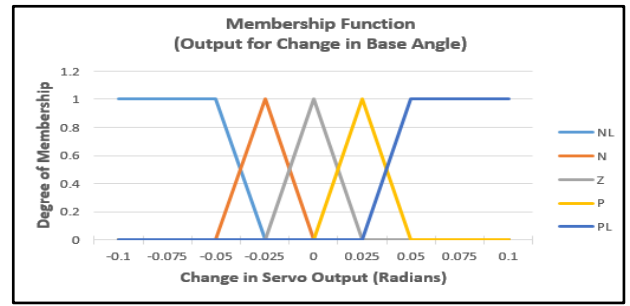


Fig. 14. Membership function for change in base joint angle.

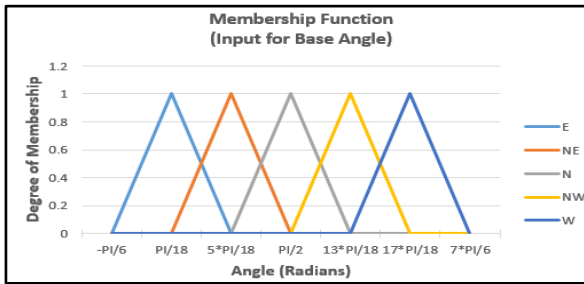


Fig. 10. Membership function for base joint angle.

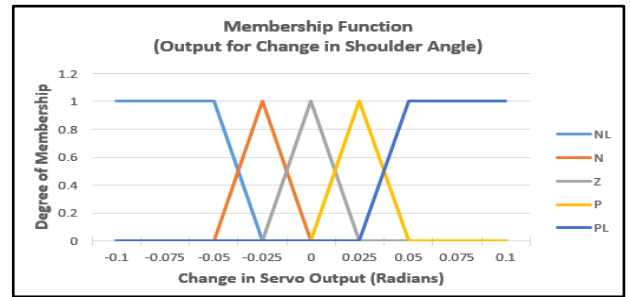


Fig. 15. Membership function for change in shoulder joint angle.

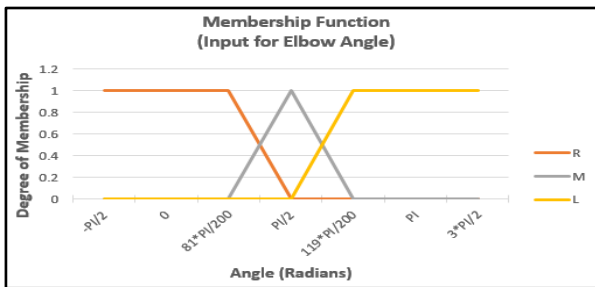


Fig. 11. Membership function for elbow joint angle.

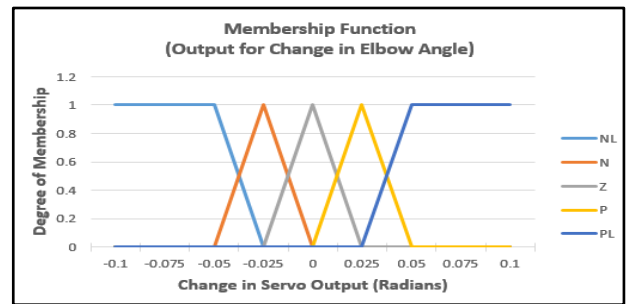


Fig. 16. Membership function for change in elbow joint angle.

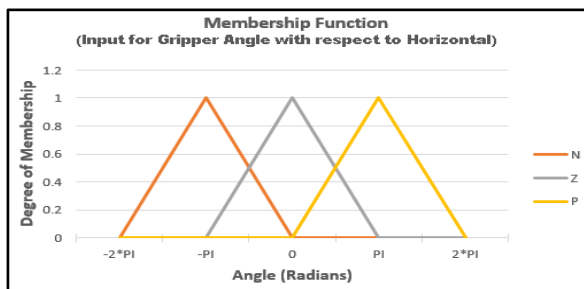


Fig. 12. Membership function for gripper angle with respect to horizontal.

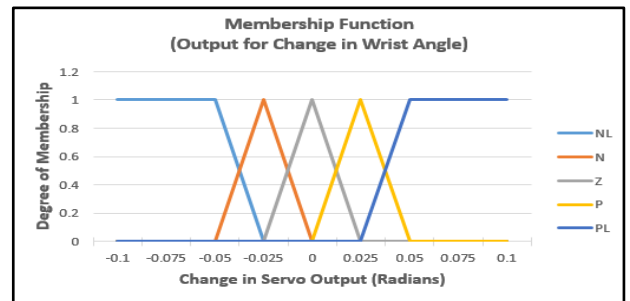


Fig. 17. Membership function for change in wrist joint angle.

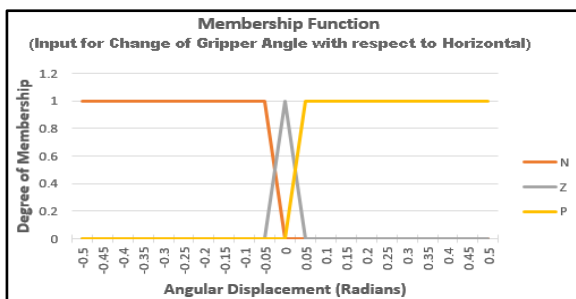


Fig. 13. Membership function for rate of change of gripper angle with respect to horizontal.

Similarly, the output membership functions are tuned by considering the sensitivity of the robotic arm as the joint angles were changed. All output joint angles are specified in units of radians. The defuzzification process used the weighted average method to reduce calculation time in calculating the crisp value. Shown in Fig. 14 through Fig. 17 are the membership functions of the four output joint angles: change in base, shoulder, elbow and wrist joint angles. The same membership labeling scheme applies as defined for the input membership functions.

C. Analysis of the Different Robotic Arm Poses

Once the input and output parameters were determined as well as their respective membership functions, the rules for inference engine are formulated. The rules can be formulated by analyzing the different robotic arm poses possible within the workspace. Of course, there are infinite arm poses that are possible within the workspace so dividing the range of possible values into subsets is necessary and it can be done by the aid of membership functions. The pose of the robotic arm is analyzed by looking at the top and side view of the robotic arm shown in Fig. 18. and 19. The different poses shall be the basis in formulating the fuzzy rules. In general, the rules are to be formulated in such a way that the input errors in x-, y- and z-coordinates are minimized in each iteration.

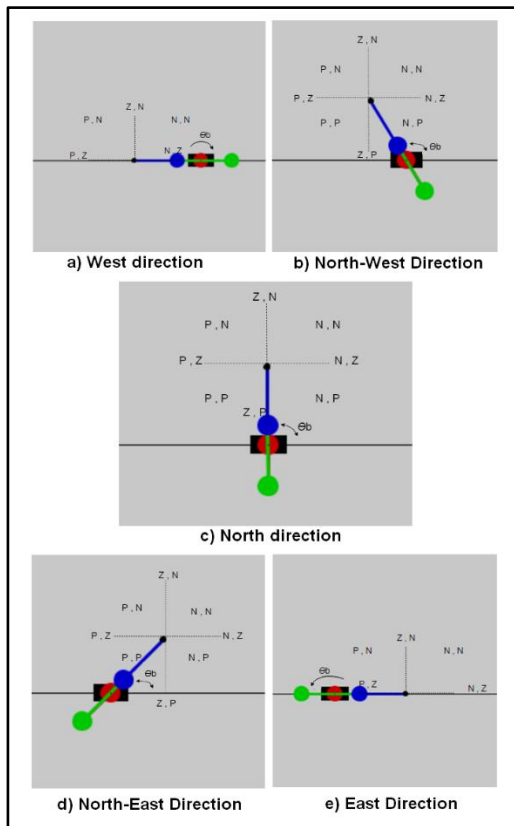


Fig. 18. Top view of possible robotic arm orientation.

The base joint angle can be oriented in five (5) different angle orientations as shown in Fig. 18: West, Northwest, North, Northeast and East. In the same figure, the black rectangle represents the base of the arm, the red link represents the shoulder-to-elbow link, the green link represents the elbow-to-wrist link and the blue link represents the gripper. The initials P, Z, and N corresponds to positive, zero and negative respectively each used to describe the position of the end-effector relative to the target. The symbol Θ_b is the base angle and a pair such as (P,N) denotes that the input errors for x and y coordinates are positive and negative respectively should the target is found at that region relative to the end-effector. Knowing the sign of the input errors will aid on formulating the fuzzy rule at which should the base angle be moved to

minimize the error. In this view, the arm can rotate clockwise or counterclockwise as well as extend or retract its links.

Shown in Fig. 19 are three possible poses when looking at the side view of the robotic arm. The three links form a coupled system that has three (3) degrees of freedom and is more than the degrees of freedom necessary to determine the radius and height of the end-effector. As such, a link can be assumed to be at fixed angle and isolate it from the other two angles. The gripper is chosen to be this link that can be fixed to maintain horizontally level with respect to the ground at all times. Effectively, we could decouple the gripper and write separate fuzzy rules for it apart from the shoulder and elbow joint angles.

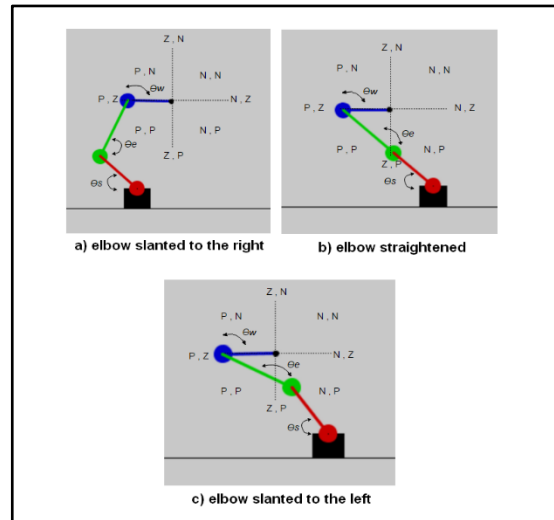


Fig. 19. Side view of possible robotic arm orientation.

D. Fuzzy Rule Formulation

By analyzing the different poses of the robotic arm, the fuzzy rules can now be facilitated by taking note of the input errors as well as their signs. In general, the rule formulation is guided by the control law that all input errors must be minimized and as close to zero as possible. From the analysis of the robotic arm, three (3) different rule blocks can be identified. For instance, if the base angle is pointing in the North direction and the target is present at the (P,N) region then the base angle must rotate counterclockwise and the robot arm must extend forward, to bring the end-effector closer to the target. The beauty of fuzzy logic controller is that you do not have to specify the magnitude explicitly but just the intuition and at which direction should the output parameters move. This analysis is applied to all enumerated poses and the rules formulated can be found on Table I through Table III(a)

Table I pertains to the fuzzy rules for the top view orientation involving the input parameters base joint angle, error in x-coordinate, error in y-coordinate and output parameter change in base joint angle. Table II contains the rules for the side view orientation involving the input parameters elbow joint angle, error in y-coordinate, error in z-coordinate and output parameters change in shoulder and elbow angles. Lastly, Table III is a list that controls how the gripper angle must maintain horizontally level at all times.

TABLE I. FUZZY RULES FOR BASE JOINT ANGLE

	Input: base joint angle (θ_b), error x (e_x), error y (e_y) Output: change in base joint angle ($\Delta\theta_b$)
1	If θ_b is E and e_x is N and e_y is N then $\Delta\theta_b$ is P
2	If θ_b is E and e_x is N and e_y is Z then $\Delta\theta_b$ is Z
3	If θ_b is E and e_x is N and e_y is P then $\Delta\theta_b$ is N
4	If θ_b is E and e_x is Z and e_y is N then $\Delta\theta_b$ is P
5	If θ_b is E and e_x is Z and e_y is Z then $\Delta\theta_b$ is Z
6	If θ_b is E and e_x is Z and e_y is P then $\Delta\theta_b$ is N
7	If θ_b is E and e_x is P and e_y is N then $\Delta\theta_b$ is P
8	If θ_b is E and e_x is P and e_y is Z then $\Delta\theta_b$ is Z
9	If θ_b is E and e_x is P and e_y is P then $\Delta\theta_b$ is N
10	If θ_b is NE and e_x is N and e_y is N then $\Delta\theta_b$ is Z
11	If θ_b is NE and e_x is N and e_y is Z then $\Delta\theta_b$ is N
12	If θ_b is NE and e_x is N and e_y is P then $\Delta\theta_b$ is N
13	If θ_b is NE and e_x is Z and e_y is N then $\Delta\theta_b$ is P
14	If θ_b is NE and e_x is Z and e_y is Z then $\Delta\theta_b$ is Z
15	If θ_b is NE and e_x is Z and e_y is P then $\Delta\theta_b$ is P
16	If θ_b is NE and e_x is P and e_y is N then $\Delta\theta_b$ is P
17	If θ_b is NE and e_x is P and e_y is Z then $\Delta\theta_b$ is P
18	If θ_b is NE and e_x is P and e_y is P then $\Delta\theta_b$ is Z
19	If θ_b is N and e_x is N and e_y is N then $\Delta\theta_b$ is N
20	If θ_b is N and e_x is N and e_y is Z then $\Delta\theta_b$ is N
21	If θ_b is N and e_x is N and e_y is P then $\Delta\theta_b$ is N
22	If θ_b is N and e_x is Z and e_y is N then $\Delta\theta_b$ is Z
23	If θ_b is N and e_x is Z and e_y is Z then $\Delta\theta_b$ is Z
24	If θ_b is N and e_x is Z and e_y is P then $\Delta\theta_b$ is Z
25	If θ_b is N and e_x is P and e_y is N then $\Delta\theta_b$ is P
26	If θ_b is N and e_x is P and e_y is Z then $\Delta\theta_b$ is P
27	If θ_b is N and e_x is P and e_y is P then $\Delta\theta_b$ is P
28	If θ_b is NW and e_x is N and e_y is N then $\Delta\theta_b$ is N
29	If θ_b is NW and e_x is N and e_y is Z then $\Delta\theta_b$ is N
30	If θ_b is NW and e_x is N and e_y is P then $\Delta\theta_b$ is Z
31	If θ_b is NW and e_x is Z and e_y is N then $\Delta\theta_b$ is N
32	If θ_b is NW and e_x is Z and e_y is Z then $\Delta\theta_b$ is Z
33	If θ_b is NW and e_x is Z and e_y is P then $\Delta\theta_b$ is P
34	If θ_b is NW and e_x is P and e_y is N then $\Delta\theta_b$ is Z
35	If θ_b is NW and e_x is P and e_y is Z then $\Delta\theta_b$ is P
36	If θ_b is NW and e_x is P and e_y is P then $\Delta\theta_b$ is P
37	If θ_b is W and e_x is N and e_y is N then $\Delta\theta_b$ is N
38	If θ_b is W and e_x is N and e_y is Z then $\Delta\theta_b$ is Z
39	If θ_b is W and e_x is N and e_y is P then $\Delta\theta_b$ is P
40	If θ_b is W and e_x is Z and e_y is N then $\Delta\theta_b$ is N
41	If θ_b is W and e_x is Z and e_y is Z then $\Delta\theta_b$ is Z
42	If θ_b is W and e_x is Z and e_y is P then $\Delta\theta_b$ is P
43	If θ_b is W and e_x is P and e_y is N then $\Delta\theta_b$ is N
44	If θ_b is W and e_x is P and e_y is Z then $\Delta\theta_b$ is Z
45	If θ_b is W and e_x is P and e_y is P then $\Delta\theta_b$ is P

TABLE II. FUZZY RULES FOR SHOULDER AND ELBOW JOINT ANGLES

	Input: elbow joint angle (θ_e), error y (e_y), error z (e_z) Output: change in shoulder joint angle ($\Delta\theta_s$), change in elbow joint angle ($\Delta\theta_e$)
1	If θ_e is R and e_y is N and e_z is N then $\Delta\theta_s$ is P and $\Delta\theta_e$ is P
2	If θ_e is R and e_y is N and e_z is Z then $\Delta\theta_s$ is P and $\Delta\theta_e$ is Z
3	If θ_e is R and e_y is N and e_z is P then $\Delta\theta_s$ is P and $\Delta\theta_e$ is N
4	If θ_e is R and e_y is Z and e_z is N then $\Delta\theta_s$ is Z and $\Delta\theta_e$ is P
5	If θ_e is R and e_y is Z and e_z is Z then $\Delta\theta_s$ is Z and $\Delta\theta_e$ is Z
6	If θ_e is R and e_y is Z and e_z is P then $\Delta\theta_s$ is Z and $\Delta\theta_e$ is N

7	If θ_e is R and e_y is P and e_z is N then $\Delta\theta_s$ is N and $\Delta\theta_e$ is P
8	If θ_e is R and e_y is P and e_z is Z then $\Delta\theta_s$ is N and $\Delta\theta_e$ is Z
9	If θ_e is R and e_y is P and e_z is P then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
10	If θ_e is M and e_y is N and e_z is N then $\Delta\theta_s$ is Z and $\Delta\theta_e$ is Z
11	If θ_e is M and e_y is N and e_z is Z then $\Delta\theta_s$ is N and $\Delta\theta_e$ is Z
12	If θ_e is M and e_y is N and e_z is P then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
13	If θ_e is M and e_y is Z and e_z is N then $\Delta\theta_s$ is P and $\Delta\theta_e$ is Z
14	If θ_e is M and e_y is Z and e_z is Z then $\Delta\theta_s$ is Z and $\Delta\theta_e$ is Z
15	If θ_e is M and e_y is Z and e_z is P then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
16	If θ_e is M and e_y is P and e_z is N then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
17	If θ_e is M and e_y is P and e_z is Z then $\Delta\theta_s$ is N and $\Delta\theta_e$ is Z
18	If θ_e is M and e_y is P and e_z is P then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
19	If θ_e is L and e_y is N and e_z is N then $\Delta\theta_s$ is P and $\Delta\theta_e$ is N
20	If θ_e is L and e_y is N and e_z is Z then $\Delta\theta_s$ is P and $\Delta\theta_e$ is N
21	If θ_e is L and e_y is N and e_z is P then $\Delta\theta_s$ is P and $\Delta\theta_e$ is N
22	If θ_e is L and e_y is Z and e_z is N then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
23	If θ_e is L and e_y is Z and e_z is Z then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
24	If θ_e is L and e_y is Z and e_z is P then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
25	If θ_e is L and e_y is P and e_z is N then $\Delta\theta_s$ is Z and $\Delta\theta_e$ is N
26	If θ_e is L and e_y is P and e_z is Z then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N
27	If θ_e is L and e_y is P and e_z is P then $\Delta\theta_s$ is N and $\Delta\theta_e$ is N

TABLE III. FUZZY RULES FOR WRIST ANGLE

	Input: gripper angle w.r.t. horizontal (θ_η), rate of change of eta ($\Delta\theta_\eta$) Output: change in wrist joint angle ($\Delta\theta_w$)
1	If θ_η is N and $\Delta\theta_\eta$ is N then $\Delta\theta_w$ is N
2	If θ_η is N and $\Delta\theta_\eta$ is Z then $\Delta\theta_w$ is N
3	If θ_η is N and $\Delta\theta_\eta$ is P then $\Delta\theta_w$ is Z
4	If θ_η is Z and $\Delta\theta_\eta$ is N then $\Delta\theta_w$ is P
5	If θ_η is Z and $\Delta\theta_\eta$ is Z then $\Delta\theta_w$ is Z
6	If θ_η is Z and $\Delta\theta_\eta$ is P then $\Delta\theta_w$ is N
7	If θ_η is P and $\Delta\theta_\eta$ is N then $\Delta\theta_w$ is Z
8	If θ_η is P and $\Delta\theta_\eta$ is Z then $\Delta\theta_w$ is P
9	If θ_η is P and $\Delta\theta_\eta$ is P then $\Delta\theta_w$ is P

IV. MACHINE VISION SYSTEM

The machine vision system is composed of the camera as its sensory vision input and the computer as an image processing unit. The camera is the Kinect sensor [26] capable of providing not only colored images as well as image depth data. The image depth data was used to properly determine the height of the detected objects and consequently the stacking level of the cylinder objects. Furthermore, the depth data was used to filter out the white platform background by exploiting the fact that its distance is farther away from the camera itself. This method is referred to as depth masking.

The computer uses the Java-based Processing [30] software environment that provides the interfacing between the devices attached to it such as the Kinect sensor and the Arduino. Processing-based OpenCV [31] and SimpleOpenNI [32] libraries were used for the software development. The OpenCV library provided the tools to filter the image based on Hue-Saturation-Value (HSV) as well as detect the presence of blobs and their respective coordinates. On the other hand, the SimpleOpenNI library allows the system to communicate with the Kinect sensor to get the RGB image and depth data. Fig. 20. shows how the raw RGB image is eventually filtered out to keep the blue cylinder objects. The binarized image on the

right was processed further to detect the blobs present and store their coordinates for sorting purposes.

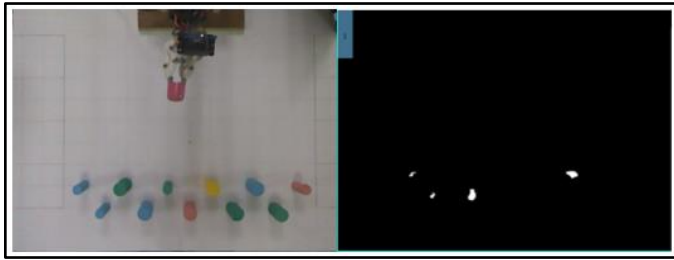


Fig. 20. On the left: actual gripper and cylinder objects as seen by the camera atop; on the right: filtered image showing the blue cylinder objects.

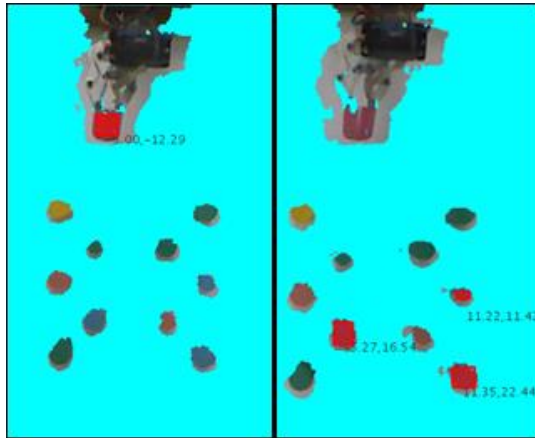


Fig. 21. Detection of gripper (left) and blue cylinder objects (right).

This filtering, detection and coordinate acquisition processes were done for other colors as well as for the gripper. The gripper is colored differently from the possible colors of the cylinder objects to properly recognize and locate the coordinates of the gripper itself. Shown in Fig. 21. is the results of detection and extraction of the coordinates of the gripper and the blue cylinder objects. It is worth mentioning that the gripper has two distinguishable shapes if it is wide open. In such a case, the reported coordinates of the gripper are found by calculating the centroid of the two separately detected gripper objects.

V. GRAPHICAL USER INTERFACE AND SORTING ALGORITHM

The coordinates of the gripper and the objects are measured relative to the origin point at the center of the region of interest as well as their respective heights relative to the white background platform. Together with the measured joint angles of the robotic arm, these values constitute the input variables for the sorting algorithm. The process begins by the user configuring the system through a designed graphical user interface (GUI) shown in Fig. 22. The user selects which object colors are to be sorted first according to the priority the user wishes. The user can also choose how the objects are to be sorted and decided where each color should land on pre-determined locations. After configuration, the “Sort” button can be pressed to begin the sorting process.

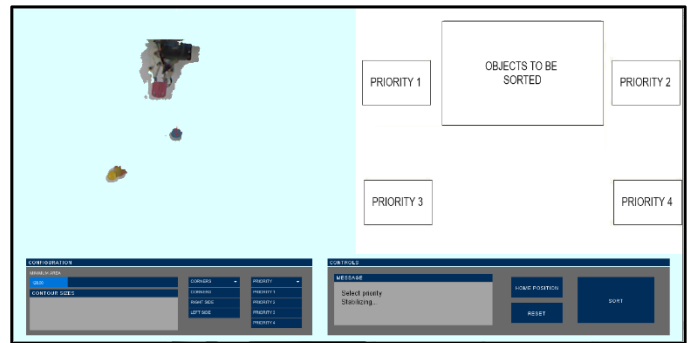


Fig. 22. The graphical user interface for the color-based sorter.

Pressing the “Sort” button invokes the sorting algorithm. The sorting algorithm first applies depth masking to differentiate depth levels among the platform, the objects and the gripper. The next process applies HSV-based color filtering to detect the presence of the objects in a specific color. The sequence at which colors are detected is determined by the priority configuration set by the user. As objects are detected for each color, the coordinates of the objects are acquired and stored to their respective buffers. After all object colors are found, the coordinates of the gripper are acquired. A hysteresis function is applied to the obtained coordinates to eliminate the sudden changes in coordinates due to noise.

Now that the coordinates of all objects are found, the gripper coordinates and the first object to be sorted or targeted is fed to the FLJC. There are several rules that determine which object should be fed to the FLJC: 1) based from the priority set by the user, and 2) the distance of the target object from the gripper. Since it is possible that there are multiple objects of the same color, the sorting algorithm would pick the target object with the minimal distance from the gripper. Should objects have same color and same distance from the gripper, the object with least change in base angle needed to reach will be picked up first. In this manner, the priority of which object should be picked up is resolved. Now that the target object is determined, the FLJC applies the appropriate changes to the joint angle needed to further minimize the distance between the gripper and the target. The machine vision will then locate the gripper coordinates and feed it back to the FLJC. This process repeats until the gripper coordinates is sufficiently coincident with that of the target object. It is also worth mentioning that a parallax error will be imposed upon the gripper coordinates relative to the platform depending on its height and location. To mitigate the parallax shift, a proper coordinate transformation is applied to the gripper coordinates before fed to the FLJC. For the coordinates (x, y) and height z as seen by the camera, the actual coordinates (x', y') are found to be:

$$x' = (x - x_c)(1 - p(z + B)/h) + x_c \quad (1)$$

$$y' = (y - y_c)(1 - p(z + B)/h) + y_c \quad (2)$$

where (x_c, y_c) are the coordinates of the center of ROI relative to the base of the robotic arm, p is the parallax factor, h is the camera height and B is the base height of the robotic arm.

As the gripper closes in to the object, the gripper is closed to grip the target object. The attached tactile sensors will determine if the object was indeed grasped. If the object is found to be grasped, a predetermined sequence of robotic arm

movements will place the object on its specific location according the configuration set by the user. Once the object is placed, the robotic arm goes into its home position. The sorting algorithm will feed the gripper coordinates and the next target object. This process will repeat until all such objects are placed to their appropriate destinations. During the sorting process, the user cannot reconfigure the sorter until it is finished but may press the emergency button to terminate.

VI. DATA AND RESULTS

A. Robotic Arm Simulator

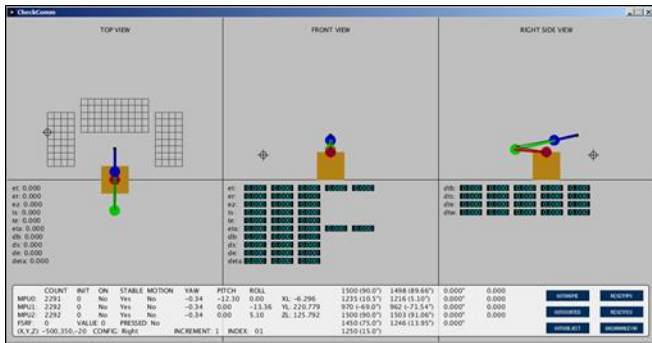


Fig. 23. Robotic arm simulator.

The robotic arm simulator shown in Fig. 23. is an improvement over that shown in [4]. The simulator can function as a monitoring tool to show the actual robotic arm pose in real time as well as provide a visualization of the controller in action. Furthermore, important parameters are indicated below to guide in the tuning of the fuzzy membership functions and several buttons that are programmed to move the robotic arm in a pre-determined sequence such as going to its home position and placing an object to designated areas. Once the actual robotic arm's response is satisfactory, the tuned fuzzy membership functions are transferred to a final program to be integrated with the sorting algorithm and a designed GUI.

B. Robotic Arm Movement

To test the accuracy of the fuzzy logic controller, the robotic arm was stretched forward along the Y direction. The plots of actual robotic arm end-effector coordinates plotted against the desired y-coordinate are shown in Fig. 24. through Fig. 26. A comparison was made against the inverse kinematic implementation. From Fig. 24. the fuzzy logic controller had lesser sideway excursions as compared to inverse kinematics implementation. The fuzzy logic controller was able to follow closely the ideal y-coordinate value as compared to inverse kinematic implementation shown in Fig. 25. The inverse kinematic implementation is found below the required y-value because of the weights of the robotic arm links. This effect is more pronounced as the height has significantly drooped shown in Fig. 26. Again, the moment due to the weight of extending arm is increasing as the y-coordinate increases. The fuzzy logic controller on the other hand managed to maintain a satisfactory level that is within 5 mm from the ideal height of 100 mm.

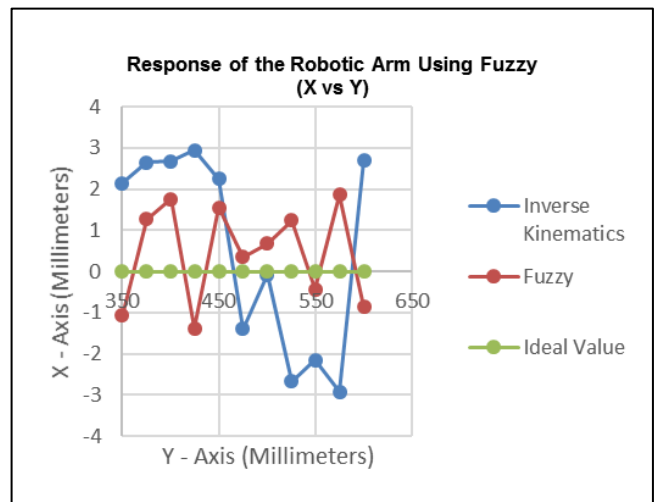


Fig. 24. X-Y movement response.

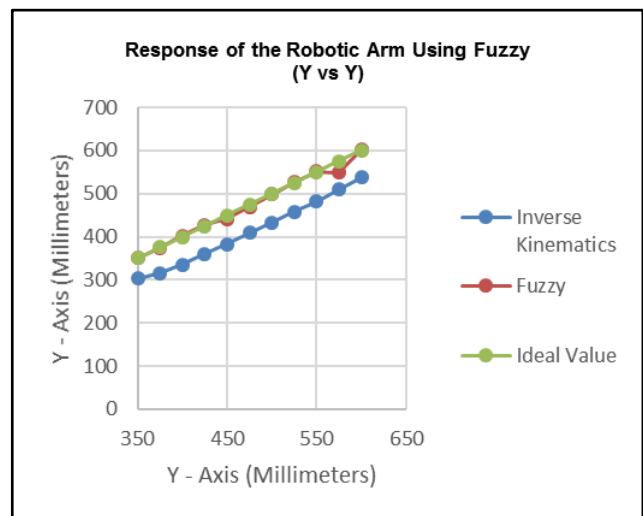


Fig. 25. Y-Y movement response.

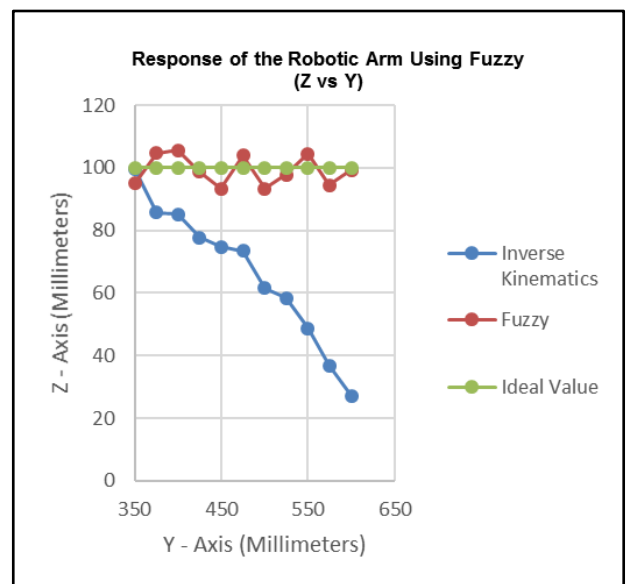


Fig. 26. Z-Y movement response.

C. Accuracy of the End Effector

The end-effector’s accuracy is tested by feeding the FLJC with an ideal coordinate coincident to the intersection of gridlines on the platform. By marking of the platform beneath the gripper, the distance between x-, y- and z-coordinates are obtained for at least 30 trials.

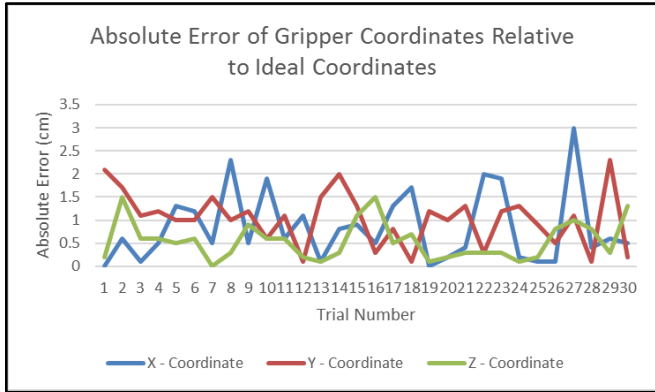


Fig. 27. Measure of end-effector’s absolute error before pickup.

These coordinates are randomly picked from the workspace area where the objects to be sorted are placed. The differences are measured and plotted as shown in Fig. 27. The average values of 0.8, 1 and 0.6 cm for absolute errors in x-, y- and z-coordinates were calculated for the end-effector relative to origin, respectively. On average, the end-effector coordinates were accurate enough to allow tolerance of 2 cm radius.

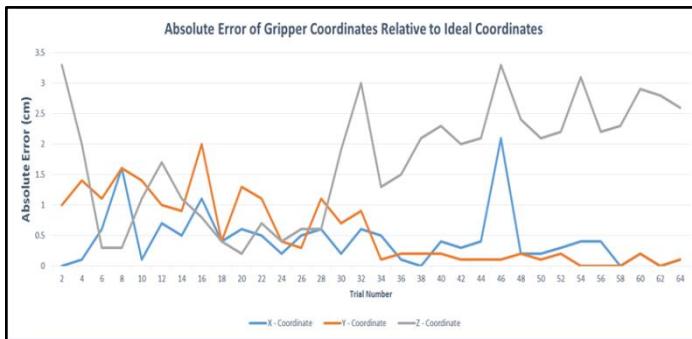


Fig. 28. Measure of end-effector’s absolute error after placement.

Fig. 28 shows the plot of absolute error of gripper coordinates in x-, y- and z-coordinates as it moved towards the pre-determined coordinates as destination for sorted objects. The average values of 0.45, 0.92 and 1.77 cm for absolute errors in x-, y- and z-coordinates were calculated for the end-effector relative to origin respectively. On average, the end-effector coordinates were accurate enough to allow tolerance of 2 cm radius.

D. Reliability of Machine Vision System

Two tests are performed to determine the reliability of the machine vision system: 1) accuracy in acquisition of coordinates, and 2) accuracy in color discrimination of objects both on the platform level or the second stack level. Twelve (12) trials each containing at least 10 colored objects are to be detected and the coordinates acquired. The average error for each trial is shown in Table IV. On average, the absolute error

for the overall test of the vision system was about 0.19 cm for x-coordinates and about 1.41 cm for the y-coordinates, well within 2 cm tolerance value of accuracy.

TABLE IV. AVERAGE ERROR FOR OBJECT COORDINATES

Trial No.	Average Error		Trial No.	Average Error	
	x	y		x	y
1	0.16	0.76	7	0.59	1.83
2	0	1.16	8	0.25	1.38
3	0.22	1.06	9	0.22	1.33
4	0.13	1.73	10	0.08	1.3
5	0.14	1.67	11	0.07	1.45
6	0.15	2.03	12	0.21	1.22

For the second test, the ability of the machine vision system to properly discriminate colors are tested for the same number of trial with same number of objects. A summary of confusion matrix was constructed as shown in Table V for all tested objects on the first level. The data shows that the colors blue, green and yellow were detected 100% accurately. Notice also that there is a 100% precision for colors blue, green and yellow, and 96% for red. The red color was found to have the least among them all because of the proximity of the red color to the gripper’s color, making it hard to delineate in HSV space. Nevertheless, the gripper is never mistakenly detected as an object.

TABLE V. SUMMARY OF CONFUSION MATRIX FOR 1ST LEVEL OBJECTS

Color	Blue	Green	Red	Yellow
Accuracy	100.00%	100.00%	81.54%	100.00%
True Positive Rate	100.00%	100.00%	100.00%	100.00%
False Positive Rate	0.00%	0.00%	0.94%	0.00%
True Negative Rate	100.00%	100.00%	99.06%	100.00%
False Negative Rate	0.00%	0.00%	0.00%	0.00%
Precision	100.00%	100.00%	96.00%	100.00%

Additional objects were stacked on top of the first level making it a second level stacked object. Similar test for the first level were conducted to test the ability of the system to properly discriminate stacked colored objects. Table VI shows the summary of confusion matrix for detection of objects on the second stack level. The data shows that the green and red were detected 100% accurately. Notice also that there is a 100% precision for colors green and red and about 97% for red and yellow.

TABLE VI. SUMMARY OF CONFUSION MATRIX FOR 2ND LEVEL OBJECTS

Color	Blue	Green	Red	Yellow
Accuracy	96.97%	100.00%	100.00%	96.88%
True Positive Rate	100.00%	100.00%	100.00%	100.00%
False Positive Rate	4.17%	0.00%	0.00%	4.55%
True Negative Rate	95.83%	100.00%	100.00%	95.45%
False Negative Rate	0.00%	0.00%	0.00%	0.00%
Precision	90.00%	100.00%	100.00%	90.91%

VII. CONCLUSION AND RECOMMENDATION

The study was successful in integrating the autonomous robotic arm with fuzzy logic-based joint controller (FLJC) with a machine vision system capable of accurate color discrimination into a color-based sorter system. An improved robotic arm simulator made it possible to tune the membership functions and see the actual effect on the robotic arm's response. Additionally, the end-effector is well accurate enough to have less than 2 cm absolute error. The coordinates of the different target objects with different colors and stacking levels of up to second level as well as the coordinates of the gripper were successfully acquired by means of Processing with SimpleOpenNI and OpenCV library. The overall accuracy of the machine vision system shows that it has the same precision as the end-effector and is at least 95% accurate in properly discriminating colored objects. This extended study has demonstrated that it is capable of sorting even second level stacked color objects. The utilization of the depth data made it possible to determine the height of the colored object in question.

As for improvement, the researchers aim to introduce different controllers such as the hybrid neuro-fuzzy system and genetic algorithm to aid in fine tuning the membership functions and fuzzy rule formulation. Furthermore, the machine vision system can be further improved by applying more advanced color clustering techniques which will eventually allow more colors to be discriminated without ambiguity.

ACKNOWLEDGMENT

The authors would like to thank their family and friends, De La Salle University – Manila (DLSU), De La Salle University – Laguna Campus – Biñan Laguna (DLSU-Laguna) and Department of Science and Technology – Engineering Research and Development for Technology (DOST-ERDT) for funding and helping us to finish this study.

REFERENCES

- [1] J. Rifkin, *The end of work: the decline of the global labor force and the dawn of the post-market era*. New York: Jeremy P. Tacher, 2004.
- [2] Adelhard Beni Rehiara (2011). "Kinematics of AdeptThree Robot Arm," Robot Arms, Prof. Satoru Goto (Ed.), InTech, DOI: 10.5772/17732. Available from: <https://www.intechopen.com/books/robot-arms/kinematics-of-adeptthree-robot-arm>
- [3] OSHA, "OSHA Technical Manual (OTM) | Section IV: Chapter 4 - Industrial Robots and Robot System Safety", 2016. [Online]. [Accessed: 13- Feb- 2016].
- [4] D. D. Ligutan, L. J. S. Cruz, M. C. D. P. Del Rosario, J. N. S. Kudhal, A. C. Abad and E. P. Dadios, "Design and implementation of a fuzzy logic-based joint controller on a 6-DOF robot arm with machine vision feedback," 2017 Computing Conference, London, 2017, pp. 249-257.
- [5] E. P. Dadios, D. J. Williams, "Multiple fuzzy logic systems: A controller for the flexible pole-cart balancing problem," Proc. of the IEEE Robotics and Automation International Conference, Minneapolis, Minnesota USA, April 24-26, 1996. ICRA 1996: 2276-2281.
- [6] A. A. Khalate, G. Leena, G. Ray, "An Adaptive Fuzzy Controller for Trajectory Tracking of Robot Manipulator," Intelligent Control and Automation, 2011, 364-370
- [7] A. R. F. Quiros, A. C. Abad, and E. P. Dadios, "Object locator and collector robotic arm using artificial neural networks," International Conference on Humanoid, Nanotechnology, Information Technology, Communication and Control, Environment and Management (HNICEM), Cebu City, December 2015.
- [8] X. Li, H. Wang, X. Lu, Y. Liu, Z. Chen, and M. Li, "Neural network method for robot arm of service robot based on D-H model," 2017 Chinese Automation Congress (CAC), 2017.
- [9] Elmer P. Dadios, Patrick S. Fernandez, and David J. Williams, "Genetic Algorithm On Line Controller for the Flexible Inverted Pendulum Problem," Journal of Advanced Computational Intelligence and Intelligent Informatics, Vol.10 No.2, 2006
- [10] S. Števo, I. Sekaj, and M. Dekan, "Optimization of Robotic Arm Trajectory Using Genetic Algorithm," IFAC Proceedings Volumes, vol. 47, no. 3, pp. 1748–1753, 2014.
- [11] A. R. F. Quiros, A. Abad, R. A. Bedruz, A. C. Uy, and E. P. Dadios, "A genetic algorithm and artificial neural network-based approach for the machine vision of plate segmentation and character recognition," 2015 International Conference on Humanoid, Nanotechnology, Information Technology, Communication and Control, Environment and Management (HNICEM), Cebu City, December 2015.
- [12] R. R. P. Vicerra, K. K. A. David, A. R. D. Cruz, E. A. Roxas, K. B. C. Simbulan, A. A. Bandala, and E. P. Dadios, "A multiple level MIMO fuzzy logic based intelligence for multiple agent cooperative robot system," TENCON 2015 - 2015 IEEE Region 10 Conference, November 2015.
- [13] E. Maravillas, E.P. Dadios, "Hybrid Fuzzy Logic Strategy for Soccer Robot Game," Journal of Advanced Computational Intelligence and Intelligent Informatics, Vol 8 No. 1, pp 65-71, FUJI Technology Press, January 2004.
- [14] O. F. D. Astilla, J. S. Guerrero, R. S. S. Mendoza, M. T. P. Roxas, A. C. T. Sy, R. R. P. Vicerra, E. P. Dadios, A. R. D. Cruz, E. A. Roxas, and A. A. Bandala, "Obstacle avoidance of hybrid mobile-quadrotor vehicle with range sensors using fuzzy logic control," 2015 International Conference on Humanoid, Nanotechnology, Information Technology, Communication and Control, Environment and Management (HNICEM), Cebu City, December 2015.
- [15] A. Abad, G. Abulencia, W. Pacer, E. Dadios, N. Gunay, "Soccer Robot Shooter Strategy," National Electrical, Electronics, and Computing Conference 2009, Science Discovery Center SM Mall of Asia – December 9-11, 2009
- [16] K.G. B. Leong, S. W. Licarte, G. M. S. Oblepias, E. M. J. Palomado, and E.P. Dadios N. G. Jabson, "The Autonomous Golf Playing Micro Robot: With Global Vision And Fuzzy Logic Controller," International Journal on Smart Sensing and Intelligent Systems, vol. 1, no. 4, pp. 824-841, December 2008.
- [17] E. P. Dadios, R. Baylon, R. De Guzman, A. Floren Lee, and Z. Zulueta, "Vision guided ball-beam balancing system using fuzzy logic," in Industrial Electronics Society, 2000. IECON 2000. 26th Annual Conference of the IEEE 2000, pp. 1973-1978 vol.3.
- [18] Zheng Li, Ruxu Du, Haoyong Yu and Hongliang Ren, "The Inverse Kinematics Solution of a 7 DOF Robotic Arm Using Fuzzy Logic". 7th IEEE Conference on Industrial Electronics and Application. 2012.
- [19] Vishank Bhatia., V. Kalaichelvi, Karthikeyan R., "Application of a Novel Fuzzy Logic Controller for a 5-DOF Articulated Anthropomorphic Robot," IEEE international conference on Research in Computational Intelligence and Communication Networks, Kolkotta, Nov 20-22, 2015.
- [20] M. Mirzadeh, M. Khezri, J. Mahmoodi, H. Karbasi, "Design Adaptive Fuzzy Inference Controller for Robot Arm", IJITCS, vol.6, no.9, pp.66-73, 2014.
- [21] J. R. Sanchez-Lopez, A. Marin-Hernandez and E. R. Palacios-Hernandez, "Visual Detection, Tracking and Pose Estimation of a Robotic Arm End Effector", in Proc. of ROSSUM 2011, Xalapa, Ver., Mexico, June 27-28, pp 41-48, 2011.
- [22] "RobotShop | Robot Store | Robots | Robot Parts | Robot Kits | Robot Toys," RobotShop Blog. [Online]. Available: <http://www.robotshop.com/>. [Accessed: 13-May-2016].
- [23] "Store | Robots | 3D Printers | CNC | Telepresence Robots | R&D," RoboSavvy. [Online]. Available: <https://robosavvy.com/store/>. [Accessed: 23-June-2016].
- [24] "MPU-6050 | TDK," InvenSense. [Online]. Available: <https://www.invensense.com/products/motion-tracking/6-axis/mpu-6050/>. [Accessed: 13-October-2017].

- [25] "Arduino - Home," *Arduino Reference*. [Online]. Available: <https://www.arduino.cc/>. [Accessed: 13-May-2016].
- [26] "Kinect Sensor," *About Processes and Threads (Windows)*. [Online]. Available: <https://msdn.microsoft.com/en-us/library/hh438998.aspx>. [Accessed: 07-March-2018].
- [27] L. A Zadeh, "Fuzzy Sets," *Information and Control*, Vol. 8, pp. 338-353, 1965.
- [28] E. H. Mamdani. "Application of fuzzy algorithms for control of simple dynamic plant". *Proceedings of the Institution of Electrical Engineers*. 121 (12): 1585-1588, 1974.
- [29] S. Sarkar and A. Basu, "Reasoning with uncertainty-Fuzzy Reasoning," Department of Computer Science & Engineering, Indian Institute of Technology, Module 11, Version 1 CSE IIT, Kharagpur, 2009.
- [30] C. Reas and B. Fry, "Processing: programming for the media arts", *Journal AI & Society*, volume 20(4), pp. 526-538, 2006.
- [31] Atduskgreg, "atduskgreg/opencv-processing," *GitHub*, 22-May-2017. [Online]. Available: <https://github.com/atduskgreg/opencv-processing>. [Accessed: 07-March-2018].
- [32] Wexstorm, "wexstorm/simple-openni," *GitHub*. [Online]. Available: <https://github.com/wexstorm/simple-openni>. [Accessed: 07-March-2018].

Mixed Profile Method of Speed and Location for Robotic Arms Motion used for Precise Positioning

Liliana Marilena Matica¹, Cornelia Györödi²,
Helga Silaghi³
Faculty of Electrical Engineering and Information
Technology, University of Oradea
Oradea, Romania

Andrei Silaghi⁴
Faculty of Electronics and Telecommunications
Engineering, University Politehnica Timișoara,
Timișoara, Romania

Abstract—The paper describes a new real-time computation method named **Mixt Profile of Speed (MPS)**, which is used to obtain the value of speed, at every sampling period of time, during the acceleration and deceleration stage, whereas the motion has three stages: 1) acceleration, 2) motion with imposed constant speed, and 3) deceleration. The method will determinate the location of a robotic arm for every sampling period of time. The originality of this new computation method refers to the deceleration stage; it determines an accurate positioning at the end of the motion in a well determinate interval of time. During the forced constant motion stage, the trajectory is imposed and it is linear or circular. The ADNIA algorithm (numerical differential analysis interpolation algorithm) can be implemented at this stage (during the motion with imposed constant speed of the robotic arm) in order to ensure the maximum precision of the computation for the waypoints Cartesian coordinates.

Keywords—Sampling period of time; waypoints; location matrix for a robotic arm; acceleration; deceleration; motion stage; mixt profile of speed; trapezoidal profile; parabolic profile

I. INTRODUCTION

The robotic technology is used in various industry sectors [1], [2], [6], [7]. The problems that arise are complexity of planning a robot motion and real-time computation of the speed and location for the robotic arm motion; these can be computationally intensive and time-consuming. One approach described in textbooks [4], [5] generates a trajectory that satisfies acceleration and speed constraints from a list of waypoints and use linear segments with parabolic blends. The approach is not applicable to automatically generated paths with potentially dense waypoints. In [2] was presented a method to generate the time optimal trajectory along a given path within given bounds on acceleration and speed. The method assumes that the acceleration and speed of individual coordinates are limited. In [3] was proposed to reduce the computation time for path planning of motion for a robotic arm by use to the techniques of caching frequent arm trajectories. In Fig. 1 is presented the block diagram of the robotic system.

The challenges of motion for robotic arm, involve finding the best precision for reaching the end point of motion, also it is important to obtain precise value about required motion time. Both conditions are very difficult to be obtained. The proposed real-time computation method in this paper, namely, **Mixt Profile of Speed** (named shortly **MPS**), defined and

described in this paper, accomplished those two conditions, [10].

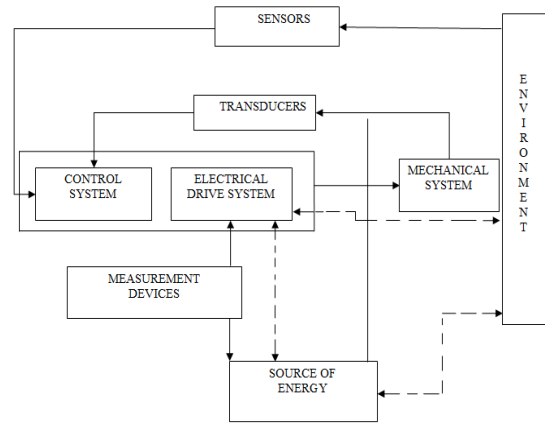


Fig. 1. Block diagram of the robotic system.

Usually, regarding robotic arm motion [1]-[3], the speed variation profile may have a trapezoidal profile, as in Fig. 2:

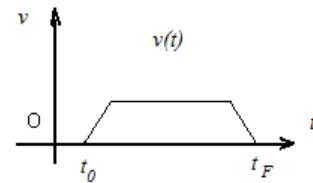


Fig. 2. Trapezoidal profile of speed.

or a parabolic profile, as presented in Fig. 3:

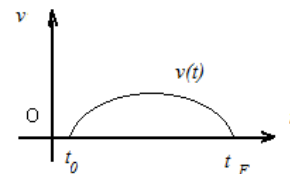


Fig. 3. Parabolic profile of speed.

In previous figures, index O shows the start of motion and index F shows the motion stop.

This paper describes another profile of motion speed, named **MPS (Mixt Profile of Speed)**, as in Fig. 4:

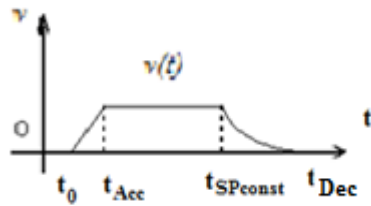


Fig. 4. Mixt profile of speed.

In previous figure, index *Acc* defines the end of acceleration process, index *SPconst* defines the end of motion with constant imposed speed and index *Dec* defines the end of deceleration process.

The motion with a mixt profile of speed supposes three motion stages:

- a) acceleration stage (the first stage of motion);
- b) motion with constant imposed speed (the second stage of motion);
- c) deceleration stage (the third stage of motion).

The MPS profile of speed variation ensures a precise positioning, at the end of the motion. It may be implementing about motion with a well-defined constant speed, on a linear or circular imposed trajectory [9].

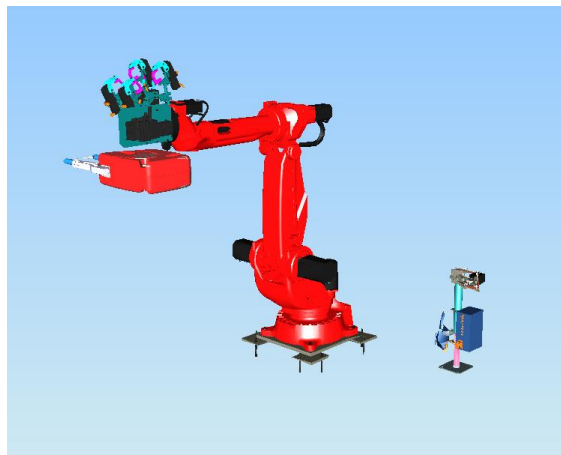


Fig. 5. Industrial robotic arm.

About a robotic arm as shown in Fig. 5, the motion speed value determines the axle components, named: $\vec{p}_x; \vec{p}_y; \vec{p}_z$ of position vector (named: \vec{p} presented in Fig. 5), in the location matrix of a robotic arm [8]:

$$G = \begin{bmatrix} n_x & o_x & a_x & p_x \\ n_y & o_y & a_y & p_y \\ n_z & o_z & a_z & p_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (1)$$

About the location matrix, the others values are the axle components of the orientation versors: $\vec{n}; \vec{o}; \vec{a}$ (a versor is a vector having module equal with 1, because its size is not important; only its orientation is important), Fig. 6 [1]-[3].

This real-time computation method was implemented about positioning pieces, in a flexible manufacturing cell for welding industrial process. The computation method has not been mentioned about robotic arms motion; this paper adapts the computation method for robotic arms motion.

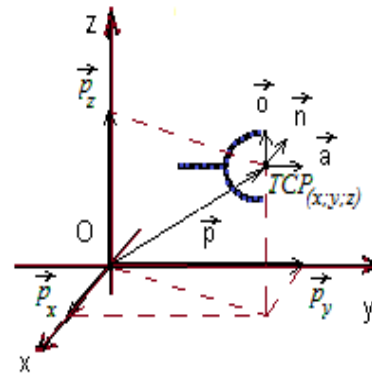


Fig. 6. Position vector for a robotic arm; orientation versors.

The next section presents a new real-time computation method named Mixt Profile of Speed (MPS), which is used to obtain the value of speed, at every sampling period of time, during the acceleration and deceleration stage. The motion with a mixt profile of speed has three stages: acceleration stage, motion with constant imposed speed and deceleration stage. Based on the results of the tests performed, several conclusions are presented in the last section.

II. ACCELERATION STAGE, STAGE OF MOTION WITH CONSTANT IMPOSED SPEED AND DECELERATION STAGE OF MOTION

The real-time computation method named mixt profile of speed (MPS) is described in this paragraph [1].

A. Acceleration Stage

Let consider the maximum value of acceleration for a sampling period of time, named: acc_{Max} ; the initial value of speed at the motion start: v_0 . At every sampling period of time, named T, the speed increases with acc_{Max} value, till reach the imposed value, named: v_{impoz} . It result the value of motion speed, at every sampling period of time, indexed m, (during the acceleration stage):

$$v_m = v_{m-1} + T \cdot acc_{Max} = v_0 + m \cdot T \cdot acc_{Max} \quad (2)$$

This computation must be considered for every axle component of speed:

$$v_{x,m} = v_{x,m-1} + T \cdot acc_{Max} = v_{x,0} + m \cdot T \cdot acc_{Max}$$

$$v_{y,m} = v_{y,m-1} + T \cdot acc_{Max} = v_{y,0} + m \cdot T \cdot acc_{Max} \quad (3)$$

$$v_{z,m} = v_{z,m-1} + T \cdot acc_{Max} = v_{z,0} + m \cdot T \cdot acc_{Max}$$

It result the variation of position vector, for every sampling period of time (named T), during acceleration stage of motion:

$$\begin{aligned}\bar{p}_{x,m} &= \bar{p}_{x,m-1} + T \cdot v_{x,m} \cdot \vec{i} \\ \bar{p}_{y,m} &= \bar{p}_{y,m-1} + T \cdot v_{y,m} \cdot \vec{j} \\ \bar{p}_{z,m} &= \bar{p}_{z,m-1} + T \cdot v_{z,m} \cdot \vec{k}\end{aligned}\quad (4)$$

For example, let consider this example of computation: the location matrix of the robotic arm, at the motion start, index 0 (the values are expressed in millimetres, [mm]) is:

$$G_0 = \begin{bmatrix} n_{x,0} & o_{x,0} & a_{x,0} & p_{x,0} \\ n_{y,0} & o_{y,0} & a_{y,0} & p_{y,0} \\ n_{z,0} & o_{z,0} & a_{z,0} & p_{z,0} \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 3 \\ 0 & 1 & 0 & 5 \\ 0 & 0 & 1 & 7 \\ 0 & 0 & 0 & 1 \end{bmatrix}\quad (5)$$

Let consider the value of sampling period of time $T = 0.01$ seconds [s]. Let consider the initial speed having 0 value (the start motion value of speed) and the maximum value of acceleration: $acc_{Max} = 200$ [mm/s] (identical for every axle); about the first three sampling periods of time, the speed has those OX axle components values:

$$\begin{aligned}v_{x,1} &= v_{x,0} + 1 \cdot T \cdot acc_{Max} = 0 + 1 \cdot 0.01 \cdot 200 = 2 = v_{y,1} = v_{z,1} \\ v_{x,2} &= v_{x,0} + 2 \cdot T \cdot acc_{Max} = 0 + 2 \cdot 0.01 \cdot 200 = \\ &= 4 = v_{y,2} = v_{z,2} \\ v_{x,3} &= v_{x,0} + 3 \cdot T \cdot acc_{Max} = 0 + 3 \cdot 0.01 \cdot 200 = 6 = v_{y,3} = v_{z,3}\end{aligned}\quad (6)$$

The OX axle components of position vector, for the first three sampling period of time, during acceleration stage of motion have those values:

$$\begin{aligned}\bar{p}_{x,1} &= \bar{p}_{x,0} + T \cdot v_{x,1} \cdot \vec{i} = (3 + 0.01 \cdot 2) \cdot \vec{i} = 3.02 \cdot \vec{i} \\ \bar{p}_{x,2} &= \bar{p}_{x,1} + T \cdot v_{x,2} \cdot \vec{i} = (3.02 + 0.01 \cdot 4) \cdot \vec{i} = 3.06 \cdot \vec{i} \\ \bar{p}_{x,3} &= \bar{p}_{x,2} + T \cdot v_{x,3} \cdot \vec{i} = (3.06 + 0.01 \cdot 6) \cdot \vec{i} = 3.12 \cdot \vec{i}\end{aligned}\quad (7)$$

In previous relations (rel.7), versor \vec{i} is the OX axle versor (it has the module equal with 1 value and the orientation along the positive sense of this axle); also, in this paper, versor \vec{j} is the OY axle versor and versor \vec{k} is the OZ axle versor.

Similar computation (as rel.7) must be implemented about OY and OZ axle components of position vector; it results next values: $\bar{p}_y = 5.12 \cdot \vec{j}$ and $\bar{p}_z = 7.12 \cdot \vec{k}$; so, after 3 periods of time, the location matrix of the robotic arm is:

$$G_3 = \begin{bmatrix} 1 & 0 & 0 & p_{x,3} \\ 0 & 1 & 0 & p_{y,3} \\ 0 & 0 & 1 & p_{z,3} \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 3.12 \\ 0 & 1 & 0 & 5.12 \\ 0 & 0 & 1 & 7.12 \\ 0 & 0 & 0 & 1 \end{bmatrix}\quad (8)$$

In the acceleration stage, about the last sampling period of time, the considered acceleration may be less then maximum value of acceleration, in order to reach the imposed constant speed. Let consider the number of sampling period of time required for acceleration stage, named: N_{Acc} ; it results the acceleration for the last sampling period of time:

$$acc_{N_{Acc}} = v_{impoz} - (N_{Acc} - 1) \cdot acc_{Max}\quad (9)$$

For example, if the maximum acceleration value is 200 mm/s; the value of sampling period of time $T = 0.01$ s. The speed at the motion start has 0 value and the constant imposed speed is 9 mm/s, it results the speed values, during the acceleration stage, at every sampling period of time, having those values: 0; 2; 4; 6; 8; 9 mm/s (for the last period of time, during acceleration stage).

The acceleration has the value of 100mm/s, which is the maximum value possible). It results the number of sampling period of time for acceleration stage (named: N_{Acc}), about this example is 5.

B. Motion with Imposed Constant Speed

The next motion stage is performed with the constant imposed value of speed, named: v_{impoz} . Usually, motion upon an imposed well-defined trajectory (linear or circular) is executed with an imposed constant speed.

Let consider an imposed linear trajectory and the axle components of the imposed constant speed: $v_{X,impoz}$; $v_{Y,impoz}$; $v_{Z,impoz}$. It results the axle steps, executed by the robotic arm, named: δ_x ; δ_y ; δ_z , at every sampling period of time (named T), during this motion stage:

$$\begin{aligned}\delta_x &= T \cdot v_{X,impoz} \\ \delta_y &= T \cdot v_{Y,impoz} \\ \delta_z &= T \cdot v_{Z,impoz}\end{aligned}\quad (10)$$

It results the axle components of position vector, indexed l (index l starts with 1 value), for every sampling period of time, during this motion stage:

$$\begin{aligned}\bar{p}_{x,l} &= \bar{p}_{x,l-1} + \delta_x \cdot \vec{i} = \bar{p}_{x,N_{Acc}} + l \cdot \delta_x \cdot \vec{i} \\ \bar{p}_{y,l} &= \bar{p}_{y,l-1} + \delta_y \cdot \vec{j} = \bar{p}_{y,N_{Acc}} + l \cdot \delta_y \cdot \vec{j}\end{aligned}\quad (11)$$

$$\vec{p}_{z,m} = \vec{p}_{z,m-1} + \delta_z \cdot \vec{k} = \vec{p}_{z,N_{Acc}} + l \cdot \delta_z \cdot \vec{k}$$

In previous relation (rel.11), $\vec{p}_{x,N_{Acc}}$; $\vec{p}_{y,N_{Acc}}$; $\vec{p}_{z,N_{Acc}}$ are the initial values of axle components of position vector, about the second stage of motion; those values are identical with position vector axle component values, at the end of the first stage of motion, the acceleration stage.

For example, let considers those values: $\vec{p}_{x,N_{Acc}} = 10 \cdot \vec{i}$ [mm]; $\vec{p}_{y,N_{Acc}} = 100 \cdot \vec{j}$ [mm]; $\vec{p}_{z,N_{Acc}} = 1000 \cdot \vec{k}$ [mm]. Let considers the axle steps having those values of millimetres: $\delta_x = 1.11$ [mm]; $\delta_y = 1.22$ [mm]; $\delta_z = 1.33$ [mm]. It results the axle components of position vector, at every sampling period of time (index l), during second motion stage:

$$\begin{aligned} \vec{p}_{x,l} &= \vec{p}_{x,N_{Acc}} + l \cdot \delta_x \cdot \vec{i} = 10 \cdot \vec{i} + l \cdot 1.11 \cdot \vec{i} \\ \vec{p}_{y,l} &= \vec{p}_{y,N_{Acc}} + l \cdot \delta_y \cdot \vec{j} = 100 \cdot \vec{j} + l \cdot 1.22 \cdot \vec{j} \\ \vec{p}_{z,l} &= \vec{p}_{z,N_{Acc}} + l \cdot \delta_z \cdot \vec{k} = 1000 \cdot \vec{k} + l \cdot 1.33 \cdot \vec{k} \end{aligned} \quad (12)$$

Considering rel.12, it may be computed the position vector, at every sampling period of time, during the motion with constant imposed speed (the second stage of motion). For example, after 200 sampling period of time, the axle component of position vector have the values:

$$\begin{aligned} \vec{p}_{x,200} &= 10 \cdot \vec{i} + 200 \cdot 1.11 \cdot \vec{i} = (10 + 222) \cdot \vec{i} = 121 \cdot \vec{i} \\ \vec{p}_{y,200} &= 100 \cdot \vec{j} + 200 \cdot 1.22 \cdot \vec{j} = 344 \cdot \vec{j} \\ \vec{p}_{z,200} &= 1000 \cdot \vec{k} + 200 \cdot 1.33 \cdot \vec{k} = 1266 \cdot \vec{k} \end{aligned} \quad (13)$$

It result the location matrix for every sampling period of time. In order to work with different index value about location matrix during all three motion stages, it must be considered: N_{Acc} ; so, during motion with constant speed, (index l), the location matrix is:

$$G_{l+N_{Acc}} = \begin{bmatrix} 1 & 0 & 0 & p_{x,l} \\ 0 & 1 & 0 & p_{y,l} \\ 0 & 0 & 1 & p_{z,l} \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 121 \\ 0 & 1 & 0 & 344 \\ 0 & 0 & 1 & 1266 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (14)$$

In previous relations (rel.11; rel.12; rel.14), the index l goes from 1 (one) value to $N_{SPconst}$. The value of $N_{SPconst}$ results from conditions defined about others two motion stage (as it follows the explanations). So, the distance performed during the acceleration stage (the first stage of motion) and deceleration stage (the third and the last stage of motion) must be computed.

The start point and the stop point, about motion upon the imposed linear trajectory, must be defined; it results the entire

distance of motion. Deducting from entire distance of motion, it result the distance of the second motion stage (motion with constant imposed speed). Considering this distance value and the value of imposed constant speed, it results the time required for the second motion stage. It must be divided this time value to the value of sampling period of time; so, it results $N_{SPconst}$ (last index value about second stage of motion).

During first and second motion stage, the required number of sampling period of time is: $N_{Acc} + N_{SPconst}$.

C. Deceleration Stage of Motion

About the third stage, the deceleration stage of motion (the last motion stage), the variation of speed (indexed q) is:

$$v_q = v_{q-1} - T \cdot dec_q \quad (15)$$

The deceleration value, named dec_q , has not a constant value, it has a decreasing value, in purpose to ensure a precise positioning, at the motion end (the b value is a constant and adjusts the decreasing of speed with several others characteristics of motion):

$$dec_q = \frac{b}{q^2} \quad (16)$$

For example, let consider $b=50$, it results the values of deceleration: $dec_1=50$ mm/s²; $dec_2=12.5$ mm/s²; $dec_3=5.55$ mm/s²; $dec_4=3.11$ mm/s²; $dec_5=2$ mm/s².

During the deceleration stage, the speed decreases, till it reaches the 0 value.

The time lapse, during the deceleration stage, is very well defined, linked with the number of sampling period of time required for the deceleration stage, named N_{Dec} (index q goes from 1 value to N_{Dec} value). It results from those conditions: it begin with v_{impoz} and end with 0 value of speed, about deceleration stage:

$$v_{impoz} = v_{q=0} \quad (17)$$

$$v_{N_{Dec}} = 0 \quad (18)$$

The computation method must consider the different values of speed axle components; so, the previous conditions must be applied for the maximum value of speed axle component:

$$v_{q=1} = \max(v_{x,impoz}; v_{y,impoz}; v_{z,impoz}) \quad (19)$$

Considering the previous computation example, this value is the maximum from: 111 mm/s; 122 mm/s and 133 mm/s.

The number (named N_{Dec}) of sampling period of time, required for deceleration stage, is:

$$\max(v_{x,impoz}; v_{y,impoz}; v_{z,impoz}) = \frac{b}{(N_{Dec})^2} \quad (20)$$

The resulting number must be the next integer value, greater then: $\sqrt{\frac{b}{v_{impoz}}}$ because the end point approach must be done with a very small deceleration.

Considering the previous example of deceleration for OX axle component of speed, the computation of this axle component starts with those relations:

$$\begin{aligned} v_{x,2} &= v_{x,impoz} - T \cdot dec_{x,1} = v_{x,impoz} - T \cdot 50mm/s^2 \\ v_{x,3} &= v_{x,2} - T \cdot dec_{x,2} = v_{x,2} - T \cdot 12.5mm/s \\ v_{x,4} &= v_{x,3} - T \cdot dec_{x,3} = v_{x,3} - T \cdot 5.55mm/s \end{aligned} \quad (21)$$

The initial value of speed for deceleration stage (index $q=0$) is the imposed constant speed named: v_{impoz} (for motion on the imposed trajectory).

The similar computation, about OX axle component of speed, must be applied about OY and OZ axle component of speed, so, the computation of axle components of speed is:

$$\begin{aligned} v_{x,q} &= v_{x,q-1} - T \cdot dec_{x,q} \\ v_{y,q} &= v_{y,q-1} - T \cdot dec_{y,q} \\ v_{z,q} &= v_{z,q-1} - T \cdot dec_{z,q} \end{aligned} \quad (22)$$

Index q goes from 1 value to N_{Dec} value, (computed with rel.21). For each speed axle component, the last value of deceleration must be adjusted to the proper value, in order to obtain 0 value of speed.

It results the computation of position vector (index q):

$$\begin{aligned} \bar{p}_{x,q} &= \bar{p}_{x,q-1} + T \cdot v_{x,q} \cdot \bar{i} \\ \bar{p}_{y,q} &= \bar{p}_{y,q-1} + T \cdot v_{y,q} \cdot \bar{j} \\ \bar{p}_{z,q} &= \bar{p}_{z,q-1} + T \cdot v_{z,q} \cdot \bar{k} \end{aligned} \quad (23)$$

The forth column of location matrix for robotic arm contains those axle components of position vector (as it was explained previously), the location matrix during deceleration stage is:

$$G_{q+N_{Acc}+N_{SPconst}} = \begin{bmatrix} 1 & 0 & 0 & p_{x,q} \\ 0 & 1 & 0 & p_{y,q} \\ 0 & 0 & 1 & p_{z,q} \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (24)$$

The motion with mixt profile of speed needs $N_{Acc} + N_{SPconst} + N_{Dec}$ sampling periods of time (as it was explained).

III. CONCLUSIONS

The paper defines and describes the real-time computation method named mixt profile of speed (speed variation), for motion on an imposed linear trajectory.

The motion implementing mixt profile of speed has three stages: acceleration stage; motion with constant imposed speed; deceleration stage.

The paper shows the real-time computation of mixt profile of speed, for a robotic arm motion, on a linear trajectory. For each of the three stages of motion, the position vector may be computed with relations (4), (12) and (23), thus it results the location matrix of the robotic arm relations (5), (14) and (24).

The proposed MPS method implements the maximum computation precision, for robotic arm motion, upon an imposed linear trajectory, with a constant imposed speed.

The method offers the best precision for reaching the end point of motion; also it obtains a precise value for required motion time.

REFERENCES

- [1] L.M.Matica, H. Oros, "Speed Computation for Industrial Robots Motion Followed by Accurate Positioning. International Journal of Computers", Communications & Control, ISSN 1841-9836, 12(1):76-89, February 2017.
- [2] T. Kunz, M. Stilman, "Time-Optimal Trajectory Generation for Path with Bounded Acceleration and Velocity", Robotics: Science and Systems 8th Conference, July 9-13, 2012, Sydney, Australia, ISBN 978-0-262-51968-7, pp. 209-216.
- [3] K. Krishnaswamy, J. Sleeman, T. Oates, "Real-Time Path Planning for Robotic Arm", Proceedings of the 4th International Conference on Pervasive Technologies Related to Assistive Environments, May 25-27, 2011, ISBN: 978-1-4503-0772-7, doi:10.1145/2141622.2141635.
- [4] J.J. Craig, "Introduction to Robotics: Mechanics and Control (3rd Edition)", Prentice Hall, 2004, ISBN-10: 0133489795.
- [5] B. Siciliano, L. Sciavicco, L. Villani, G. Oriolo, "Robotics: modelling, planning and control", Springer 2009, ISBN 978-1-84628-641-4.
- [6] Z. Shiller, H.H. Lu. "Computation of path constrained time optimal motions with dynamic singularities", Journal of Dynamic Systems, Measurement, and Control, Vol 114, pp. 34 - 40, 1992.
- [7] M. Zucker, N. Ratliff, A.D. Dragan, M. Pivtoraiko, M. Klingensmith, C. M. Dellin, J.A. Bagnell, S.S. Srinivasa, "CHOMP, Covariant hamiltonian optimization for motion planning", Int. Journal of Robotics Research, 2012.
- [8] L.M.Matica, C.Gyorodi, H.Silaghi, S.Abrudan Caciora, "Real Time Computation for Robotic Arm Motion upon a Linear or Circular Trajectory", International Journal of Advanced Computer Science and Applications, Vol.9, Issue 2, 2018.
- [9] S.Dale, Helga Silaghi, D. Zmaranda, U. Rohde, "Intelligent Design Environment for Second-Order Positioning Systems", International Journal of Computers Communications and Control, ISSN 1841-9836, Vol.10, No 1, 2015.
- [10] C.R. Costea, Helga Silaghi, D. Zmaranda, M.A. Silaghi, "Control Systems Architecture for a Cement Mill based on Fuzzy Logic", International Journal of Computers Communications and Control, ISSN 1841-9836, Vol.10, No 2, 2015.

A Novel E-Mail Network Evolution Model based on user Information

Lejun ZHANG

College of Information Engineering
Yangzhou University, Yangzhou, 225127
Yangzhou, China

Chunhui ZHAO

College of Information Engineering
Yangzhou University, Yangzhou, 225127
Yangzhou, China

Tongxin ZHOU

College of Computer Science and Technology
Harbin Engineering University
Harbin, China

Zilong JIN

School of Computer and Software
Nanjing University of Information Science and Technology
Nanjing, China

Abstract—E-mail is one of the main means of communication in society today, and it is a typical social network. Studying the evolution of the social network structure by constructing an e-mail network evolution model is of great significance to the literature. In this paper, we first analyze the e-mail network by constructing an e-mail network communication model; this mainly includes analysis of the structure of the e-mail network and analysis of the user information in the e-mail network; then, we propose an e-mail network evolution model based on the characteristics of user information and give the specific evolutionary steps; finally, the simulation experiments are carried out to analyze the characteristics of the model. Experiments show that the nodes are characterized by a power-law distribution, and compared with other models; the model is closer to the real network, so it has important practical significance.

Keywords—Information characteristics; e-mail; network evolution; complex network

I. INTRODUCTION

With the rapid development of information technology, people's lives have been fully integrated into a complex network world, tangible, intangible, various, ubiquitous. Networks are like large systems; each node in the network is a different element in the system, and the relationship between different elements forms the edge between nodes. So scientists want to find a certain law to construct the network topology and thus to better understand the network, finally determining the value of the network. Application of a complex network analysis method can better reveal the characteristics of the network; it has important significance for network formation and expansion, information dissemination and other research. With the development of computer technology and the Internet, many scholars at home and abroad have studied the network model in many ways. In different fields, the particularity of the structure of the network model is also different.

With the rapid development of networks and computers, people have studied networks in a more profound and comprehensive way. They have found that regular networks

are not applicable to the universality of real networks. Then, the first to initiate change were Watts and Strogatz [1], who proposed a WS network named after them. In the process of network generation, the edge will be randomized to reconnection, when the network reaches a certain scale; the average path of the network is small, and the clustering coefficient is high; the network has the characteristics of a "small world"; subsequently, Newman and Watts [2] found that the defects of the WS network, random network rewiring, may lead to more independent nodes appearing; in response to this phenomenon, they used adding edges instead of randomized reconnection. Later, Barabasi and Albert [3], [4] found that nodes in a network have the power law characteristic through a large number of actual networks; then they put forward a scale-free network, which is known as the BA network; then, more and more models were proposed; Flammini [5] propose a criterion of network growth that explicitly relies on the ranking of nodes according to any prestige measure; Sun et al. propose a model driven by events and interests [6]; Alireza et al. validate the significance of betweenness centrality in the evolution of research collaboration networks [7]; Barrat et al. study the complex weighted network [8], [9]; Sun et al. propose a topological evolution to simulate social activities [10]; Song et al. build up a class of edge-growing network models and provide an algorithm for finding spanning trees of edge-growing network models [11]; Huang et al. use spanning trees and other graphs to illustrate some results and phenomenon and try expressing mathematically key notions from researching scale-free networks [12]; a likelihood analysis is provided about evaluation models by Wang et al [13]; Zou et al. propose an evolving network model growing fast in units of module [14]; article [15] finds that with the increment of network interdependence, the evolution of cooperation is dramatically promoted on the network playing Prisoner's Dilemma, and the cooperation level of the network playing Snowdrift Game reduces correspondingly; Zhuang et al. study the problem of maximizing influence diffusion in a dynamic social network [16]; Lhan et al. present a framework for modeling community evolution prediction in social networks [17].

This paper proposes an e-mail communication network named “User-Information-keyword” through network analysis of an Enron dataset and analyzes the structure and user information of an e-mail network; finally, an e-mail network evolution model based on user information is proposed. Through the above methods, the research on e-mail networks for a certain range of people has a certain significance and value in reflecting the communication between people in real society.

II. E-MAIL NETWORK ANALYSIS: A CASE STUDY OF AN ENRON DATASET

A. Introduction to the Enron Dataset

In this paper, we use the Enron data set; the Enron dataset is from the former energy giant Enron Corp in the United States, which went bankrupt because of bad management and a bribery scandal. The United States Department of Justice conducted an investigation of Enron Corp, including their email. Later, MIT purchased the dataset in order to implement the CALO project, SRI Laboratory Start to sort out the mail; each mail will be stored in the format of the SMTP protocol, and the mail attachments are removed. The email mainly includes the following: mail ID, post time, send mailbox, receive mailbox, mail theme, mail content and so on. The paper uses the Enron_20150507 version; the dataset contains 150 users, with a total of 517374 emails.

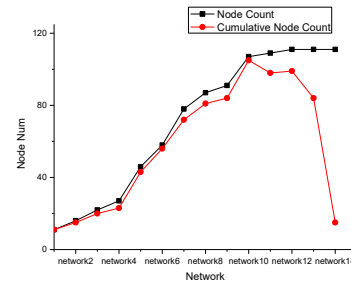
The Enron dataset is huge; after treatment by researchers, the dataset still has many mistakes or useless information. For example, the mail records email information from January 1999 to June 2002, but some emails are from 1980 and 2044; we delete these; in addition, the dataset contains 150 user mail folders; each folder corresponds to a user, but some folders do not correspond to the right email address; for example, for the folder crandell-s, the e-mail address is *.crandall@enron.com; the user name and e-mail address do not match, and we need to manually change them; it is not possible to avoid by getting the corresponding email address in later; and there are other situations like message ID repeats or sending their own e-mail. The above examples show the vulnerabilities of email, which must be addressed.

Because the main analysis contains the internal mail information records of 150 employees of Enron, we delete email addresses that do not belong to the company’s internal e-mail address and keep the communication records of the 150 Enron employees. The messages contain many mass emails or copied emails; because the communication record should be transformed into a graph or matrix later, the paper separates the mass addresses and copy addresses; if a message corresponds to multiple e-mails, this will only increase the total number of communication records, but the communication between employees remains unchanged, and does not affect the subsequent analysis.

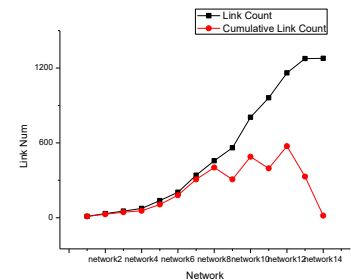
B. E-mail Network Structure Analysis

In this paper, we use the ORA software developed by the CASOS Laboratory of Carnegie Mellon University, which can transform data in a corresponding format to the network structure diagram, and we can perform dynamic analysis of the data.

In order to analyze the generation process of the network, this paper divides the network into different time periods according to the sending time. The Enron dataset is from January 1999 to June 2002; we take three months for a period of time; to establish a network of dynamic books, each book time corresponds to a communication network, for a total of 14 networks. For example, the time from January 1999 to April 1999 corresponds to network1, followed by analogy.



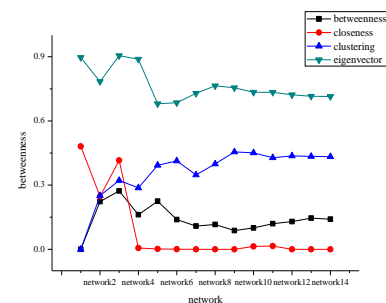
(a) Enron network nodes graph



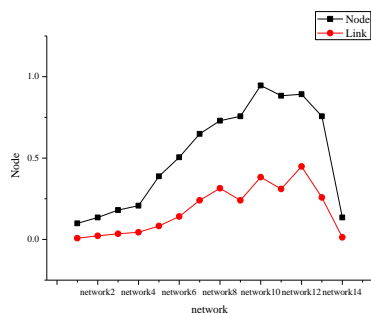
(b) Enron network edges graph

Fig. 1. Enron network diagram.

Fig. 1(a) and (b) show the distribution of the number of nodes and edges in the whole network, respectively, and in each period network as time goes on in the e-mail network; it can be seen from the graph that the number of nodes and the number of edges in the whole network is an increasing process; in the case of network1-network11 in the figure, there are new nodes joining the network and the activity of nodes in each time book is very high; then, the trend of the nodes in the network tends to be gentle; as shown in figure network11-network14, the total number of nodes and edges changes slowly, and the activity of each period decreases sharply. It is during the time of the outbreak of the Enron Corp crisis.



(a) Enron network centrality analysis



(b) Enron information occupancy

Fig. 2. Enron network parameter analysis.

As shown in Fig. 2(a), betweenness, closeness, clustering, eigenvector express the four centralities of network analysis. The process of the whole network from growth to smoothness can be seen. Fig. 2(b) shows the occupancy of nodes and edges in each period network, which represents the activity of 111 users in the network. It can be seen from the figure that user activity exceeds 50% from network5 to network13; the activity of users is very high at this time.

Based on the above analysis, this section summarized the following points: (1) the process of mail network growth; in the case of a certain network size, the node change is increased rapidly at first and then the rate of increase tends to slow until the number of nodes does not change; (2) in the case of a certain network size, traffic in the network begins to increase, reaches a threshold, and the amount of communication is maintained at a fluctuating value; (3) the structure of the mail network will stabilize after the traffic is stable and the network structure has been formed.

C. Email Information Analysis

In order to avoid the problem of sparse data in the analysis of mail messages, the paper selects the data from network5 to network13 to analyze. Analysis from the previous section shows that communication traffic is very stable, and user activity is higher; the network structure is stable; the messages of this period can reflect the basic features of each user.

1) "User-Information feature-Keyword" Model

In order to analyze the message sending behavior of email users, this paper establishes a model of "User-Information feature-Keyword". As shown in Fig. 3, User1 sends an email to User2; then, a directed edge is generated from User1 to the Feature Network; the Feature Network generates a feature vector $(F_1, F_2, F_3 \dots)$; F_1 represents the weight of this type feature in the current mail; the higher the proportion of the weight in the whole feature, the more the text tends to this kind of characteristic; then a directed edge is generated between the feature network and User2. The relationship model between the user and information feature is constructed. In the paper, the feature class is composed of the key words, and the different feature classes are composed of key words with different attributes; keywords with similar attributes are regarded as the same feature class. Thus, the mail communication model of "User-Information feature-Keyword" is constructed.

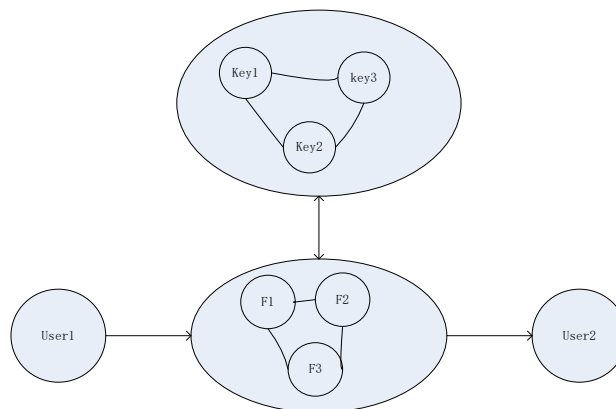


Fig. 3. Schematic diagram of mail communication network

Due to the large number of messages, how to extract the user information effectively is a problem. Because the user information characteristics are excavated from the emails, the problem of extracting the user information characteristics is changed to how to extract the text feature, that is, keyword. In the mail message, the subject and the text are the two main parts of the message content. The subject can often show the main content of the text, but the subject is generally composed of a short sentence; the number of words are usually not more than 10; it cannot express the whole message; in this paper, through the integration of subject and text information, the most representative keywords in an email are extracted, and the feature class is formed by the similarity between these keywords.

2) Keyword Extraction Method

The paper uses the DF (Document Frequency) method to select features. The DF algorithm is a relatively simple feature selection algorithm; it refers to the amount of text including words in the dataset. The general document frequency algorithm sets the threshold to remove the feature according to which the document frequency is particularly high or the document frequency is particularly low; these two features represent the two extremes, that is, "no representation" and "no use". The evaluation function of the DF is a text correlation method, which makes the establishment of stoplist critical, and this will directly affect the classification feature. Because the dataset is English, the stop word is (<http://jmlr.csail.mit.edu/papers/volume5/lewis04a/al1-smart-stop-list/english.stop>). The DF method is also flawed; for example, some scarce words can reflect the good performance of certain characteristics, but they are removed because the amount is too small. In order to reduce the interference caused by these types of factors, in this paper, in addition to considering the text feature, we also take into account the key words of the subject to achieve a complementary and to maximize the expression of e-mail information characteristics.

Based on the above methods, we select 150 key words as text feature vectors; although the stop word is considered in the algorithm, for the problem of data, there are a lot of useless words, such as "pm", "st" and some numbers; by deleting these useless words, we get a total of 92 key words, and the following Table I shows the top of the list of words:

TABLE I. KEYWORD LIST

KeyWord	Count	KeyWord	Count
meeting	2056	report	887
credit	1521	master	882
trader	1427	california	859
presentations	1285	power	855
gas	1238	capacity	655
enron	1225	trading	648
responses	1062	notification	500
agreement	1027	company	489
new	908	summary	446
energy	903	draft	445

Enron was the biggest energy giant in the United States; it was one of the most active stocks traded in the 2000-2002 the period; from the first few keywords, “meeting”, “credit”, “trader” and “gas” fit the company’s image; taking as an example “presentations”, from the words “enron”, “responses” and “agreement”, we can see there is a great deal of work communication in the emails. These keywords can be used to express the email information characteristics to a certain extent.

There are many keywords that will impact the analysis of sparse data, so the K-means method is used to cluster the keywords; the parameters of the algorithm are as follows: initialize type = Random cluster, cluster method = Euclidean distance, cluster num = 8; we classify all keywords into 8 feature classes by clustering, which eliminates the interference caused by sparse data. After clustering, we name the eight kinds of clusters as the feature class, and the following Table II shows the key words of the eight feature classes:

TABLE II. FEATURE CLASSES

Feature	Size	Member	Primary Key
Feature1	5	america,received,north,entity,master	america
Feature2	12	thanks,price,energy,deals,deal,day,group,need,know,sara,gas,time	energy
Feature3	18	texas,type,approved,isda,contract,distributed,transaction,executed,products date,Stephanie,copies,border,susan,exchange,financial,effective,counterparty	products
Feature4	2	enron,hou	enron
Feature5	11	work,forwarded,john,trading,power,business,subject,new,market,pm,mark	trading
Feature6	4	meeting,office,west,mike	meeting
Feature7	50	eoI,carol,let,fax,going,want,tomorrow,morning,phone,following,jones,don,does Email,Monday,think,forward,change,report,information,through,mail,smith,use,question,available,contact,legal,credit,taylor,like,list,sent,look,regards,regarding,attached,week,houston	credit

		street,tana,today,Friday,working,company,david,just,make	
Feature8	2	agreement,corp	corp

3) Analysis of User Participation in Feature Class

In the paper, we analyze the situation of different users participating in the feature class, from which we can find that users have different preferences for different classes. As shown in Table III, the paper compares several users with a large amount of communication.

TABLE III. COMMUNICATION RATIO OF DIFFERENT FEATURE CLASSES

	1	2	3	4	5	6	7	8
jones-t	5.16 %	7.39 %	13.6 7%	28.2 5%	11.4 8%	0.54 %	28.7 1%	4.80 %
shacklet on-s	5.33 %	20.7 6%	7.50 %	16.7 0%	11.8 7%	2.22 %	31.7 8%	3.84 %
grigsby-m	0.63 %	21.5 1%	4.96 %	14.4 1%	17.2 3%	15.2 5%	24.6 9%	1.32 %
stclair-c	1.19 %	10.2 3%	8.01 %	16.8 6%	12.4 8%	1.65 %	46.5 5%	3.04 %
williams -w3	0.38 %	47.0 0%	2.67 %	2.51 %	9.94 %	3.22 %	33.9 9%	0.30 %
phanis-s	23.3 5%	5.16 %	40.5 7%	7.81 %	2.13 %	0.17 %	8.75 %	12.0 3%
delainey -d	1.04 %	11.2 1%	1.33 %	22.8 0%	24.6 1%	4.16 %	31.6 6%	3.18 %
taylor-m	1.84 %	9.94 %	9.00 %	20.5 2%	22.5 3%	2.67 %	29.9 5%	3.56 %
keiser-k	0.49 %	10.2 0%	1.11 %	11.7 7%	4.19 %	31.6 8%	31.5 9%	8.97 %
symes-k	0.21 %	31.3 2%	4.75 %	4.50 %	11.2 8%	4.10 %	42.5 5%	1.30 %
whalley-l	0.00 %	18.5 3%	3.44 %	10.1 2%	41.3 6%	0.48 %	24.9 3%	1.15 %
perlingie re-d	13.3 6%	6.95 %	8.01 %	15.8 4%	5.39 %	0.87 %	43.9 3%	5.65 %
heard-m	23.3 7%	7.98 %	26.6 2%	7.27 %	5.14 %	0.62 %	16.5 4%	12.4 7%
skilling-j	3.18 %	5.98 %	2.73 %	9.41 %	8.71 %	13.5 0%	52.7 3%	3.78 %
haedicke -m	2.18 %	8.07 %	5.92 %	22.7 1%	27.8 7%	3.08 %	28.1 2%	2.06 %

From the table, we can see that most users show big differences in feature classes in communication. We can see that stclair-c, symes-k, perlingiere-d and skilling have great interest in the seventh type of features, and the seventh feature class occupies more than 40% of the proportion; similarly, the communication traffic of williams-w3 occupies 47% of the proportion in the second feature class; the communication traffic of whalley-l occupies 41.36% of the proportion in the fifth feature class; for other users, such as jones-t, delainey-d and keiser-k, the communication traffic occupies about 30%; other users can reach more than 20%.

Fig. 4 shows the communication ratio of different feature classes by users with large communication traffic. We can see that several users in the linear graph have obvious uplift in a certain feature class while having a low proportion in other feature classes; in addition, some users’ line drawings have two distinct bumps, indicating that the users are interested in more than one category, while they also have a low interest in other feature classes. Fig. 4 is a very good response to the feature of email messages.

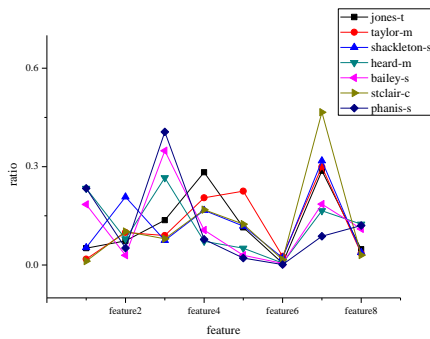


Fig. 4. Communication ratio of different users in different feature classes.

The paper selects the user with the largest communication traffic to analyze, for example, user jones-t and stclair-c; his communication table and corresponding feature relation graph is as Fig. 5 and 6, and Tables IV and V.

TABLE IV. COMMUNICATION RELATIONSHIP OF STCLAIR-C

name	shackleton-s	taylor-m	Jones-t	bailey
stclair-c	299	235	227	269

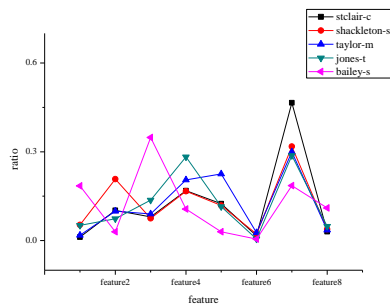


Fig. 5. Comparison between user Stclair-C and its communication characters.

TABLE V. COMMUNICATION RELATIONSHIP OF JONEST-T

name	Taylor	shackleton-s	heard-m	bailey-s	stclair-c
jones-t	402	313	302	269	261

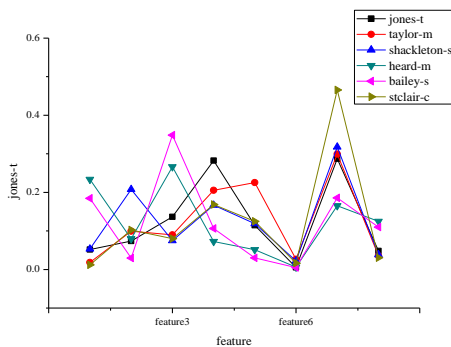


Fig. 6. Comparison between user jones-t and its communication characters.

As shown in Fig. 5 and 6, which show the communication relationship of users with larger communication traffic, the black line represents the user being compared in the figure; other colors represent the users involved in the comparison.

From the figure, we can find that there is more than one feature class that takes a larger proportion between users participating in the comparison and users being compared.

Through the above analysis, we find that the “User-Information features-keyword” model can effectively analyze the e-mail dataset; we can find the user information characteristics and what they are interested in to classify the users. These findings can provide a theoretical and data basis for the identification, recommendation and evolution of the network in the future.

III. E-MAIL NETWORK EVOLUTION MODEL BASED ON USER INFORMATION CHARACTERISTICS

In the last chapter, we find that there is a correlation between the user’s communication intensity and the user’s information characteristics through the “User-Information Feature-keyword” model, so this paper proposes an E-mail Network Evolution Model Based on User Information, the abbreviated UIEM model. The UIEM model is proposed based on Undirected Weighted Network Model (BBV) and Local World Model; it considers the node selection idea of BBV and the group of LocalWorld.

A. Related Knowledge Theory

Define abbreviations and acronyms the first time they are used in the text, even after they have been defined in the abstract. Abbreviations such as IEEE, SI, MKS, CGS, sc, dc, and rms do not have to be defined. Do not use abbreviations in the title or heads unless they are unavoidable.

1) BBV Network Model

The BBV [18] model considers the edge weight based on the scale-free network; it can be represented by the adjacency matrix of the network; w_{ij} represents the weight of the edge between node i and node j ; if the network is undirected, the matrix is symmetric, $w_{ij} = w_{ji}$; if the network is directed, w_{ij} and w_{ji} are considered separately. Here, we introduce the BBV model.

In the BBV model, the degree of the node is called the strength or tendency of the node, in which the strength of the node s_i is defined as:

$$s_i = \sum_{j \in Neighbor(i)} w_{ij} \tag{1}$$

The $Neighbor(i)$ represents the neighbor node set of node i ;

The evolution rules of BBV model are as follows:

a) The initial network: the network is a unity coupling network with m_0 nodes; each edge is given the initial weight w_0 .

b) The growth of the network: In each time step, there is a new node N with $m(m < m_0)$ edges; the new node selects a node from the original network to connect according to a certain probability; the probability that the node is selected is as follows:

$$\Pi = \frac{s_i}{\sum_j s_j} \tag{2}$$

c) The dynamic evolution of edge weights: in the BBV model, each time the new edge (n, i) is given the initial weights of w_0 , and to facilitate the analysis, the new edges will only impact the weight between node i and its neighbor j , so the weight of the edge between node i and node j readjusts as:

$$w_{ij} \rightarrow w_{ij} + \Delta w_{ij} \quad (3)$$

$$\Delta w_{ij} \rightarrow \delta_i \frac{w_{ij}}{s_i} \quad (4)$$

It can be seen from the above equation that while node i adds a new edge, the edge between node i and its neighbor will increase. δ_i represents the rate of change; the edge weight between node i and node j increases Δw_{ij} , which shows that the neighbors of node i allocate additional traffic according to the edge weight. Therefore, the weight of node i finally adjusts to:

The δ_i said the rate of change, the weight of node i and its neighbor nodes and the edges of j increased Delta w_{ij} ; Delta w_{ij} said i node and its neighbors through the weights of edges to allocate additional traffic. Therefore, the weight of the node i is finally adjusted to:

$$s_i \rightarrow s_i + \delta_i + w_0 \quad (5)$$

Compared with the scale-free model, the BBV model considers the edge weights, which makes the network strength distribution obey the power-law distribution; by adjusting the δ_i values, it can change the same feature measurement, which makes great progress compared with scale-free network.

2) Local World network model

In the real world, many people have a specific circle, and they only live in this circle, which is the origin of the local world network; the local world is only a part of the entire network. The local world model [19], [20] is used to describe the situation.

The evolution rules of the local world network model are as follows:

a) The initial network: a network with m_0 nodes and e_0 edges;

b) The growth of the network: each time step a new node N joins into a network, node N with $m(m < m_0)$ edges. $M(m < M)$ nodes are selected randomly from the network as the local world of the new node N , and the new node N selects the m nodes from the local world network according to the node priority probability formula.

$$\Pi_{local} = \frac{M}{m_0 + t} \cdot \frac{k_i}{\sum_{j \in local} k_j} \quad (6)$$

The local world network model is suitable for some specific networks, when the network size is large enough, the cluster coefficient is close to 0.

B. Basic Concepts of the Model

Based on the research and analysis of the Enron e-mail network in the third chapter, the paper proposes an e-mail

network model based on the user information characteristics. Prior to this, we first give some definitions of basic concepts.

Definition 1: Feature vector of node: The feature attributes of nodes are represented by the tendency of nodes to fall under different feature classes. The information feature vector of the node is represented mathematically as:

$$f_i = \{F_1, F_2, F_3, F_4 \dots F_k\} \quad (7)$$

f_i represent the feature vector of node i , F_k represents the weight of the first K feature class of node i .

Definition 2: The similarity between nodes: The degree of similarity between the nodes, the higher similarity expresses the nodes are more likely to belong to the same class, and they have a high possibility of connecting with each other. The similarity between nodes is expressed mathematically as:

$$Similarity(i, j) = \frac{F_{i1} \cdot F_{j1} + F_{i2} \cdot F_{j2} + \dots + F_{ik} \cdot F_{jk}}{\sqrt{F_{i1}^2 + F_{i2}^2 + \dots + F_{ik}^2} \cdot \sqrt{F_{j1}^2 + F_{j2}^2 + \dots + F_{jk}^2}} \quad (8)$$

In this paper, we take the cosine similarity of vector space, which does not take the distance between two vectors into account.

Definition 3: Feature similarity network: A collection of M nodes with higher degree of features similarity from the original network after the new node is added. Mathematical representation:

$$V_{Feature}(i) = \{v_1, v_1, v_1 \dots v_M\} \quad (9)$$

Definition 4: Node strength: In a directed weighted network, the in strength of node i is the sum of edge weights, while node i is the in node; the out strength of node i is the sum of edge weights, while the node i is the out node.

The in strength of node i is:

$$s(in)_i = \sum_{j \in Neighbor(i)} w_{ji} \quad (10)$$

The out strength of node i is:

$$s(out)_i = \sum_{j \in Neighbor(i)} w_{ij} \quad (11)$$

The strength of node i is:

$$s_i = s(in)_i + s(out)_i \quad (12)$$

C. Model Evolution Rule

In the paper, according to the email transmission, sending, forwarding and replying, we consider the characteristics of the user in the message communication process; when a new node is added, selecting a certain number of nodes from the original network to form a feature similarity network, the new node selects a node from the feature network to connect; at the same time, there is internal evolution in the original network.

The specific construction algorithm of the e-mail network model based on the characteristics of user information is as follows:

1) Initial Network

The initial network contains m_0 nodes, and each node initializes a feature vector. For the sake of simplicity, this paper defines the m_0 class in the feature vector, and the initial value of each feature class in the feature vector of m_0 nodes is:

$$F_i(k) = \begin{cases} 0, & i \neq k; \\ 1, & i = k; \end{cases} \quad (13)$$

$F_i(k)$ represent the value of the first k class of node i ; the m_0 nodes form a fully coupled network, and the initial edge weight is w_0 .

2) The Growth of the Network

in each time step, a new node n joins the network, and the node n is randomly assigned a feature vector; at the same time, $M(M < m_0)$ nodes is selected from the original network according to the similarity of the feature vector to form the feature similarity network and with a certain probability to proceed as follows:

a) The new node with M ($m < m_0$) edge joins the feature similarity network with probability $p1$; some edges are out edges with probability q , and others are edges with probability q . The probability of node i as in node is:

$$\Pi = \frac{s(in)_i}{\sum_j s(in)_j} \quad (14)$$

The probability of node i as the out node is:

$$\Pi = \frac{s(out)_i}{\sum_j s(out)_j} \quad (15)$$

Among them, j is a set node that forms the feature similarity network.

b) The evolution in a feature similar network: to add m edges into a feature similarity network with probability $p2$ to achieve internal growth; in the feature similarity network, the new edge is $\langle i,j \rangle$; if there is a connection between node i and node j , we increase their weight; or, we establish a new edge and assign the initial weights w_0 . The probability of choosing node i is (15); the probability of choosing node j is (14).

c) Connection between the feature similarity network and the external network: to add m edges with probability $p3$ between the feature similarity network and the external network. In terms of the feature similarity network, m edges is as out edges with probability q , and m edges is as in edges with probability $1-q$. The choice of node is according to the operation 2).

3) The dynamic Evolution of the Weight

The generation of the new edge will trigger the readjustment of the weight between the node and the neighbor node. If the new edge is the in edge, the weight associated with the node i is changed to:

$$w_{ji} = w_{ji} + \Delta w_{ji} \quad (16)$$

$$\Delta w_{ji} = \delta_i \frac{w_{ji}}{s(in)_i} \quad (17)$$

Parameter δ_i is the additional traffic burden while new edge $\langle n,i \rangle$ is added; the neighbor nodes of node i share this traffic. So, the strength of node i is adjusted to:

$$s(in)_i = s(in)_i + w_0 + \delta_i \quad (18)$$

If the new edge is the out edge, the weight of the node i changes to:

$$w_{ij} = w_{ij} + \Delta w_{ij} \quad (19)$$

$$\Delta w_{ij} = \delta_i \frac{w_{ij}}{s(out)_i} \quad (20)$$

The out strength of node i is adjusted to:

$$s(out)_i = s(out)_i + w_0 + \delta_i \quad (21)$$

The growth of the network in steps 2) and 3) belongs to the evolution of the inner network, and the model does not consider additional traffic burden, so the corresponding weight of the edge and node strength are increased by w_0 .

4) Node Feature Vector Adjustment

While a new edge $\langle i,j \rangle$ is added, the node i delivers information to the node j , and it only causes the feature vector adjustment of node j ; that is, each node that changes is an in node. The paper uses vector space cosine similarity, and it only considers the direction gap but not the distance gap, so the feature class in the feature vector of node j is adjusted to:

$$F_{jk} = F_{jk} + F_{ik} \quad (22)$$

After t time steps, the network contains $m_0 + t$ nodes.

IV. EXPERIMENT AND ANALYSIS

According to the evolution model based on user information characteristics, we use Java programming to achieve the evolution of the network; then we get a network topology matrix; finally, we use MATLAB to calculate the parameters distribution. In the following, we analyze the model from two aspects, that is, average path length and cluster coefficient.

A. Average Path Analysis

Network path and diameter are important parameters of network transmission delay, and network transmission delay is an important factor of network performance and information dissemination. In order to represent the performance of the whole network, the concept of average path is introduced. First, the shortest path for each node to other nodes is obtained, each node is only allowed access once. After finding the shortest path between all the nodes, the average path length of the current network can be calculated.

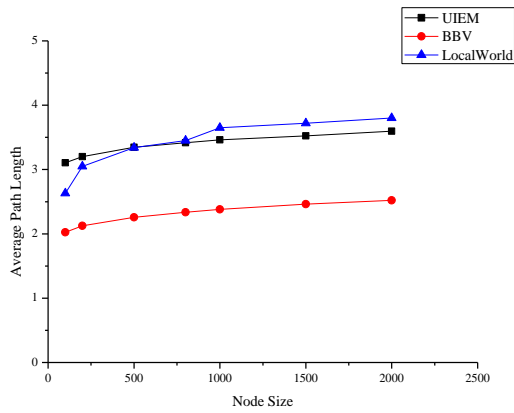


Fig. 7. Average path comparison of different models when the network size is 2000.

Fig. 7 shows the changes in network average path while the network size ranges from 1 to 2000. We can see that the average path of UIEM is similar to the local world network. Compared to the BBV network model, the average path of UIEM is larger, but the growth of the average path length is slower than that of the Local World network.

B. Network Cluster Coefficient

Cluster coefficient is the relationship between the node and its neighbors; in general, it is used to express the possibility that people’s friends are also friends. Because our network is a weighted network, we need to calculate the weighted clustering coefficient of the network [21]-[23]. The cluster coefficient represents the clustering degree of nodes in the network and is an important feature of a network [24]-[26]. A large number of studies have shown that real networks have high clustering characteristics.

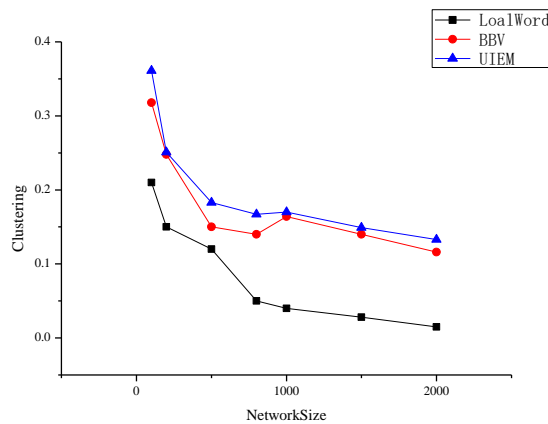


Fig. 8. Comparison of cluster coefficients of different models when the network size is 2000.

As shown in Fig. 8, the abscissa represents the size of the network, and the ordinate represents the cluster coefficient of the network. It can be seen that the cluster coefficient of UIEM is similar to the BBV model, but it is higher than the Local World network. The cluster coefficient is the degree of the network group. The cluster coefficient of the Local World

network tends to 0 when the network size is large enough, and it is clearly inconsistent with the actual network. The cluster coefficient of UIEM declines slowly with the growth of the network; this is consistent with the different actual networks to better reflect the authenticity of the network

C. Contrast Experiment

In this paper, in addition to using the Enron dataset, we also downloaded a mail data uploaded by the Department of Automation, Shanghai University to analyze the two data sets and perform a comparison with UIEM. For practical reasons, the Enron dataset contains 150 nodes, and the mail dataset from the Department of Automation, Shanghai University contains 1133 nodes, so this paper uses the corresponding number of nodes to compare with UIEM to ensure fairness.

TABLE VI. COMPARISON OF ENRON NETWORK WITH THREE MODELS

	Average Path	Cluster Coefficient
Enron	6.3	0.433
BBV	2.053	0.302
LocalWorld	2.82	0.18
UIEM	3.125	0.36

As shown in Table VI, when the scale is small, the average path of UIEM is longer, which is closer to the real Enron network; and the cluster coefficient of UIEM is higher; it is close to the actual network model.

The following is a comparison between the e-mail dataset and the three network models, which are shared by the Department of automation, Shanghai University.

TABLE VII. COMPARISON OF E-MAIL NETWORK FROM DEPARTMENT OF AUTOMATION, SHANGHAI UNIVERSITY WITH THE THREE MODELS

	Average Path	Cluster Coefficient
Email From Department of Automation, Shanghai University	3.606	0.22
BBV	2.381	0.155
Localworld	3.802	0.042
UIEM	3.482	0.168

From Table VII, we can see that when the network size is 1133, the average path length of the Local World network is growing too fast, more than the average path of real email network, but our UIEM is closer to a real email network than the other two models.

V. CONCLUSION

In this paper, we take the idea of using the Local World network and the dynamic evolution of the BBV model; then, according to the relationship between user information characteristics and communication that is found in chapter three, we present an e-mail network evolution model based on the characteristics of user information and give the construction rules and related definitions. Finally, realizing the

evolution of the network by programming, we find that the strength and degree of nodes are in accordance with the power law distribution. And compared with the BBV model and the Local World, UIEM is closer to the actual network and has practical significance.

ACKNOWLEDGMENT

The authors would like to thank the reviewers for their detailed reviews and constructive comments, which have helped improve the quality of this paper. This work is sponsored by the Natural Science Foundation of Heilongjiang Province of China under Grant No. LC2016024, the Natural Science Foundation of the Jiangsu Higher Education Institutions Grant No. 17KJB520044.

REFERENCES

- [1] D.J. Watts, S.H. Strogatz. "Collective dynamics of 'small-world' networks". *Nature*, vol. 393, pp. 440-442, 1998
- [2] M.E.J. Newman, D.J.Watts. "Renormalization group analysis of the small-world network model". *Physics Letters A*, vol. 263, pp. 341-346, 1999.
- [3] Barabasi A L, Albert R. "Emergence of scaling in random networks" *Science*, vol. 286, pp. 509. 1999.
- [4] R. Albert, A.L. Barabasi. "Topology of evolving networks: local events and universality". *Physical Review*, vol. 85, pp. 5234, 2000.
- [5] S. Fortunato, A. Flammini, F. Menczer. "Scale-free network growth by ranking". *Physical Review*, vol. 96, 2006
- [6] X.L. Sun, H.F. Lin, K Xu. "A social network model driven by events and interests". *Expert Systems With Applications*, vol. 42, pp. 4229-4238, 2015.
- [7] A. Abbasi, L. Hossain, L. Leydesdorff. "Betweenness centrality as a driver of preferential attachment in the evolution of research collaboration networks". *Journal of Informetrics*, vol. 6, pp. 403-412, 2012.
- [8] A. Barrat, M. Barthélemy, A. Vespignani. "Weighted evolving networks: coupling topology and weight dynamics". *Physical Review*, vol. 92, 2004
- [9] A. Barrat, M. Barthélemy, R. Pastor-Satorras, A. Vespignani. "The architecture of complex weighted networks". *Proceedings of the National Academy of Sciences of the United States of America* .vol. 101, pp. 3747-3452, 2004
- [10] X. Sun, J.Y. Dong, R.C. Tang, M.T. Xu, L. Qi, Y. Cai. "Topological evolution of virtual social networks by modeling social activities". *Physica A: Statistical Mechanics and its Applicat*, vol. 433, pp. 259-267, 2015.
- [11] B. Yao, M. Yao, X.E. Chen, X. Lin, W.J. Zhang. "Applied Mechanics and Materials. Research on Edge-Growing Models Related with Scale-Free Small-World Networks". *Applied Mechanics & Materials*, vol. 513, pp. 2444-2448, 2014.
- [12] B. Yao, C. Yang, M. Yao, et al. "Graphs as Models of Scale-Free Networks". *Applied Mechanics & Materials*, vol. 380, pp. 2034-2037, 2013.
- [13] W.Q. Wang, Q.M. Zhang, T. Zhou. "Evaluating network models: A likelihood analysis" *EPL (Europhysics Letters)*, vol. 98, pp. 28004-28009, 2012.
- [14] Z.Y. Zou, P. Liu, L. Lei, J.Z. Gao. "An evolving network model with modular growth". *Chinese Physics B*, vol. 21, pp. 603-609, 2012.
- [15] B.K. Wang, Z.H. Pei, L. Wang. "Evolutionary dynamics of cooperation on interdependent networks with the Prisoner's Dilemma and Snowdrift Game". *EPL (Europhysics Letters)*, vol. 107, pp. 58006, 2014.
- [16] H. Zhuang, Sun Y, Tang J, et al. "Influence maximization in dynamic social networks". *2013 IEEE 13th International Conference on Data Mining (ICDM)* .pp. 1313-1318, 2013.
- [17] N. Ilhan, I.G. Oğuducu. "Community Event Prediction in Dynamic Social Networks". *2013 12th International Conference on Machine Learning and Applications (ICMLA)* , pp. 191-196. 2013.
- [18] A. Barrat, M. Barthélemy, A. Vespignani. "Modeling the Evolution of Weighted Networks". *Physical Review E*, vol. 70, pp. 1-13, 2004.
- [19] X. Li, G. Chen. "A local world evolving network model". *Physical A*. vol. 328, pp. 274-286, 2003.
- [20] Z.f. Pan, X. Li, X.F. Wang. "Generalized local-world models for weighted networks". *Physical review. E, Statistical, nonlinear, and soft matter physics*. Vol. 73. 2006
- [21] X. Xue, S. Wang, B. Gui, et al. "A computational experiment-based evaluation method for context-aware services in complicated environment". *Information Sciences*, vol. 373, pp. 269-286, 2016.
- [22] X. Xue, Y.M. Kou, S. Wang, et al. "Computational experiment research on the equalization-oriented service strategy in collaborative manufacturing". *IEEE Transactions on Services Computing*, vol. 11, pp. 369-383, 2018.
- [23] X. Xue, H. Han, S. Wang, et al. "Computational Experiment-based Evaluation on Context-aware O2O Service Recommendation". *IEEE Transactions on Services Computing*, 2016
- [24] T. Wang, .Y. Wu, X. He, et al. "A Cross Unequal Clustering Routing Algorithm for Sensor Network". *Measurement Science Review*, vol. 13, pp. 200-205, 2013.
- [25] T. Wang, Y. Cao, Y. Zhou, et al. "A Survey on Geographic Routing Protocols in Delay/Disruption Tolerant Networks (DTNs)". *International Journal of Distributed Sensor Networks*, vol.6 , 2016.
- [26] Y. Cao, T. Wang, O. Kaiwartya, et al. "An EV Charging Management System Concerning Drivers' Trip Duration and Mobility Uncertainty". *IEEE Transactions on Systems Man & Cybernetics Systems*, vol. 48, pp. 596-607, 2018.

Design of Traffic Flow Simulation System to Minimize Intersection Waiting Time

Jang, Seung-Ju

Department of Computer Engineering,
Donggeui University

Abstract—This paper designs a traffic simulation system for minimizing intersection waiting time. We use SUMO simulator which is widely used as a traffic flow simulation tool for traffic flow simulation. Using the SUMO simulator to set the route from the source to the destination and measuring the time required when using the existing intersection signal system. Through this simulation, we want to measure how much the proposed system can minimize the waiting time. In order to minimize the intersection waiting time, it is assumed that there is a loop sensor that can recognize whether there is a waiting vehicle in each direction of the intersection. Using this information, a signal lamp is used as a waiting signal in the case of a direction in which there is no waiting vehicle, and a driving signal is given in the case of a waiting vehicle or an entering vehicle. In this paper, we try to reduce the time required for vehicles to arrive at their destination by making the traffic flow smoothly without any expense such as road expansion through the limited system.

Keywords—Traffic flow simulation; SUMO simulator; reduce traffic time; intersection traffic flow; simulation design

I. INTRODUCTION

Advances in vehicle technology have provided people with convenient and safe transport. However, the rapid increase in the number of vehicles has intensified traffic congestion, and physical solutions such as road extension are no longer a good solution. It is no longer possible to physically construct roads and extend buildings.

In recent years, Intelligent Transportation Systems (ITS) have been studied in an effort to solve these problems by using existing facilities more efficiently through advanced IT technology. ITS is a convergence of IT technology and transportation. It is a next-generation transportation system that integrates intelligent advanced technologies such as electronics, control, and communication with components of existing transportation systems such as roads, vehicles, and signal systems.

In advanced foreign metropolises, a traffic control system that manages only urban highways and safety management measures are introduced and operated separately from general highways. Recently, the importance of intelligent transportation system as a strategic target facility has been increasing. The actual situation of the traffic congestion including the expressway and the main road is not a mutually independent system but an organically integrated system such as system operation and influence due to individual control strategy.

Research on the integrated control model, which is an approach to this system worldwide, is actively being conducted. In Korea, there are no cases that have been studied from this point of view. Due to the development of Intelligent Transportation System (ITS), it is possible to collect data, and dynamic and intelligent traffic signal control becomes possible. Despite these technological advances, the problem of traffic congestion in urban areas is still not resolved.

The urban area consists of a network of multiple intersections with a very high traffic volume. Therefore, if a part of the traffic network is congested, it can affect not only the traffic flow of the following road but also the traffic flow of the other intersection. ITS is a system that maximizes the efficiency of transportation facilities and provides transportation convenience and safety, and the infrastructure is being established under the leadership of the government and local governments. The ITS system is converged or integrated with Geographic Information System (GIS), Global Positioning System (GPS), LBS, and telematics element technologies to provide traffic information to users.

II. RELATED RESEARCH

A related study for minimizing the intersection waiting time has been to efficiently schedule the green signal to reduce the average waiting time and the total travel time at the intersection, assuming that the intersection signal has the final destination information of the vehicle. Algorithms for controlling the traffic lights to provide services to the vehicles with the shortest route remaining from the intersection to the destination are being studied [1, 2].

Another signal control method to improve traditional signaling using fixed-time scheduling is being studied to analyze the pattern of vehicle flow through an intersection during the day. A study on the algorithm that adjusts the signal pattern for each signal cycle by controlling the vehicle flow at the intersection according to the predicted vehicle pattern or by using the statistical value of the traveling direction of the vehicle leaving the intersection during the previous signal cycle [3, 4, 5].

However, research to reduce the waiting time of intersections is not very active. Most of them are operated in a simple form in which the signal pattern is firstly determined in consideration of the traffic conditions of the surrounding roads [6, 7, 8].

In recent years, research has been shifting from fixed signal control to active control to reduce waiting time at

intersections in urban areas. The active control method collects the flow of the vehicle in real time and performs traffic control based on this information. Recent studies have been conducted to control signals using reinforcement learning algorithms. It is possible to see the waiting queue length of the vehicle waiting at the intersection and set the signal flexibly [9].

Reinforcement learning is a method in which a defined agent recognizes the current state and selects a behavior or sequence of actions that maximizes compensation among selectable behaviors. A study on signal control using Q-learning algorithm, which is one of the reinforcement learning algorithms, is being conducted. The reinforcement learning algorithm is one of the research fields of machine learning. It is a learning theory that accumulates the feedback obtained from the surrounding environment through repetitive search and takes the optimal selection based on this feedback. The following is the operation process of the reinforcement learning algorithm (Fig. 1).

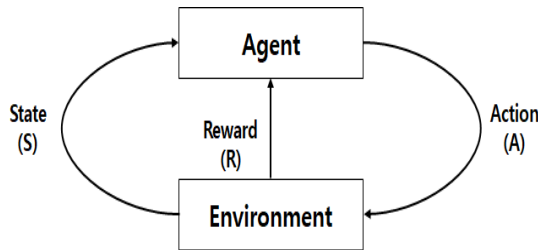


Fig. 1. Process of the reinforcement learning algorithm

III. TRAFFIC FLOW SIMULATION DESIGN

This paper designs a traffic flow simulation system to minimize the intersection waiting time on the road. Simulation of Urban Mobility (SUMO) simulator is used for traffic flow simulation system design.

The SUMO simulator tool is an inter- and multi-modal, space-continuous and time-discrete traffic flow simulation platform. The SUMO simulator tool was developed in 2002 for anyone to use in open source form. The SUMO simulator is a publicly available traffic simulator tool that follows the GPL policy. The SUMO simulator also supports the ability to use it in conjunction with existing simulators.

SUMO is a traffic simulator dealing with a wide range of road networks based on open source developed since 2000 at ITS of German Aerospace Center.

The main features of SUMO are as follows.

- Free collision avoidance of vehicle nodes
- Various vehicle characteristics information applicable
- Multiple lane and lane change function
- Interacting with other applications
- Application of intersection characteristics such as actual road environment

SUMO is capable of handling traffic network node information of 10,000 large-scale environments, and has the

advantage of generating node topology using files of various formats such as Visum, Vissim, ArcView, and XML. The following Fig. 2 shows the GUI interface of the SUMO traffic simulator.

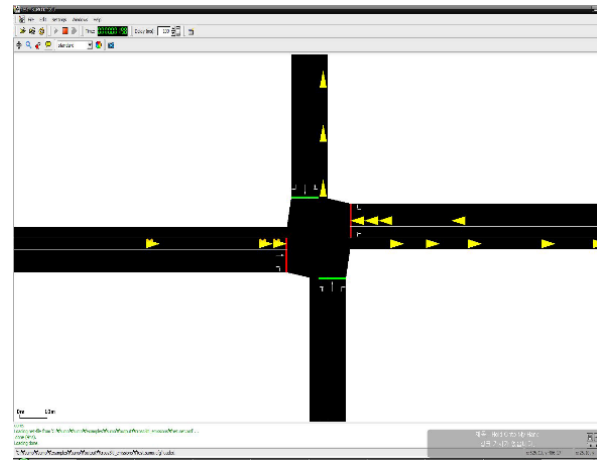


Fig. 2. SUMO Traffic simulator GUI

In this paper, we simulate traffic flow of intersection system using SUMO simulator. First, it is necessary to identify the problems with the existing traffic flow system. Most intersection signaling systems provide a certain amount of waiting time at the intersection, and when this time passes, they can pass through the intersection with green signals. These systems have advantages, but they also have several disadvantages. The advantage of a system that controls the flow of vehicles at regular time intervals can be the most optimal method for busy intersections. However, in the case of a no busy intersection, there is a fatal disadvantage of waiting for a certain period of time, even though the vehicle is not in the other lane. This causes a problem of delaying the running time of the vehicle at an intersection which is not much troublesome.

This paper simulates actual road traffic situation by using SUMO simulation tool to find improvement direction of existing intersection traffic signal system. In order to practice the road traffic situation using the SUMO simulation tool, we set the starting point and the end point of the actual road. In this paper, the actual starting point for the simulation of the traffic situation is the Busan Metropolitan City Dong Eui University. The terminal point is set at the entrance to the Hwangryung Tunnel of Busan Metropolitan City. To establish the actual road configuration for these two points, we construct road information linked with eWorld.

3.1 Link with SUMO Simulator and eWorld

In this paper, the actual starting point for the simulation of the traffic situation is the Busan Metropolitan City Dong Eui University. The terminal point is set at the entrance to the Hwangryung Tunnel of Busan Metropolitan City. To establish the actual road configuration for these two points, we construct road information linked with eWorld.

The following Fig. 3 shows the result of linking with the SUMO simulator tool on the starting point location and the ending point location using eWorld.

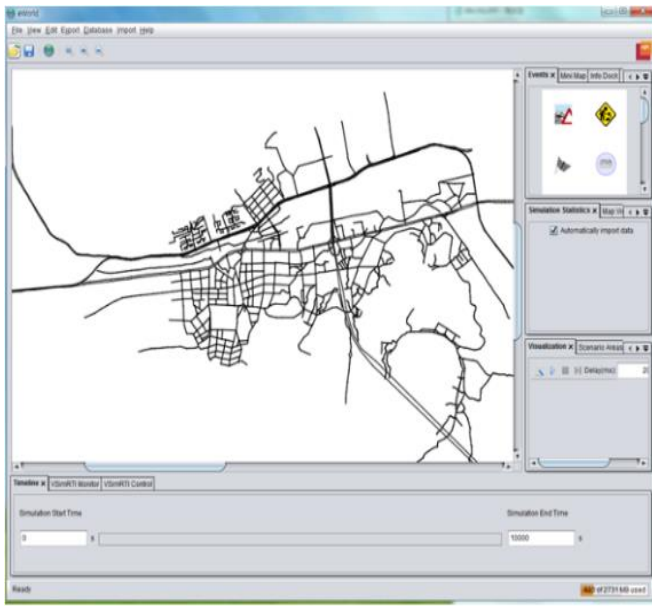


Fig. 3. Road map created using eWorld

Fig. 3 shows the path from the start to the end of the road to be simulated. In this way, an environment similar to the actual road situation is constructed to simulate the traffic flow.

3.2 Simulation Environment Design using SUMO Simulator

We design environment to measure intersection waiting time in actual road environment using SUMO simulator. We use the SUMO simulator to design an eWorld map for the same environment as the actual road situation. To simulate actual road conditions using the SUMO simulator, set it from Busan Metropolitan City Dong Eui University to the entrance of the Hwangryung Tunnel. After designing this environment, the designed files are as follows. : deuu123.add.xml, deuu123.edg.xml, deuu123.evt.xml, deuu123.flows.xml, deuu123.net.xml, deuu123.nod.xml, deuu123.poi.xml, deuu123.sumo.cfg

The deuu123.rou.xml file is described as follows. The deuu123.rou.xml file contains the following declarations for simulating the source to destination. Fig. 4 shows the declaration of these attributes.

```

<?xml version="1.0" encoding="TUF-8"?>
<routes>
<vType accel="3.0" decel="6.0" id="CarA" length="5.0" minGap="4.5"
maxSpeed="50.0" sigma="0.5"/>
<vType accel="2.0" decel="6.0" id="CarB" length="7.5" minGap="4.5"
maxSpeed="50.0" sigma="0.5"/>
<vType accel="1.0" decel="5.0" id="CarC" length="5.0" minGap="4.5"
maxSpeed="40.0" sigma="0.5"/>
<vType accel="1.0" decel="5.0" id="CarD" length="7.5" minGap="4.5"
maxSpeed="30.0" sigma="0.5"/>
    
```

Fig. 4. Declaration of the simulation attributes

Table I shows the contents of the attribute definition for Fig. 4

TABLE I. PROPERTIES AND FUNCTIONS FOR SUMO SIMULATION

Attribute Name	Value Type	Default	Description
id	id(string)	-	The name of the vehicle type
accel	float	2.6	The acceleration ability of vehicles of this type(in s/m ²)
decel	float	4.5	The acceleration ability of vehicles of this type(in s/m ²)
sigma	float	0.5	Car-following model parameter
tau	float	1.0	Car-following model parameter
length	float	5.0	The vehicle's netto-length(length)(in m)
minGap	float	2.5	Empty space after leader[m]
maxSpeed	float	70.0	The vehicle's maximum velocity(in m/s)

The deuu123.rou.xml source code above is designed to define automotive properties for simulation. You can specify the name of the car and the length of the car, as well as the distance between the front and back of the car and the maximum speed. The SUMO simulator supports the ability to simulate through the specification of these attributes. You can specify a route here. This function specifies the path the car should go during the simulation. The route id specifies the name of the route through which the car to be simulated passes.

Next is the information configuration for edges. This is the edge where the car travels from one junction to another. If the junction is not correctly set between one edge and the other, an error will occur. Detailed configuration of edges is shown in the following file deuu123.edg.xml.

```

<edge id="0034560-0447-4677-9b0c-e630d0a077" from="46e28113-b0f9-404e-e035-370a004694F" to="31c3d10e-8f19-40e0-91c1-d0d0de3226f" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="1618101-0402-470c-80da-802c4401c26f" from="50c3d10e-8f19-40e0-91c1-d0d0de3226f" to="e035d0e6-e044-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="1894027e-0376-451c-145-01195110e4eF" from="e035d0e6-e044-452c-805f-2d46464646F" to="c1f1e020-950a-4308-915f-40147474664F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="7042267f-702c-4407-027f-030e0750e40a" from="c1f1e020-950a-4308-915f-40147474664F" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="4309036-0704-420a-0119-4670c00c22c" from="0a0e027e-2208-452c-805f-2d46464646F" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="23e0e0e-0100-401a-0119-4670c00c22c" from="0a0e027e-2208-452c-805f-2d46464646F" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="195468f-0114-401c-0119-4670c00c22c" from="55334e68-014a-401c-0119-4670c00c22c" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="195468f-0114-401c-0119-4670c00c22c" from="55334e68-014a-401c-0119-4670c00c22c" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="195468f-0114-401c-0119-4670c00c22c" from="55334e68-014a-401c-0119-4670c00c22c" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="720e1001-4f70-444d-805c-071280e0e55f" from="57324240-504e-420a-0119-4670c00c22c" to="55334e68-014a-401c-0119-4670c00c22c" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="7f00040-01f5-444d-805c-071280e0e55f" from="55334e68-014a-401c-0119-4670c00c22c" to="55334e68-014a-401c-0119-4670c00c22c" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="7b40310c-0177-420c-0154-10f0e030e070f" from="01544054-e044-410b-0a0e-070b0704e003" to="57324240-504e-420a-0119-4670c00c22c" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="1693004-0814-401a-0154-10f0e030e070f" from="57324240-504e-420a-0119-4670c00c22c" to="57324240-504e-420a-0119-4670c00c22c" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="5510070-066f-404e-0109-21f10726000a" from="4c78021f-447b-401b-070b-030c035e000a" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="0353714-004a-440a-0109-21f10726000a" from="4c78021f-447b-401b-070b-030c035e000a" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="0353714-004a-440a-0109-21f10726000a" from="4c78021f-447b-401b-070b-030c035e000a" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="0a140010-0800-440a-0109-21f10726000a" from="0a0e027e-2208-452c-805f-2d46464646F" to="0a0e027e-2208-452c-805f-2d46464646F" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="7010000-0802-420c-01f1-0342007080f" from="e035d0e6-e044-452c-805f-2d46464646F" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="1601000-0141-420c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
<edge id="70e1404-030f-4f0c-01f1-0342007080f" from="20f43107-447b-401b-070b-030c035e000a" to="20f43107-447b-401b-070b-030c035e000a" priority="4" numLanes="1" speed="22.2222222222" />
    
```

Fig. 5. deuu123.edg.xml Configuration File

In Fig. 5, the edge id is set, and from to is also set. This indicates that, from one junction to the next junction set "From" and "To".

In this paper, we are designing a system for minimizing the intersection waiting time. I explained the current research contents. In the future research direction, we try to implement similar to the actual environment with the designed contents.

IV. EXPERIMENT

In this paper, we design a simulation method to reduce intersection waiting time. Experiments were conducted on the proposed design contents. Experiments were conducted using the SUMO simulator, linking the eWorld map to specify the starting and destination locations to be similar to the actual road conditions.

The following figure is a map of Busan Metropolitan City with Dong Eui University as the starting point and the destination with the Hwangryung Tunnel. The following shows how to simulate the actual start and destination through SUMO settings.

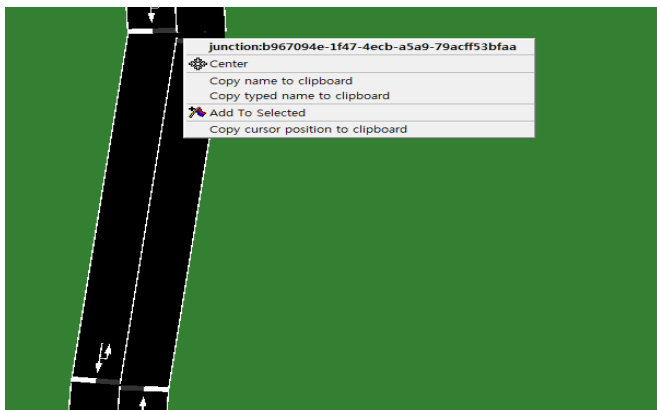


Fig. 6. SUMO Environment Setting

As shown in Fig. 6, simulation can be performed similar to the actual environment through SUMO setting. Fig. 7 shows the process of initial simulation starting place of Dong Eui University.

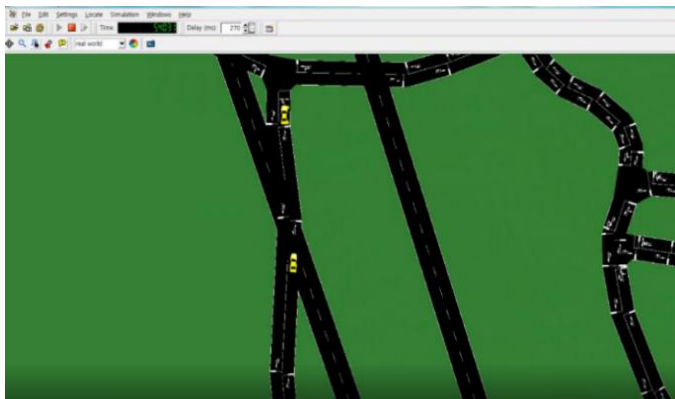


Fig. 7. Simulation initial process

Fig. 7 shows the initial process of simulation using SUMO. And the road situation is not complicated. Fig. 8 shows the process of the middle step to some extent using the SUMO simulator.

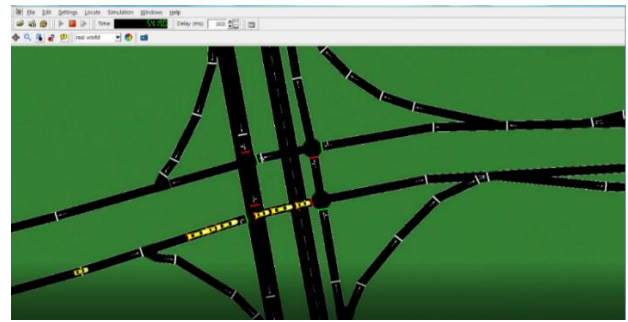


Fig. 8. Simulation middle process

In the case of Fig. 8, it can be seen that the vehicle is gradually increasing and proceeding to complicated road conditions. In this way, the simulation proceeds to the destination. As a result of this experiment, it can be confirmed that the design works normally. We will extend the simulation environment constructed in this paper to estimate the actual time required for future research.

V. CONCLUSIONS

This paper designs a traffic simulation system for minimizing intersection waiting time. We use SUMO simulator which is widely used as a traffic flow simulation tool for traffic flow simulation. Using the SUMO simulator to set the route from the source to the destination and measuring the time required when using the existing intersection signal system.

In this paper, we use SUMO simulator to simulate intersection waiting time. We try to find various solutions through simulation. We use the SUMO simulator to design in conjunction with actual road conditions. We work with eWorld maps to connect with real roads. In this paper, the actual starting point for the simulation of the traffic situation is the Busan Metropolitan City Dong Eui University. The terminal point is set at the entrance to the Hwangryung Tunnel of Busan Metropolitan City.

We design a system to simulate the time required for the route to the destination using the SUMO simulator on actual roads. To design such a system, the SUMO environment is set and constructed. Once the SUMO environment is established, we will experiment with the environment in which the cars run for the simulation. As a result of this experiment, it can be confirmed that the design works normally. We will extend the simulation environment constructed in this paper to estimate the actual time required for future research.

REFERENCES

- [1] Sung Yoon Young, Lee Sook Young, Lee Mee Jeong, "A Survey on traffic light control algorithm for vehicular traffic flow optimization based on real-time traffic information and contraflow lane", Proceedings of Symposium of the Korean Institute of communications and Information Sciences , pp. 618-619, Jan. 2017.
- [2] Y. K. Chin, N. Bolong, A. Kiring, S. S. Yang, K. T. K. Teo, "Q-Learning Based Traffic Optimization in Management of Signal Timing Plan," International Journal of Simulation Systems Science & Technology, vol.12, no.3, pp.29-35, 2011.
- [3] Y. K. Chin, W. Y. Kow, W. L. Khong, M. K. Tan, K. T. K. Teo, "Q-Learning Traffic Signal Optimization within Multiple Intersections

- Traffic Network," UKSim-AMSS 6th European Modelling Symposium, pp.343-348, Valetta, Malta, 2012.
- [4] Jinseop Cho, et al., "Multiple-Intersection Traffic Signal Control based on Traffic Pattern Learning", Journal of KIISE, JOK, No. 20, Vol. 3, Mar.2014.
- [5] YoungTae Jo, et al., "Intersection Traffic Signal Control based on Traffic Pattern Learning for Repetitive Traffic Congestion", Journal of KIISE, JOK, No. 20, Vol. 8, Aug.2014.
- [6] Kwang-back Kim, "Intelligent Traffic Light Control using Fuzzy Method", Journal of the Korea Institute of Information and Communication Engineering, No. 16, Vol. 8, May 2012.
- [7] Won-Kee Hong and Woo-Seok Shim, "Traffic Signal Control Scheme for Traffic Detection System based on Wireless Sensor Network", Journal of Institute of Control, Robotics and Systems (2012) 18(8):719-724
- [8] Jin-Tae Kim, "Diagnosis of Local Traffic Controller for Effective", JOURNAL OF THE KOREAN SOCIETY FOR RAILWAY VOL.18
- [9] Hun Choi, et al., "Business Model of U-Intelligent Traffic Information and Control Services in U-City Environment", The Journal of the Korea Contents Association, Vol. 10 No. 5, Apr.2010.

Application of the Hierarchy Analysis Method at the Foodstuffs Quality Evaluation

Nikitina Marina Aleksandrovna¹

Center of Economic and Analytical Research and
Information Technologies
V.M. Gorbатов Federal Research Center for Food Systems
of Russian Academy of Sciences, Russia

Nikitin Igor Alekseevich²,
Semenkina Natalya Gennadiyevna³

Department “Technology of Grain Processing, Bakery,
Pasta and Confectionery Industries”
K.G. Razumovsky Moscow State University of
Technologies and Management
(the First Cossack University), Russia

Zavalishin Igor Vladimirovich⁴

Department “Rectorate”
K.G. Razumovsky Moscow State University of
Technologies and Management
(the First Cossack University), Russia

Goncharov Andrey Vitalievich⁵

Department “Automation and Control in Technical
Systems”
K.G. Razumovsky Moscow State University of
Technologies and Management
(the First Cossack University), Russia

Abstract—In Russia as well as in the other countries of the world national programs are implemented to improve the health of the population. An integral part of those programs are measures of improvement of food processes structure as well as the quality of food itself. New types of functional and specialized food products that meet the physiological needs of specific groups of the population with a therapeutic and therapeutic-prophylactic action spectrum are becoming more widespread. The article proposes the concept of determining the quality of food products through the indicator of “effective functionality” on the basis of a multicriteria approach using the hierarchy analysis method. On the example of gluten-free flour confectionery products, the determination of the organoleptic evaluation of the quality of a food product is shown, as a particular solution for finding one of the complex indicators of the first level. The use of T. Saaty’s method in making technological decisions on a large number of criteria is substantiated. The analysis of the obtained data allows to draw a conclusion that the greatest weight among alternatives was possessed by the sample containing three kinds of flour: buckwheat, amaranth and linen in the ratio 60:30:10.

Keywords—Effective functionality; hierarchy analysis method; gluten-free flour confectionery products; organoleptic evaluation of the quality; food product quality

I. INTRODUCTION

Health of the population is the most important indicator of the well-being of the nation. The constant impact on the population of various environmental factors in combination with psychoemotional loads leads to a decrease in the adaptive capacity of the human body. Today the number of alimentary-dependent diseases continues to increase the leading position among which is occupied by diseases of the digestive system. The leading role in the prevention and treatment of these diseases belongs to metabolic therapy, which is based on diet therapy. The therapy considered now as one of the most

important adaptation-protective factors that promote the maintenance of good health, normal growth and development of the organism, preservation of working capacity and adaptation of the organism to adverse environmental factors [1].

The problem of nutrition correction is also relevant for Russia. The policy of the state is aimed at solving the problems connected with the organization of healthy nutrition of the population. “The fundamentals of the state policy in the area of healthy nutrition of the population of the Russian Federation for the period until 2020” define the increase of production of specialized products including flour confectionery as a priority task.

Nowadays the production of food products free of certain ingredients is rapidly developing, because these ingredients can be not recommended for certain medical indications (allergens, some types of proteins, oligosaccharides, polysaccharides, etc.). Taking into account the successes of nutrigenomics and nutrigenetics the trend towards the personalization of diets will increase and as a result contribute to an increase in the volume of the market for functional and specialized food products [2], [3]. The first key direction in the development of such products is a scientifically based selection of functional food ingredients with the required sanitary and hygienic, medical and biological indicators, therapeutic and prophylactic properties. And the second key direction is the development of new technological solutions that allow not only to influence the organoleptic and physicochemical parameters of raw materials and finished products increasing their nutritional value, but also to give them directed functional properties [4], [5].

An objective assessment of the increase in nutritional value and the imparting of functional properties to the finished product should be based on the principles of qualimetry.

For a more effective description of the evaluation characteristic and the possibility of comparing different functional and specialized products it is advisable to introduce a quantitative indicator of functionality that will allow to speak about the “effective functionality” of a food product and to determine it in a dimensionless quantity called the generalized complex efficiency index of the top (or zero) level K_0 .

The structure of the complex indicator is considered by the authors as a multilevel hierarchical set of properties, among which it is necessary to single out such basic indicators of the first level such as the chemical composition, organoleptic characteristics, physical and chemical properties, safety indicators and microbiological indicators as well as cost.

The scheme of this approach is shown in Fig. 1.

In addition to finding a complex zero-level indicator this scheme involves the definition of first-level indicators as the finding of particular solutions in assessing the quality of food. For example, safety indicators, microbiological indicators or organoleptic characteristics of a product can also be determined using this scheme and the hierarchy analysis method and implemented at a selection of the optimal formulation of the final product.

The purpose of this paper is to demonstrate the use of the hierarchy analysis method of T. Saaty to make decisions in the field of a limited study - research and management in assessing the quality of products of the food industry.

On the example of a specialized food product (gluten-free gingerbread, which contains non-traditional types of raw materials), a definition of the organoleptic evaluation of product quality is shown as a particular solution for finding one of the complex indicators of the first level.

II. OBJECTS AND METHODS OF RESEARCH

The research was carried out at the laboratory of the “Technology of processing grain, bakery, macaroni and confectionery productions” chair of the K.G. Razumovsky Moscow State University of technologies and management (the First Cossack University) in conjunction with the Information Technology Department of the Center for Economic and Analytical Research and Information Technologies of the “Federal State Research Institution of V.M. Gorbatov” of the Russian Academy of Sciences.

The object of the study was model samples of brewed gluten-free gingerbread with different ratios of non-traditional types of flour and protein concentrates.

For this purpose the following types of flour were chosen as the main raw material: amaranth, rice, buckwheat and corn, used in the control of celiac disease, a multifactorial disease that disrupts digestion caused by damage to villi in the small intestine by certain foods containing certain proteins: gluten (gluten) and close to it proteins of cereals (avenin, hordein, etc.) [6]-[8]; linseed flour, sesame, pumpkin seeds and milk thistle seeds served as protein concentrates.

The samples were sent to determine organoleptic quality indicators: taste, aroma, color, shape, appearance in the fracture.

The data was processed using the hierarchy analysis method of T. Saaty using the developed model for the effective evaluation of food quality indicators.

III. RESULTS AND DISCUSSION

To develop a methodology for predicting the quality of food products, authors analyzed methods that are used to solve similar problems in adjacent areas.

The quality of food products is always evaluated by some determining indicator. Since the degree of significance of the individual quality indicators is not the same, a weight coefficient is introduced [9].

Thus, the quality assessment is related to the task of quantitative evaluation by constructing its complex indicator. There is a dynamic, hierarchical, value and quantitative approaches [10].

As a result, the authors used the hierarchy analysis method to assess the organoleptic quality indicators of flour confectionery products.

The top of the hierarchy is the main goal. Elements of the lower level are many options for achieving the goal. Elements of intermediate levels meet the criteria or factors that connect the goal with the alternative. Having built a food system as a hierarchy, it is necessary to determine the priorities of all the nodes [11].

Priorities are the relative weights of the elements in each group. Like probabilities, priorities are dimensionless quantities that can take values from 0 to 1. The higher the priority value, the more significant is the element corresponding to it.

K_0 - complex upper-level indicator; K_1 - chemical composition; K_2 - organoleptic characteristics; K_3 - physical and chemical properties; K_4 - safety indicators; K_5 - microbiological indicators; K_6 - cost price; K_{11} - protein content; K_{12} - fat content; K_{13} - carbohydrate content; K_{14} - mineral content; K_{15} - vitamin content; K_{16} - energy value; K_{21} - taste; K_{22} - aroma; K_{23} - colour; K_{24} - form; K_{25} - appearance in the fracture; K_{31} - humidity; K_{32} - acidity; K_{33} - porosity; K_{34} - specific volume; K_{35} - deformation of crumb compression; K_{41} - pesticides; K_{42} - radionuclides; K_{43} - toxic elements; K_{44} - mycotoxins; K_{51} - content of the number of mesophilic aerobic and facultative anaerobic microorganisms (NMAFAnM); K_{52} - content of the *colibacillus* group bacteria (CGB); K_{53} - content of *S. aureus*; K_{54} - content of *Proterus* bacteria; K_{55} - content of pathogenic. incl. *Salmonella*; K_{56} - the content of mold; K_{61} - cost more than 50% higher than the average cost of analogue of this product is not a functional purpose; K_{62} - the cost price is not more than 50% higher than the average cost price for the analogue of this product is not a functional purpose; K_{63} - prime cost as in the analogue. K_{64} - cost of not more than 50% of the lower average cost of the analogue of this product is not functional; K_{65} - cost more than 50% of the lower cost of the analogue of this product is not functional; K_{111} - the content of

essential amino acids; K_{112} - the content of interchangeable amino acids; K_{121} - content of SFA; K_{122} - content of LSFA; K_{123} - content of PUFA; K_{131} - the content of digestible carbohydrates; K_{132} - the content of dietary fiber; K_{141} - the content of water-soluble vitamins; K_{142} - the content of fat-soluble vitamins; K_{151} - the content of macroelements; K_{152} - the content of trace elements (Fe-iron); K_{1111} - the amino acid content of lysine (Lys); K_{1112} - the amino acid content of methionine + cystine (Met + Cys); K_{1113} - the amino acid content of tryptophan (Trp); K_{1114} - the amino acid content of

isoleucine (Ile); K_{1115} - the amino acid valine (Val); K_{1116} - amino acid content of phenylalanine + tyrosine (Phe + Tyr); K_{1117} - amino acid content of threonine (Thr); K_{1118} - amino acid content of leucine (Leu); K_{1411} - vitamin B1 (thiamine) content; K_{1412} - vitamin B2 (riboflavin) content; K_{1413} - vitamin PP content (niacin, nicotinic acid); K_{1421} - vitamin E content; K_{1422} - β -carotene content; K_{1511} - the content of calcium (Ca); K_{1512} - the content of magnesium (Mg); K_{1513} - the content of phosphorus (P).

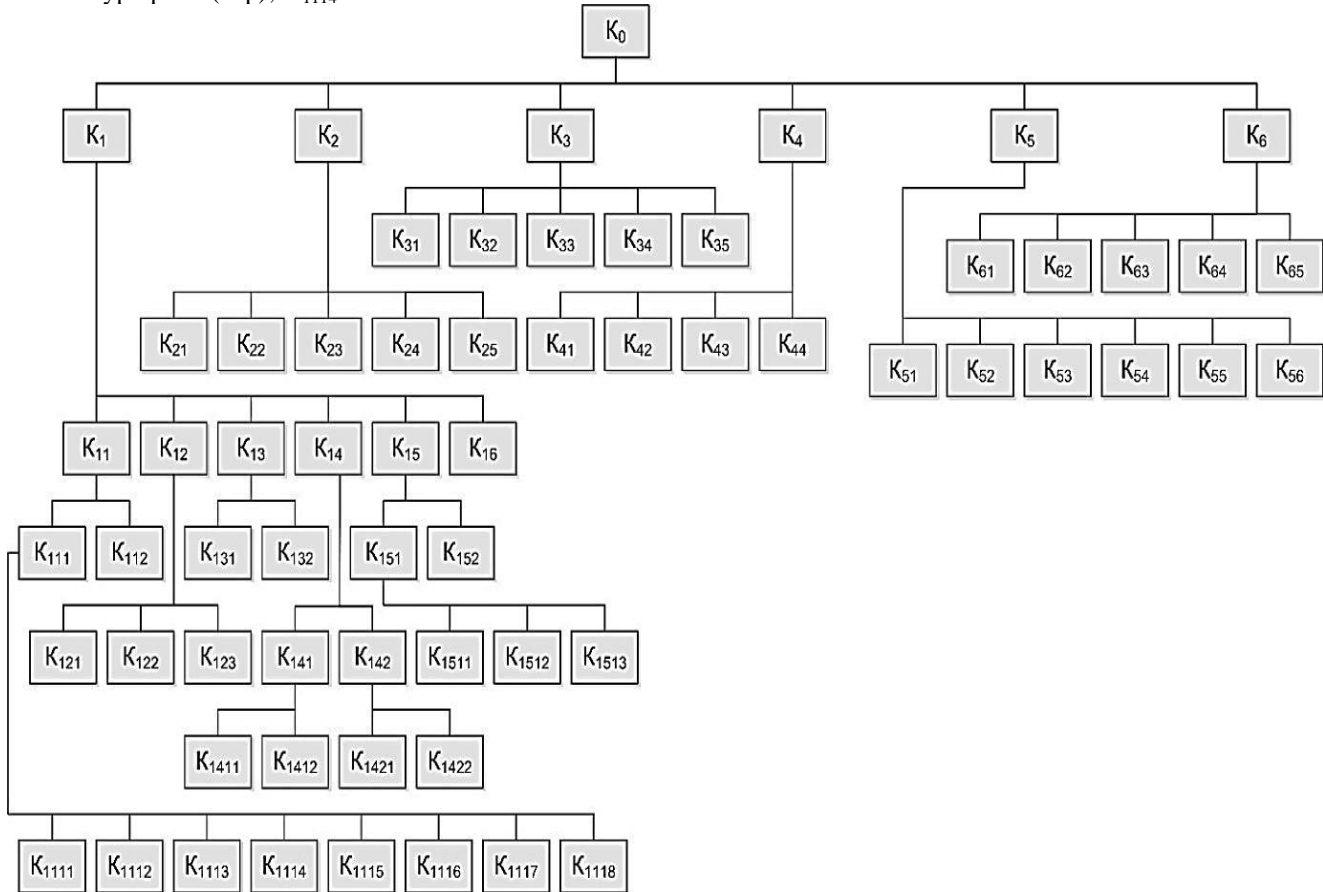


Fig. 1. Multi-level structure of a complex indicator of the quality of a food product (on the example of a bakery product).

IV. CONSTRUCTION OF A MATRIX OF RELATIONS BETWEEN ALTERNATIVE SOLUTIONS OF OBJECTS ON THE EXAMPLE OF GLUTEN-FREE GINGERBREAD

Authors of the article suggest using the hierarchy analysis method (or the Saaty method) [12] to study the weight of each individual parameter in assessing the quality of a food product. The founder of the decision-making process Analytic Hierarchy Process (AHP), known in Russia as a “hierarchy analysis method”, is the American scientist T. Saaty from the University of Pittsburg (www.pitt.edu) (www.business.pitt.edu/katz/faculty/saaty.php) [13]-[17].

The method developed by the American mathematician T. Saaty is a more justified means of solving multicriteria problems in a complex situation with hierarchical structures involving both tangible and intangible factors than approaches based on linear logic. As T. Saati said [18], the hierarchy analysis method is a closed logical construction that provides,

through simple rules, the analysis of complex problems in all their diversity and leading to the best answer. In addition, the application of the method makes it possible to include in the hierarchy all the knowledge and imagination available to the researcher on the problem under consideration. This, from the authors’ point of view, is a balanced way of solving a difficult problem: leaving the math simple and letting the structure’s diversity carry the burden of complexity.

AHP is based on paired comparisons of alternatives according to various criteria using a 5-point scale and the subsequent ranking of a set of alternatives according to all criteria and objectives. The relationship between the criteria is taken into account by constructing a hierarchy of criteria and applying the pairwise comparison method to identify the importance of criteria and subcriteria (Fig. 2).

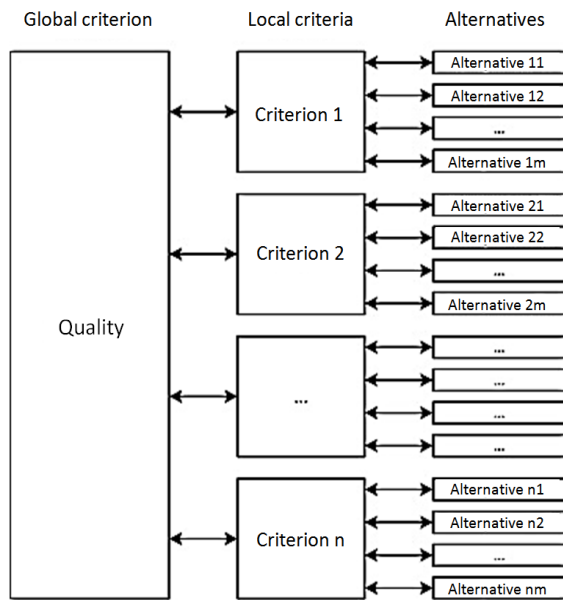


Fig. 2. The scheme of multi-criteria choice with the simplest hierarchy.

To implement the selection algorithm it is sufficient to have information about the type of relationship between each pair of objects and in particular about the existence of strict preference relations between two objects. To do this a relationship variable (1) is introduced

$$a_{ij} = \begin{cases} 1, & \text{if } i\text{-th variant is equal to } j\text{-th} \\ 3, & \text{if } i\text{-th variant moderately exceeds } j\text{-th} \\ 5, & \text{if } i\text{-th variant is more significant than } j\text{-th} \end{cases} \quad (1)$$

Values 2, 4 at the T. Saaty scale are intermediate values between adjacent values of the scale.

On the basis of the data obtained a square matrix $\|a\|$ (Table I) is constructed for the relationship between the alternatives of solutions $a_{ji} = \frac{1}{a_{ij}}, a_{ii} = 1, i, j = \overline{1, n}$.

TABLE I. THE MATRIX OF RELATIONSHIP BETWEEN ALTERNATIVE SOLUTIONS (ON THE EXAMPLE OF ORGANOLEPTIC INDICATORS OF GLUTEN-FREE GINGERBREAD)

a_{ij}	Taste	Aroma	Colour	Form	View of the fracture
Taste	1	4	2	5	3
Aroma	1/4	1	2	3	2
Colour	1/2	1/2	1	1/3	1/4
Form	1/5	1/3	3	1	1/3
View of the fracture	1/3	1/2	4	3	1
Main vector $X_j = \sum_{i=1}^n a_{ij}$	2.2833	6.3333	12.0000	12.3333	6.5833

V. DEVELOPMENT OF GLUTEN-FREE CUSTARD CAKE RECIPES BASED ON NON-TRADITIONAL RAW MATERIALS

For example, in the article the calculation of the weight coefficients of one of the criteria (organoleptic characteristic) in the assessment of the quality of model gluten-free custard cakes is given.

The custard cakes belong to the group of confectionery products and are one of the components of the diet of the population. However in diseases associated with hereditary genesis, which includes celiac disease (gluten enteropathy), not everyone can eat foods containing wheat flour [19]. Foods that do not contain gluten make one of the segments of the fast-growing market of specialized food products. The assortment of bakery and flour confectionery products for the gluten-free diet is constantly expanding. As gluten-free raw materials, starch-containing raw materials are most often used. It reduces nutritional value and gives the products worse organoleptic properties than traditional assortment [20]-[22].

Scientists of the “Technology of grain processing, bakery, macaroni and confectionery production” chair from the Razymovsky MSUTM (FCU) developed recipes for gluten-free custard cakes based on unconventional raw materials – amaranth, rice, buckwheat, corn, linseed, sesame, pumpkin seed flour and milk thistle seed flour. Each sample was assigned a serial number (Table II). The resulting gingerbread was analyzed for organoleptic characteristics (taste, aroma, color, shape and appearance in the fracture) using the hierarchy analysis method of the above mentioned algorithm.

Table I presents a weighted average of respondents’ preferences in the choice of flour confectionery products (gluten-free custard cakes). For the reliability of the results of the assessment the number of respondents was 7 people [23].

The vector of priorities was calculated from the matrix $\|a\|$. According to mathematical terms this is the main eigenvector, which after normalization becomes a vector of priorities. To calculate the analytical estimate of a given vector there are several ways. One of them is as follows. We

find the sum of the columns $X_j = \sum_{i=1}^n a_{ij}, j = \overline{1, n}$ of the matrix

$\|a\|$ in the form of a row vector {2.2833; 6.3333; 12.0000; 12.3333; 6.5833} and divide each column element by this sum. As a result we get a new matrix $\|a^*\|$ of values (Table III), which allows us to evaluate the significance of each individual indicator in the overall product perception characteristic.

Finding the average value of each i -line allows you to get the column vector of priorities {0.419; 0.196; 0.089; 0.104; 0.191}.

Thus, according to this expert ranking of the priorities between the indicators of organoleptic evaluation we obtain that the highest weight coefficient has “taste” - 41.9%, then the “aroma” - 19.6%, then the “fracture” in 19.1%, the “form” - 10.4%, the “color” - 8.9%.

TABLE II. COMPOSITION OF FLOUR COMPOSITE MIXTURES FOR GLUTEN-FREE GINGERBREAD, %

Sample Name	Sample number								
	1	2	3	4	5	6	7	8	9
Rice flour	60	60	-	60	-	60	-	-	70
Amaranth flour	20	-	20	-	-	30	20	30	20
Pumpkin seed flour	20	20	-	20	20	-	20	-	-
Buckwheat flour	-	20	60	-	60	-	-	60	-
Sesame flour	-	-	20	-	-	-	-	-	10
Corn flour	-	-	-	20	20	-	60	-	-
Schrot from the seeds of milk thistle	-	-	-	-	-	10	-	-	-
Flax, semi-fat flour	-	-	-	-	-	-	-	-	-

TABLE III. A NEW MATRIX OF VALUES

a_{ij}	Taste	Aroma	Colour	Form	View of the fracture	Vector of priorities $X_i = \frac{\sum_{j=1}^n a_{ij}}{n}$
Taste	0.438	0.632	0.167	0.405	0.456	0.419
Aroma	0.109	0.158	0.167	0.243	0.304	0.196
Colour	0.219	0.079	0.083	0.027	0.038	0.089
Form	0.088	0.053	0.250	0.081	0.051	0.104
View of the fracture	0.146	0.079	0.333	0.243	0.152	0.191

The obtained values are used to calculate the generalized complex quality indicator of gingerbread.

VI. INVESTIGATION OF ORGANOLEPTIC QUALITY INDICATORS OF GLUTEN-FREE GINGERBREAD USING THE HIERARCHY ANALYSIS METHOD

After conducting a general assessment of the perception of the food product it is necessary to perform calculations for each individual indicator of the investigated characteristic of the compared samples of gluten-free custard cakes. In our case there are nine samples.

First, respondents rated the indicators of organoleptic evaluation with each other (Table I), and then compared model samples of flour confectionery products against these five characteristics (Table IV).

Table IV presents the average weighted estimates of respondents in a pairwise comparison of nine samples of model flour confectionery products by organoleptic indicators – taste, aroma, color, shape and appearance of the fracture.

TABLE IV. MATRIX OF PAIRWISE COMPARISONS OF FLOUR CONFECTIONERY

a) Taste

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9
Sample 1	1	1	1/3	1/2	1/3	1/4	1/5	1/5	1/4
Sample 2	1	1	1/3	1/2	1/3	1/4	1/5	1/5	1/4
Sample 3	3	3	1	2	1	1/2	1/3	1/3	1/2
Sample 4	2	2	1/2	1	1/2	1/3	1/4	1/4	1/3
Sample 5	3	3	1	2	1	1/2	1/3	1/3	1/2
Sample 6	4	4	2	3	2	1	1/2	1/2	1
Sample 7	5	5	3	4	3	2	1	1	2
Sample 8	5	5	3	4	3	2	1	1	2
Sample 9	4	4	2	3	2	1	1/2	1/2	1
Sum	28.0000	28.0000	13.1667	20.0000	13.1667	7.8333	4.3167	4.3167	7.8333

b) Aroma

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9
Sample 1	1	1	1/3	1/2	1/3	1/4	1/5	1/5	1/4
Sample 2	1	1	1/3	1/2	1/3	1/4	1/5	1/5	1/4
Sample 3	3	3	1	2	1	1/2	1/3	1/3	1/2
Sample 4	2	2	1/2	1	1/2	1/3	1/4	1/4	1/3
Sample 5	3	3	1	2	1	1/2	1/3	1/3	1/2
Sample 6	4	4	2	3	2	1	1/2	1/2	1
Sample 7	5	5	3	4	3	2	1	1	2
Sample 8	5	5	3	4	3	2	1	1	2
Sample 9	4	4	2	3	2	1	1/2	1/2	1
Sum	28.0000	28.0000	13.1667	20.0000	13.1667	7.8333	4.3167	4.3167	7.8333

c) Colour

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9
Sample 1	1	1	1/2	1/2	1/2	1	1/2	1/2	1
Sample 2	1	1	1/2	1/2	1/2	1	1/2	1/2	1
Sample 3	2	2	1	1	1	2	1	1	2
Sample 4	2	2	1	1	1	2	1	1	2
Sample 5	2	2	1	1	1	2	1	1	2
Sample 6	1	1	1/2	1/2	1/2	1	1/2	1/2	1
Sample 7	2	2	1	1	1	2	1	1	2
Sample 8	2	2	1	1	1	2	1	1	2
Sample 9	1	1	1/2	1/2	1/2	1	1/2	1/2	1
Sum	14.0000	14.0000	7.0000	7.0000	7.0000	14.0000	7.0000	7.0000	14.0000

d) Form

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9
Sample 1	1	2	1	1	1/2	1	1	1/3	1/2
Sample 2	1/2	1	1/2	1/2	1/4	1/2	1/2	1/5	1/4
Sample 3	1	2	1	1	1/2	1	1	1/3	1/2
Sample 4	1	2	1	1	1/2	1	1	1/3	1/2
Sample 5	2	4	2	2	1	2	2	1/2	1
Sample 6	1	2	1	1	1/2	1	1	1/3	1/2
Sample 7	1	2	1	1	1/2	1	1	1/3	1/2
Sample 8	3	5	3	3	2	3	3	1	2
Sample 9	2	4	2	2	1	2	2	1/2	1
Sum	12.5000	24.0000	12.5000	12.5000	6.7500	12.5000	12.5000	3.8667	6.7500

e) View of the fracture

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9
Sample 1	1	1	1	1	2	1	1/2	1	1
Sample 2	1	1	1	1	2	1	1/2	1	1
Sample 3	1	1	1	1	2	1	1/2	1	1
Sample 4	1	1	1	1	2	1	1/2	1	1
Sample 5	1/2	1/2	1/2	1/2	1	1/2	1/3	1/2	1/2
Sample 6	1	1	1	1	2	1	1/2	1	1
Sample 7	2	2	2	2	3	2	1	2	2
Sample 8	2	2	2	2	3	2	1	2	2
Sample 9	1	1	1	1	2	1	1/2	1	1
Sum	10.5000	10.5000	10.5000	10.5000	19.000	10.5000	5.3333	10.5000	10.5000

The calculation of the priority vector (taste, aroma, color, shape and appearance of the fracture) of nine model samples of gluten-free custard cakes is presented in Table V and is obtained similarly to the one discussed above.

As a result we obtained a matrix of weight coefficients for each index of organoleptic evaluation (Table VI).

Multiplying the matrix of weight coefficients (Table VI) by the priority column vector (Table II) we obtain the weights of alternatives (gingerbread) {0.0538; 0.0497; 0.0906; 0.0715; 0.0899; 0.1169; 0.1925; 0.2104; 0.1247} in terms of preferences of respondents.

TABLE V. COMPARISON OF MODEL GLUTEN-FREE GINGERBREAD WITH RESPECT TO FIVE ORGANOLEPTIC PARAMETERS

a) Taste

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Average value
Sample 1	0.036	0.036	0.025	0.025	0.025	0.032	0.046	0.046	0.032	0.034
Sample 2	0.036	0.036	0.025	0.025	0.025	0.032	0.046	0.046	0.032	0.034
Sample 3	0.107	0.107	0.076	0.100	0.076	0.064	0.077	0.077	0.064	0.083
Sample 4	0.071	0.071	0.038	0.050	0.038	0.043	0.058	0.058	0.043	0.052
Sample 5	0.107	0.107	0.076	0.100	0.076	0.064	0.077	0.077	0.064	0.083
Sample 6	0.143	0.143	0.152	0.150	0.152	0.128	0.116	0.116	0.128	0.136
Sample 7	0.179	0.179	0.228	0.200	0.228	0.255	0.232	0.232	0.255	0.221
Sample 8	0.179	0.179	0.228	0.200	0.228	0.255	0.232	0.232	0.255	0.221
Sample 9	0.143	0.143	0.152	0.150	0.152	0.128	0.116	0.116	0.128	0.136

b) Aroma

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Average value
Sample 1	0.036	0.036	0.025	0.025	0.025	0.032	0.046	0.046	0.032	0.034
Sample 2	0.036	0.036	0.025	0.025	0.025	0.032	0.046	0.046	0.032	0.034

Sample 3	0.107	0.107	0.076	0.100	0.076	0.064	0.077	0.077	0.064	0.083
Sample 4	0.071	0.071	0.038	0.050	0.038	0.043	0.058	0.058	0.043	0.052
Sample 5	0.107	0.107	0.076	0.100	0.076	0.064	0.077	0.077	0.064	0.083
Sample 6	0.143	0.143	0.152	0.150	0.152	0.128	0.116	0.116	0.128	0.136
Sample 7	0.179	0.179	0.228	0.200	0.228	0.255	0.232	0.232	0.255	0.221
Sample 8	0.179	0.179	0.228	0.200	0.228	0.255	0.232	0.232	0.255	0.221
Sample 9	0.143	0.143	0.152	0.150	0.152	0.128	0.116	0.116	0.128	0.136

c) Colour

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Average value
Sample 1	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071
Sample 2	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071
Sample 3	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143
Sample 4	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143
Sample 5	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143

Sample 6	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071
Sample 7	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143
Sample 8	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143	0.143
Sample 9	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071	0.071

d) Form

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Average value
Sample 1	0.080	0.083	0.080	0.080	0.074	0.080	0.080	0.086	0.074	0.080
Sample 2	0.040	0.042	0.040	0.040	0.037	0.040	0.040	0.052	0.037	0.041
Sample 3	0.080	0.083	0.080	0.080	0.074	0.080	0.080	0.086	0.074	0.080
Sample 4	0.080	0.083	0.080	0.080	0.074	0.080	0.080	0.086	0.074	0.080
Sample 5	0.160	0.167	0.160	0.160	0.148	0.160	0.160	0.129	0.148	0.155
Sample 6	0.080	0.083	0.080	0.080	0.074	0.080	0.080	0.086	0.074	0.080
Sample 7	0.080	0.083	0.080	0.080	0.074	0.080	0.080	0.086	0.074	0.080

Sample 8	0.240	0.208	0.240	0.240	0.296	0.240	0.240	0.259	0.296	0.251
Sample 9	0.160	0.167	0.160	0.160	0.148	0.160	0.160	0.129	0.148	0.155

e) View of the fracture

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Average value
Sample 1	0.095	0.095	0.095	0.095	0.105	0.095	0.094	0.094	0.095	0.096
Sample 2	0.095	0.095	0.095	0.095	0.105	0.095	0.094	0.094	0.095	0.096
Sample 3	0.095	0.095	0.095	0.095	0.105	0.095	0.094	0.094	0.095	0.096
Sample 4	0.095	0.095	0.095	0.095	0.105	0.095	0.094	0.094	0.095	0.096
Sample 5	0.048	0.048	0.048	0.048	0.053	0.048	0.063	0.063	0.048	0.051
Sample 6	0.095	0.095	0.095	0.095	0.105	0.095	0.094	0.094	0.095	0.096
Sample 7	0.190	0.190	0.190	0.190	0.158	0.190	0.188	0.188	0.190	0.186
Sample 8	0.190	0.190	0.190	0.190	0.158	0.190	0.188	0.188	0.190	0.186
Sample 9	0.095	0.095	0.095	0.095	0.105	0.095	0.094	0.094	0.095	0.096

TABLE VI. COMPARISON OF THE ORGANOLEPTIC CHARACTERISTICS OF NINE MODEL FLOUR CONFECTIONERY PRODUCTS

	Taste	Aroma	Colour	Form	View of the fracture
Sample 1	0.034	0.034	0.071	0.080	0.096
Sample 2	0.034	0.034	0.071	0.041	0.096
Sample 3	0.083	0.083	0.143	0.080	0.096
Sample 4	0.052	0.052	0.143	0.080	0.096
Sample 5	0.083	0.083	0.143	0.155	0.051
Sample 6	0.136	0.136	0.071	0.080	0.096
Sample 7	0.221	0.221	0.143	0.080	0.186
Sample 8	0.221	0.221	0.143	0.251	0.186
Sample 9	0.136	0.136	0.071	0.155	0.096

Thus, according to experts' opinions on organoleptic indicators the first place has the 8-th sample.

Having ranked the samples of gluten-free gingerbread we get the following ranks:

- 1 place - sample 8
- 2 place - sample 7
- 3 place - sample 9
- 4 place - sample 6
- 5 place - sample 3
- 6 place - sample 5
- 7 place - sample 4
- 8 place - sample 1
- 9 place - sample 2.

VII. CONCLUSION

Thus, the best organoleptic properties got the gingerbread, which includes buckwheat flour, amaranth and linseed semi-fat in the ratio of 60:30:10, and the worst - a gingerbread based on rice flour, buckwheat and pumpkin seeds in a ratio of 60:20:20.

The results obtained correlate with the data obtained in the tasting organoleptic evaluation of samples on a point scale.

The application of mathematical approaches in the processing of expert assessments of the quality of food products on the basis of the hierarchy analysis method gives an objective final result. The constructed hierarchy of the global criterion (of quality) has flexibility. Adding new links to a well-structured hierarchy does not destroy its characteristics. Using the method when choosing alternatives for assessing the quality of food products, it is impossible to skip or ignore the feedbacks and reciprocal links between the components being investigated and the levels of the hierarchy, which minimizes the possibility of making a wrong decision.

REFERENCES

- [1] V.G. Kaishev, S.N. Seregin State and prospects of development of production of functional food products. Meat technologies, 2018. № 2. P. 54-57.
- [2] Barsukova N.V., Reshetnikov D.A., Krasilnikov V.N. Food engineering: technologies of gluten-free flour products. Scientific journal of NRU ITMO. A series of "Processes and devices of food production." - 2011

- [3] Nikitin I.A. Theoretical aspects of technology of effective functionality of food products. Collection of materials of the conference "Strengthening the Competitive Potential of Food Enterprises by Developing Effective Biotechnologies". FGANU Research Institute of the Bakery Industry. St. Petersburg Branch, 2016. P. 84-87.
- [4] Nikitin I.A., Kulakov V.G., Korovina E.S., Pyreseva A.I. Fragmentary research of the market of functional food products from gluten-free raw materials. Bread products, 2016. No. 11. P. 29-31.
- [5] Klokonos M.V., Semenkina N.G. Ways to confirm the functionality of bakery products. New in technics and technology of functional food products based on biomedical views Materials of the VI International Scientific and Technical Conference. Ministry of Education and Science of the Russian Federation. FGBOU VO "Voronezh State University of Engineering Technologies". - 2017. - P. 566-571.
- [6] Kuznetsova L.I., Dubrovskaya N.O. Polycomponent mixtures for the production of gluten-free products / Bakery production, 2014. No. 10. P. 20.
- [7] Schneider D.V. New program for the formation of recipes for gluten-free products. Bread products. - 2012. - No. 8. - P. 50-52.
- [8] Vokhmyanina N.V. The modern idea of celiac disease. S.-Pb.: Triad. 2009 -150 p.
- [9] Bazarnova Yu.G., Burova T.E., Marchenko V.I. Biochemical bases of processing and storage of raw materials of animal origin - SPb.: Prospect of Science, 2011. - 192 p.
- [10] Ivashov V.I., Andreenkov V.V., Solntseva G.L. Qualification of meat and meat products: overview information. - Moscow: AgroNIITEIMMP, 1989. - 48 p.
- [11] Nikitina M.A., Zakharov A.N., Shcherbinina E.O. Evaluation of organoleptic quality indices of meat products by statistical methods. Meat industry, 2017. № 5. P. 50-52.
- [12] Saaty T.L. Decision making with the analytical hierarchy process // International Journal services sciences, 2008. V.1. no. 1. pp. 83-98. URL: http://www.colorado.edu/geography/leyk/geog_5113/readings/saaty_2008.pdf (backs 08/08/2017).
- [13] T.L. Saaty. On the measurement of the intangible. Approach to relative measurements based on the eigenvector of the matrix of paired comparisons / Cloud of Science. - 2015. -V. 2. No. 1. - P. 5-39.
- [14] T.L. Saaty. Decision making with dependencies and feedbacks. Analytical relations. - Moscow: Lenard, 2015. - 360 p.
- [15] T.L. Saaty. Decision making with the analytical hierarchy process. - 2008. - Vol. 1. - p. 83-98.
- [16] T.L. Saaty. The seven pillars of the analytic hierarchy process // In: Multiple criteria decision making in the new millennium. - Berlin: Springer, 2001. - p. 1-15.
- [17] T.L. Saaty. Decision-making with the AHP: why is the principal eigenvector necessary // European Journal of Operational Research. - 2003. - Vol. 145. - p. 85-91.
- [18] Saaty T. Decision-making. The method of analyzing hierarchies. - M.: Radio and Communication, 1993. - 278 p.
- [19] Borisenko A.A., Sarycheva L.A., Borisenko A.A. Modeling. development and optimization of healthy food products. Ministry of Education and Science of the Russian Federation. Federal State budget educational institution of higher professional education "North Caucasus State Technical University". Stavropol, 2012. 197 p.
- [20] Dubrovskaya N.O., Kuznetsova L.I., Parahina O.I. Production of gluten-free bakery products using non-traditional plant materials / Bread products, 2016. No. 11. P. 36-37.
- [21] Kuznetsova L.I., Dubrovskaya N.O., Parakhina O.I. Improving the quality and nutritional value of gluten-free bread. Bakery of Russia, 2015. № 3. P. 19-22.
- [22] Chizhikova O.G. Dry mixtures with the addition of sea-buckthorn meal for gluten-free bakery products. Food industry. -2013. -No 3.-P. 18-19.
- [23] Klyachkin V.I. Statistical methods in quality management: computer technologies. - Moscow: Finance and Statistics, 2007. - 304 p.

3D Visualization of Sentiment Measures and Sentiment Classification using Combined Classifier for Customer Product Reviews

Siddhaling Urologin

Department of Computer Science
Birla Institute of Technology and Science, Pilani-Dubai,
Dubai, U.A.E.

Sunil Thomas

Department of Electrical and Electronics Engg.,
Birla Institute of Technology and Science, Pilani-Dubai,
Dubai, U.A.E.

Abstract—The Internet has wide reachability making many users to buy the products online using e-commerce websites. Usually, users provide their opinions, comments, and reviews about the products in social media, e-commerce websites, blogs, etc. The product review comments provided by the customers have rich information about the usage of the products they bought and their sentiments towards those products. In this research, we have collected reviews from Amazon.com and performed sentiment analysis to collect sentiment information. We have proposed 3D visualizations to represent sentiment information, such as sentiment scores and statistics about words used in the reviews. The 3D visualizations are useful to represent large sentiment related information and to have an in-depth understanding of sentiments of users. We have developed a combined classifier using Logistic Regression, Decision Tree and Support Vector Machine. From the reviews, we formed N-gram features using a bag of words and performed sentiment classification using combined classifier. On 10 fold cross-validation, a maximum classification rate for combined classifier of 90.22% is obtained for sentiment classification.

Keywords—Sentiment analysis; 3D visualization; sentiment classification; natural language processing; product reviews

I. INTRODUCTION

Many users are buying products online through e-commerce sites as they are widespread reachable by internet [1], [2]. The internet also provides an opportunity to the user to give their opinions in various forums such as social media, blogs, online e-commerce sites, etc. [2]. Many users provide reviews and opinions in natural language for the product they have come across. These reviews have wealth of knowledge [1] about the products and feeling of users towards the products. Sentiment analysis involves in mining the naturally expressed text to understand the feeling of people towards the entity of interest. Sentiment mining and analysis has found many application in areas of healthcare [3], [4], tourism [5], fraud detection [6], finance [7], politics [8], business [9], few more applications are listed in [10]. In [11] informatics, theoretic approach is used for classification of sentiments. A lexicon for sentiment analysis and concept level sentiment analysis is presented in [12]. The opinion mining of Amazon data is carried out in [13] using Support Vector Machines to summarize the unstructured data. As-LDA model is used in [14] for sentiment classification. Multimodal Naïve Bayes and

Decision tree classifiers are discussed in [15] for tweeter data sentiment analysis. A cloud integrated system for Support Vector Machines, Naïve Bayes and Neural Networks is presented in [16] for blog data sentiment classification. Several lexicon dictionaries such as [17], [18], have been used while determining sentiment analysis. In [18], Valence Aware Dictionary for Sentiment Reasoning (VADER), which is a rule-based model for sentiment analysis, is presented. VADER uses a lexicon list with sentiment measures to compute sentiment score for textual data. Sentiment analyzer examines the textual input to mine the feeling of user, which is present in the text and provides sentiment information indicative of feelings expressed by the users. Customer reviews for products are collected through e-commerce sites, blogs, social media etc., have rich knowledge about the products and their usability. Sentiment analyzer can be used to undermine sentiment information from the customer reviews. A large number of reviews being collected every day, better schemes are highly desirable to analyze sentiments of reviews and to provide visualization of sentiment information. In [19], authors have presented SentiView, which an interactive visualization system used to analyze people sentiments towards selected topic. It mines the data from online posts and uses uncertainty modeling to depict the changes in sentiments. Megan K et al. [20] have collected Twitter data during Gulf Oil Spill 2010 and they analyzed emotion of the broadcasted information in twitter. The emotion classification and analysis results are given with various visualizations. In [21], author has discussed the design, implementation of tools to extract, analyze and explore public messages. Then authors have performed sentiment analysis with a web application.

Motivated by these facts, we propose novel 3D visualization schemes to represent sentiment information obtained by sentiment analyzer. We utilize VADER [18] to mine the sentiments of electronic products reviews gathered from Amazon.com. The VADER provides the sentiment information such as compound score, positive score, negative score and neutral score, etc. The sentiment information has been used to construct 3D visualizations such as 3D surface plot, 3D column charts, 3D scatter plots, etc. in this research. The 3D visualization provides better schemes to undermine sentiment information especially for large review sets. Furthermore, we developed a new combined classifier for the sentiment classification. The combined classifier is built from

three base classifiers such as Logistic Regression, Decision Tree and Support Vector Machine. The voting technique is used to determine the resultant of combined classifier. The feature extraction is carried out in this research by taking N-gram features using a bag of words similar to method in [22]. This research paper is organized as follows. Review data preprocessing and filtering has been described in Section II. In Section III, the lexicon list and their measurements of VADER [18] are presented. Section IV elaborates the architecture on sentiment visualization and combined sentiment classifier. The experimental results are given in Section V and the conclusion is covered in Section VI.

II. DATA PREPARATION AND FILTERING

The customer review about the products have been gathered originally from Amazon.com and we have downloaded from [23], [24]. The customer review gathered from [23], [24] are in JSON format and an example for customer review has been shown in Table I. The format has various fields such as Title, Author, ReviewID, Overall, Content, and Date. Here Overall Rating is a five-star rating where the user can assign zero (lowest) to five (highest) stars. The content is textual description given by the user as their review comments for a product. Each customer review thus obtained will undergo preprocessing and filtering. All the text parts are retained and converted to lower case letter. Stop words are identified using the list of [18] and these stop words are removed from the reviews to retain only essential words. Then, we utilized VADER lexicon and sentiment analyzer [18] to compute the score of words and determine the sentiment score of the customer reviews. Then overall sentiment polarity of the review is found to be positive, negative or neutral based on the total sum of sentiment score of a review.

TABLE I. JSON FORMAT OF CUSTOMER REVIEW FROM AMAZON.COM

JSON Format	{ "Reviews": [{ "Title": "Overall a nice laptop for around \$1100", "Author": null, "ReviewID": "R2Y4WQSYMYCE24", "Overall": "4.0",
{	"Content": "... did not really review Model 2681CU1 as shown on this website. This models comes with Pentium 4 -Mobile 1.8 Ghz,...", "Date": "October 30, 2003"},
Title:	
Author:	
ReviewID:.....	
Overall:.....	
Content:.....	
Date:.....	
}	

III. LEXICON FOR SENTIMENT SCORE

In this research, we have used VADER lexicon and sentiment analyzer [18] to determine sentiment scores of a review. The VADER lexicon is a well-trained list of words with the polarity values especially for micro blogs of social media. This lexicon focuses on sentiment emphasis by different social media web sites. The list also includes many of generally used emotions such as “:), “:(, acronyms such as “LOL” and slang “nah”, etc. The lexicon contains the word with their polarity value between -4 to +4. In Table II few examples for VADER lexicon words with polarity values are presented.

TABLE II. EXAMPLE FROM VADERLEXICON

Word	Polarity
Okay	+0.9
Good	1.9
Great	3.1
Horrible	-2.5
:(-2.2

IV. SENTIMENT VISUALIZATION AND CLASSIFICATION

After collecting customer reviews in the JSON format, the content field is extracted, which is the commentary description given by a customer. The sentiment analysis of the review comments are carried out as shown in Fig. 1 using VADER sentiment analyzer of [18]. A list of lexicons related to micro blog along with their sentiment measures is used in VADER sentiment analyzer. This analyzer uses five heuristic rules, based on grammatical and syntactical conventions used in natural language while expressing human emotions. Using VADER sentiment analyzer, we obtain sentiment information such as compound sentiment score, positive score, negative score and neutral score for a review comment. Motivated by current research [19]-[21], we are proposing novel schemes for representing sentiment information of customer reviews into 3D visualization charts. The 3D representations of sentiment information are most useful to gain further insight of the customer reviews.

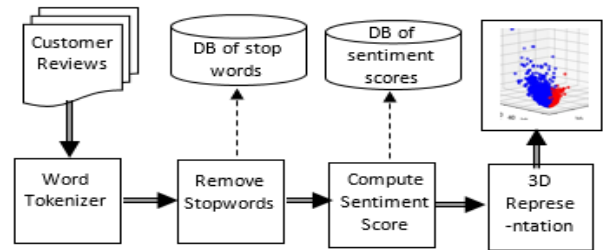


Fig. 1. Customer reviews to 3D representation.

We propose a new combined classifier of Logistic Regression, Decision Tree and Support Vector Machine classifier using voting technique to perform sentiment classification. The combined classifier using voting technique is depicted in Fig. 2. The customer reviews undergo preprocessing in which sentence tokenization and stop words removal are carried out. Thereafter N-gram feature vector is constructed for each review using bag of words similar to [22]. The sentiment classification of the reviews is performed using a combined classifier, which consists of three base classifiers such as Logistic Regression, Decision Tree and Support Vector Machine. The resultant of combined classifier is found by taking majority in voting of base classifiers.

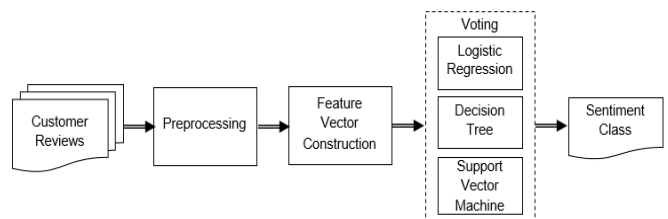


Fig. 2. Combined classifier using voting technique.

TABLE III. STATISTICS FOR THE SENTIMENT OF REVIEWS [SENTIWORDS IS SENTIMENT WORDS]

Review	Negative Sentiment Score	Neutral Sentiment Score	Positive Sentiment Score	Compound Sentiment Score	Num. Of Negative SentiWords	Num. Of Neutral SentiWords	Num. Of Positive SentiWords
looking external dvd player kids could watch dvds asus eee netbook , exactly got	0.126	0.679	0.195	0.8174	6	84	9
dont buy unit . dvd-rom needed dvd writer . read dvds read cd . absolute waste money .	0.177	0.823	0	-0.4215	1	13	0
needed good thing works , sometimes take cd input , guess dumb moments	0.217	0.592	0.191	-0.1027	1	9	1
worth every penny ! shippment speedy ! performed better expected . would recommend...	0	0.565	0.435	0.816	0	11	3
small asus netbook decided needed external dvd drive . one works fine . far 've imported cds...	0.115	0.809	0.077	-0.2263	1	19	1
mine shipped instruction manual documentation whatsoever . plugged , whirred though sensed...	0	1	0	0	0	26	0
dvd/cd drive laptop ended failing six months ago , leaving stranded laptop needed reformatted	0.043	0.912	0.046	0.0772	2	100	2
great buy inexpensive price.there n't much instructions came dvd/cd drive . think one small paragraph..	0.072	0.797	0.131	0.4295	1	25	1
received sooner expected ! part christmas gift twins go netbooks . fact came early , bonus , worrying	0.114	0.569	0.317	0.6476	1	12	2
p.o.s dead arrival . cheap chinese junk chinglish manual . power supply included . purchased...	0.203	0.797	0	-0.68	1	18	0

V. EXPERIMENTAL RESULTS

We have gathered customer reviews for electronic products from [23], [24], which are originally collected from Amazon.com. The electronic review data set contain 1,689,188 user comments. In our experiments, the sentiment visualization is presented on four electronic products and the sentiment classification is carried out on 12,500 review comments. We have carried out preprocessing by removing URL links, tags, stop words and each comment is converted to lower case letters. Each review is subjected to sentiment analysis using VADER analyzer [18]. In Table III, we have tabulated sentiment information obtained from VADER for few reviews on the electron products. In column 1, user review comments given (partial comments are shown). Next four columns show negative, neutral, positive and compound sentiment scores of a review. The number of negative, neutral and positive sentiment words found in that review is given 6, 7 and 8th columns respectively.

VI SENTIMENT VISUALIZATION

Further analysis and visualization are carried out on four electronic products reviews to reveal customer opinions in depth. We have obtained sentiment information using [18] on

reviews of products External USB DVDCD, GE 72887 Superadio III Portable AMFM Radio, NETGEAR Prosafe FS105NA, Panasonic On-Ear Stereo Headphones. There are 199 comments on External USB DVDCD with 11,630 words, for GE 72,887 Superadio III Portable AMFM Radio 265 reviews with 33,973 words, for NETGEAR Prosafe FS105NA423 comments with 25310 words and for Panasonic On-Ear Stereo Headphones 1692 comments with 1,06,284 words. The sentiment of a review is considered as positive when its compound score is greater than zero, it is neutral when the compound score is equal to zero and when the compound score is less than zero its sentiment is negative. Each word in a review contribute in compound sentiment score. It is interesting to observe the variation of compound sentiment score against number of positive and negative words.

In Fig. 3, a representation of compound score is shown as 3D surface against number of positive sentiment words and negative sentiment words. In this figure, compound score is accumulated for all the reviews of a product with respect to number of positive words and negative words. In Fig. 3(a), accumulated compound sentiment score is plotted as a 3D surface for reviews of External USB DVDCD. Similarly in Fig. 3(b), 3(c) and 3(d) accumulated compound score is plotted against number of negative and positive words.

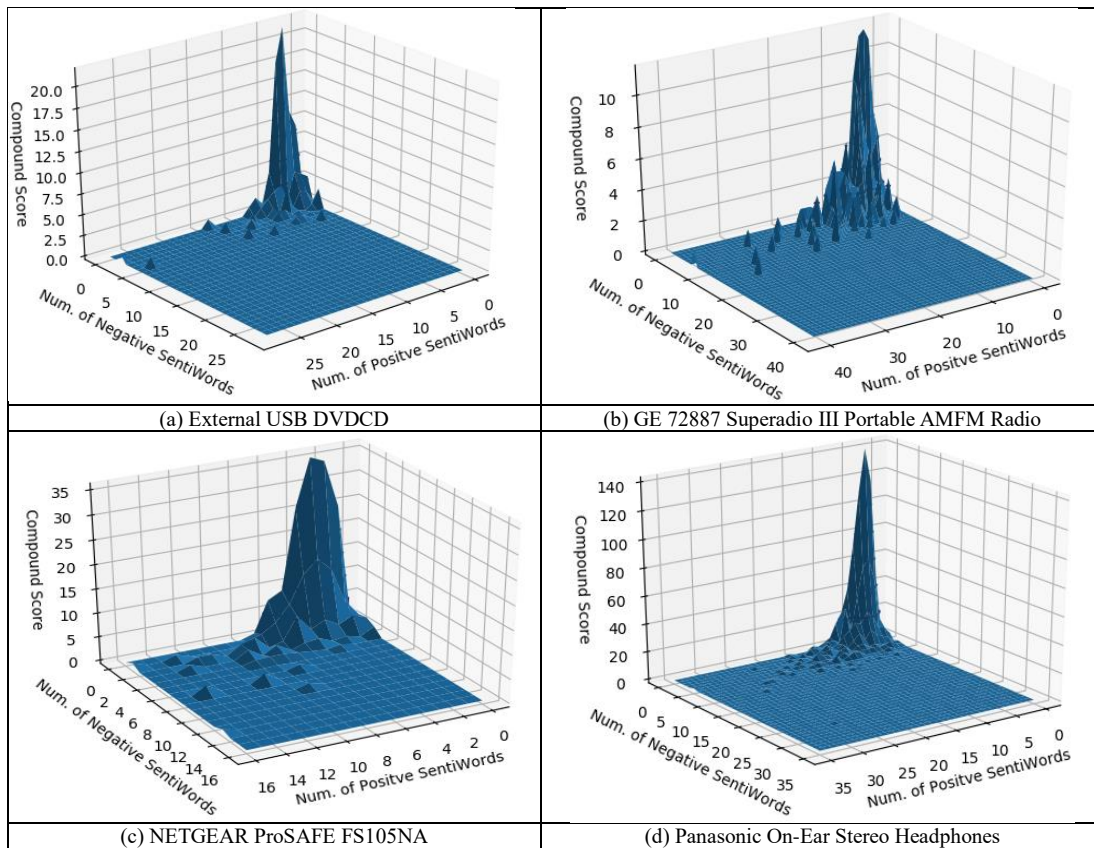


Fig. 3. (a) External USB DVDCD; (b) GE 72887 Superadio III Portable AMFM Radio; (c) NETGEAR Prosafe FS105NA; (d) Panasonic On-Ear Stereo Headphones.

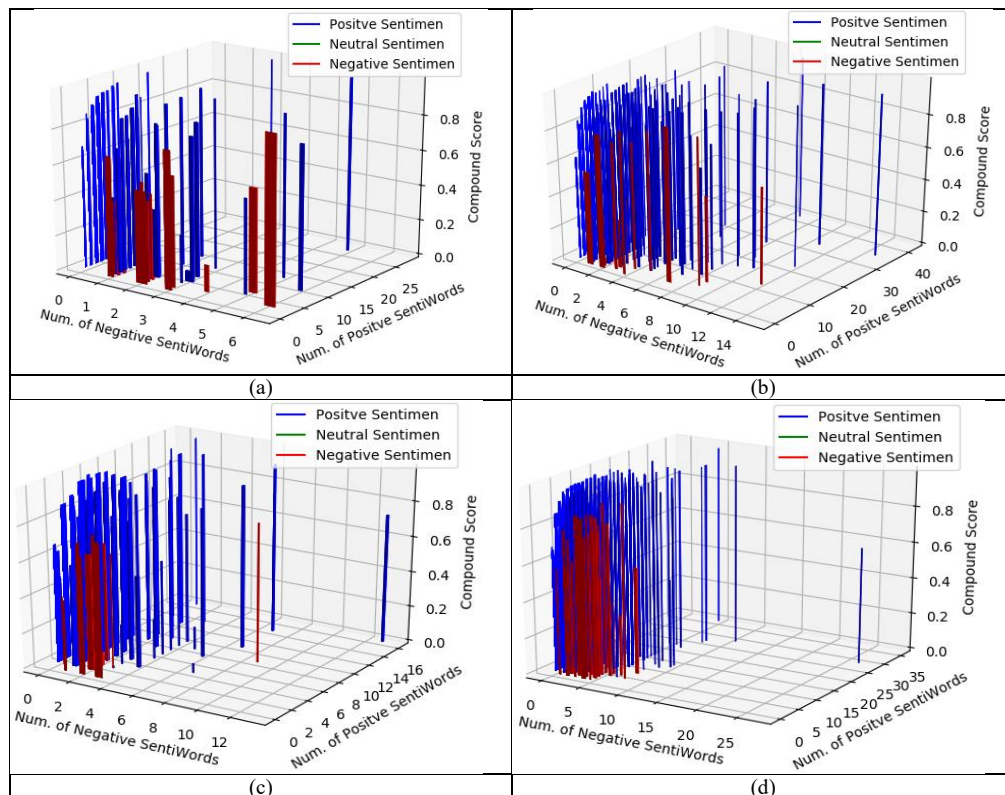


Fig. 4. (a) External USB DVDCD; (b) GE 72887 Superadio III Portable AMFM Radio; (c) NETGEAR Prosafe FS105NA; (d) Panasonic On-Ear Stereo Headphones.

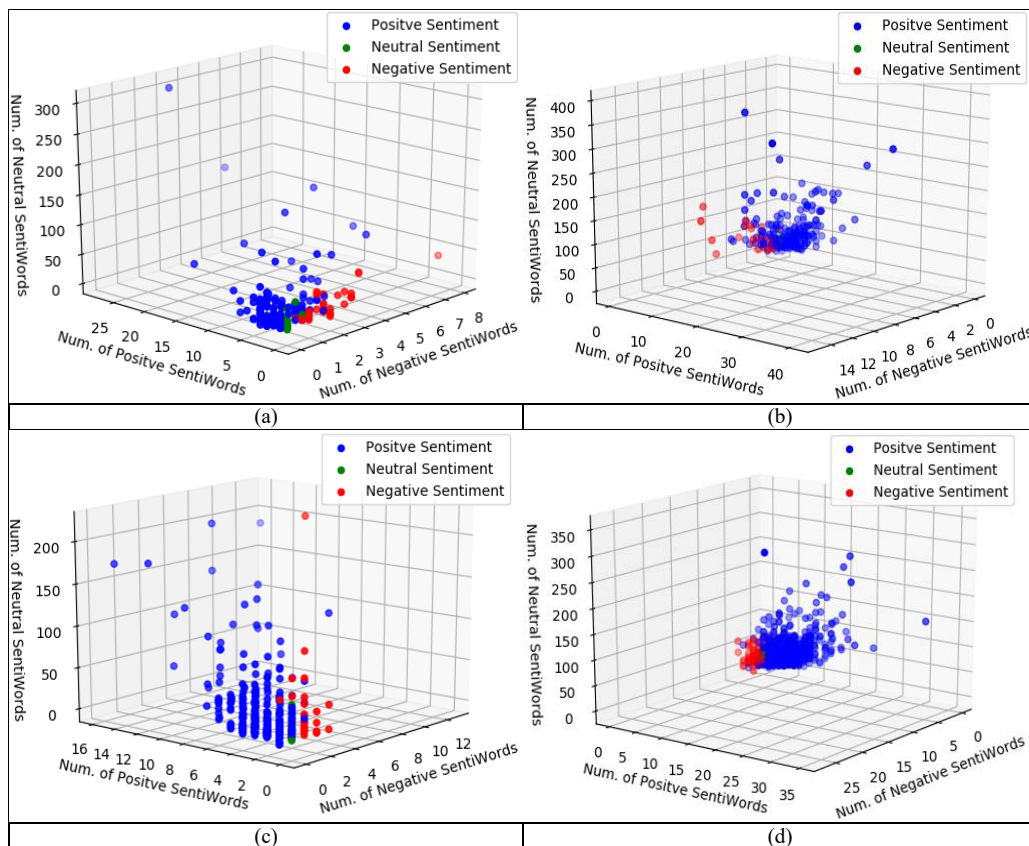


Fig. 5. (a) External USB DVDCD; (b) GE 72887 Superadio III Portable AMFM Radio; (c) NETGEAR Prosafe FS105NA; (d) Panasonic On-Ear Stereo Headphones.

These representations reveal contribution of number of positive and negative words for compound sentiment score. As depicted in Fig. 3(a), the highest score in the accumulated compound score observed as 21.64 with number of negative words as zero and positive words as three. More positive comments have been observed by customers than negative and majority of customers used around 10 positive words. There are few negative comments observed with fewer than seven words. More variation in the accumulated compound sentiment score is observed for reviews of GE 72887 Superadio III Portable AMFM Radio as shown in Fig. 3(b).

Fig. 4 shows 3D column chart of compound scores for reviews of a product. The compound review score is depicted against number of positive and negative words. In Fig. 4(a) compound scores for 199 reviews of External USB DVDCD is represented against number positive and negative words. As shown in Fig. 4(a), more positive sentiments can be observed compared to negative. Also few number of words are used to express sentiment negatively, while more number of words are used express positive opinions. The 3D surface views as in Fig. 3 and 3D column view as shown in Fig. 4 of compound sentiment score highlight usability of a number of words to express sentiment, which is important to understand the characteristics of a product reviews.

In Fig. 5, a review is represented as a point with coordinates as a number of positive sentiment words along the x-axis, negative sentiment words along they-axis and neutral sentiment words along the z-axis. Fig. 5(a) is a 3D scatter plot for all 199 reviews indicating the distribution of sentiments type with respect to the number of words present in reviews. This figure indicates that more positive reviews than negative or neutral are given by customer on External USB DVDCD. The 3D scatter plot reveals the distribution of reviews observed on sentiment type with respect the number of words.

We have collected 10 words having maximum sentiment score and minimum sentiment score for each product and represented using donut charts in Fig. 6 and 7. The top 10 words having maximum sentiment score represent positive sentiment are depicted in Fig. 6. For each word, its sentiment score along with the percentage contribution to top 10 maximum sentiment score words has been shown. Fig. 7 shows top 10 words with minimum sentiment scores representing negative sentiment. The compound sentiment score and the percentage contribution to top 10 negative sentiment words have been depicted in Fig. 7.

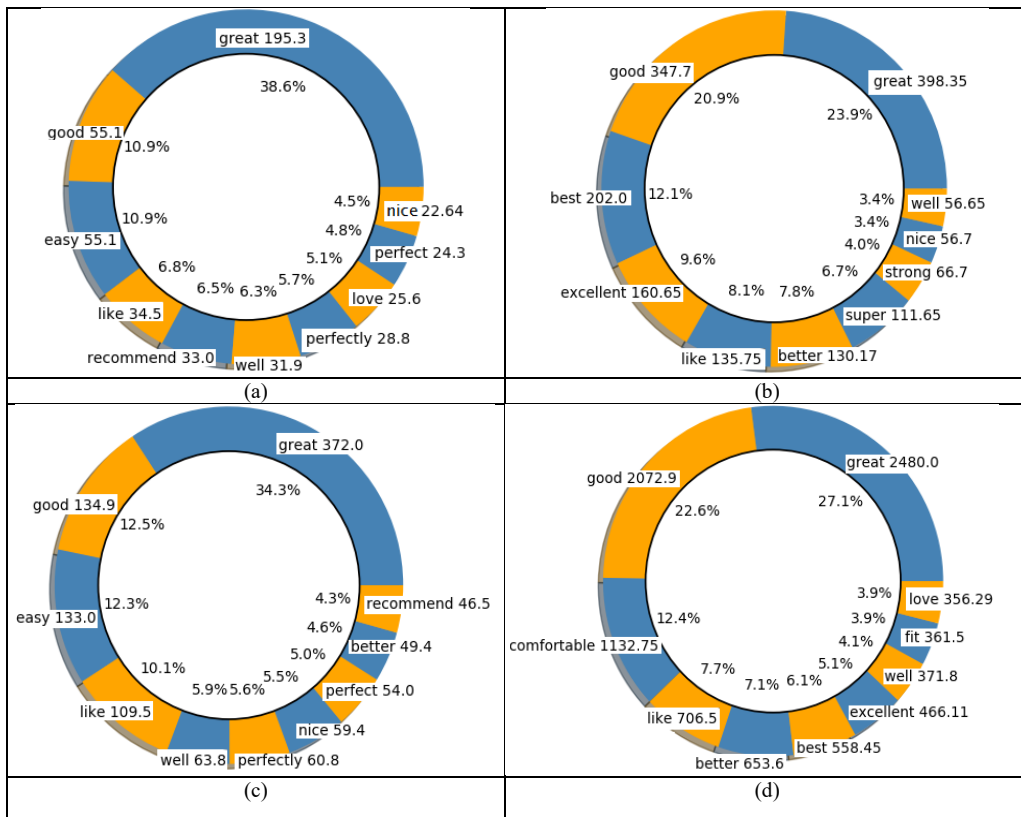


Fig. 6. (a) External USB DVDCD; (b) GE 72887 Superadio III Portable AMFM Radio; (c) NETGEAR Prosafe FS105NA; (d) Panasonic On-Ear Stereo Headphones.

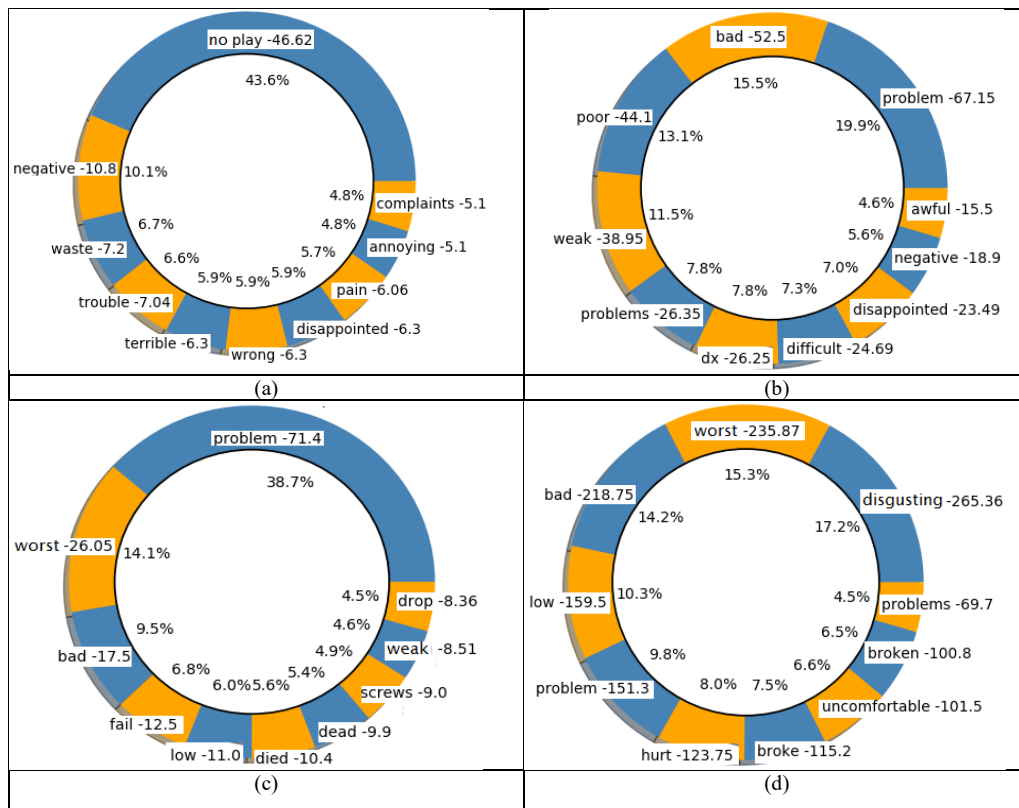


Fig. 7. (a) External USB DVDCD; (b) GE 72887 Superadio III Portable AMFM Radio; (c) NETGEAR Prosafe FS105NA; (d) Panasonic On-Ear Stereo Headphones.

VI. SENTIMENT CLASSIFICATION

We have considered 12500 review comments on electronic products from Amazon data set. Each review comment is subjected to preprocessing where in each sentence is tokenized and stop words are eliminated. Then we utilized sentiment analyzer of VADER [18] and determined the sentiment information for user comments. The sentiment analyzer VADER provides the scores such as positive score, negative score, neutral score and compound score. When compound score of a review is greater than zero then review is considered as positive and when compound score is less than zero it is taken as negative sentiment. Thereby we constructed a dataset of 12500 review comments along with their sentiment type as positive or negative and carried out sentiment classification. For a review comment, a bag of word features vector is constructed by taking N-gram size feature matrix similar to the method in [22]. To select the words from the collection of 12500 reviews, we have computed frequency of occurrence of each word in the dataset. Then the bag of words is created by using selected ‘n’ number of most frequent words. Ten most frequent important words along with their frequency are depicted in Fig. 8. In the experiments we have chosen ‘n’ as 2000, 3000, 4000, 5000 and 6000.

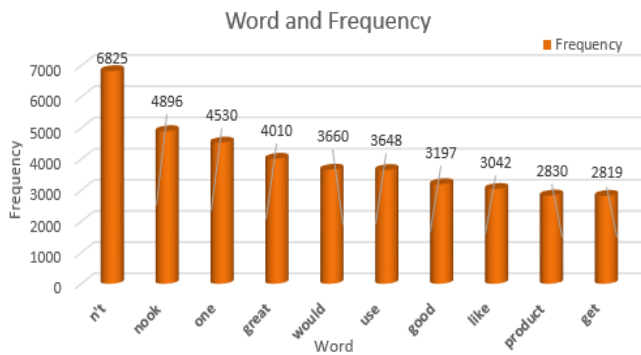


Fig. 8. Ten most frequent word in the dataset of user reviews.

TABLE IV. NAÏVE BAYES CLASSIFIER PERFORMANCE.

Sl. Num.	n	Classif. Rate in %	Misclassif. Rate in %
1	2000	82.94	17.06
2	3000	83.43	16.57
3	4000	84.26	15.74
4	5000	84.90	15.10
5	6000	85.14	14.86

Further, we have conducted experiments to predict sentiment of a review as either positive or negative sentiment using Naïve Bayes Classifier of NLTK [25]. The 10-fold cross validation is carried out on Naïve Bayes Classifier to perform classification of sentiments with ‘n’ as 2000, 3000, 4000, 5000 and 6000. The cross-validation accuracy in terms of classification rate and misclassification rate are shown in Table IV.

Next, the sentiment classification is carried out using three classifiers such as Logistic Regression, Decision Tree and Support Vector Machine independently. The 10-fold cross validation for all three classifiers with varying feature vector size n=2000, 3000, 4000, 5000 and 6000 is conducted and results are shown in Table V. In columns 3, 5 and 7 is shown the classification rate for Logistic Regression Classifier, Decision Tree Classifier and Support Vector Machine Classifier, respectively. In columns 4, 6 and 8 the misclassification rates are given. Out of three classifiers, Logistic Regression Classifier gives better classification rates ranging from 88.52% to 88.93%. The Decision Tree Classifier provides the least classification rate between 82.34% and 83.09% as tabulated in Table V.

A novel combined classifier using voting technique is developed for sentiment classification. Three classifiers such as Logistic Regression, Decision Tree and Support Vector Machine are used in the combined classifier and resultant is determined by majority voting technique as shown in Fig. 2.

TABLE V. LOGISTIC REGRESSION, DECISION TREE AND SUPPORT VECTOR MACHINE PERFORMANCE INDEPENDENTLY

Sl.Num.	n	Logistic Regression Classifier		Decision Tree Classifier		Support Vector Machine	
		Classif. Rate in %	Misclassif. Rate in %	Classif. Rate in %	Misclassif. Rate in %	Classif. Rate in %	Misclassif. Rate in %
1	2000	88.62	11.38	82.34	17.66	86.42	13.58
2	3000	88.52	11.48	83.03	16.97	86.42	13.58
3	4000	88.72	11.28	83.03	16.97	86.42	13.58
4	5000	88.63	11.37	82.99	17.01	86.42	13.58
5	6000	88.93	11.07	83.09	16.91	86.42	13.58

TABLE VI. COMBINED CLASSIFIER PERFORMANCE

Sl.Num	n	Combined Classifier in %	
		Classif. Rate	Misclassif. Rate
1	2000	89.12	10.88
2	3000	89.10	10.90
3	4000	90.22	9.78
4	5000	89.55	10.45
5	6000	89.15	10.85

The performance of the combined classifier is measured on 10-fold cross validation. The results of cross validation such as classification rate and misclassification rate are given in Table VI. Maximum classification rate of 90.22% with feature vector size $n=4000$ and minimum classification rate of 89.10% with feature vector size $n=3000$ are observed. For most of feature vector size $n=2000, 3000, 4000, 5000$ and 6000 the performance of the combined classifier is better than base classifiers such as Logistic Regression, Decision Tree and Support Vector Machine.

VII. CONCLUSION AND FUTURE WORK

In recent years e-commerce websites are gaining popularity. Users find it easier and convenient to buy products online through various e-commerce websites. Moreover, with the increase in reachability of internet, there is a substantial growth in customer provided information in terms of reviews comments for various products. Due to which, there is interest in mining useful information from customer review and comments. In this research, we performed sentiment analysis on reviews gathered from Amazon.com for electronic products. We utilized the VADER, which is a rule based sentiment analyzer and uses lexicon with sentiment measures to compute sentiment scores. VADER provides sentiment information such as compound, positive, negative, neutral score, etc. for the reviews. In a large set of data, it is important to visualize the sentiment information and related statistics for a better insight of reviews. In this research, 3D visualization such as 3D surface, 3D column and 3D scatter plot charts are presented. The 3D scatter plot, provide insight into the distribution of sentiment type with respect to a number of words in the reviews. The 3D surface and 3D column charts show variations in sentiment score with respect to a number of positive and negative words. Furthermore, we have utilized donut charts to represent most 10 positive and negative sentiment words in a set of reviews. In addition, we developed a new combined classifier for sentiment classification. The combined classifier has three base classifiers such as Logistic Regression, Decision Tree and Support Vector Machine. The combined classifier utilizes voting technique to determine resultant class. A bag words for reviews is created to collect N-gram features from reviews and performance of the classifiers is evaluated on 10-fold cross validation. For the combined classifier a maximum classification rate of 90.22% for 10-fold cross validation was obtained. The future research work will utilize the 3D visualizations such as, 3D surface, 3D column and 3D scatter plots and determine crucial features to perform sentiment classification.

REFERENCES

- [1] Z. Hu, J. Hu, W. Ding and X. Zheng, "Review Sentiment Analysis Based on Deep Learning," 2015 IEEE 12th International Conference on e-Business Engineering, Beijing, 2015, pp. 87-94. doi: 10.1109/ICEBE.2015.24.
- [2] Youcef Baghdadi, "A framework for social commerce design, In Information Systems, Volume 60, 2016, pp. 95-113, ISSN 0306-4379, <https://doi.org/10.1016/j.is.2016.03.007>.
- [3] Vincenza Carchiolo, Alessandro Longheu, and Michele Malgeri, "Using Twitter Data and Sentiment Analysis to Study Diseases Dynamics", 2015 Proceedings of the 6th International Conference on Information Technology in Bio- and Medical Informatics - Volume 9267, pp. 16-24. doi:10.1007/978-3-319-22741-2_2
- [4] Denecke K, "Sentiment Analysis from Medical Texts. In: Health Web Science. Health Information Science".2015 Springer, Cham. doi:https://doi.org/10.1007/978-3-319-20582-3_10
- [5] D Gräbner, M Zanker, G Fliedl, M Fuchs, "Classification of Customer Reviews based on Sentiment Analysis",2012 Information and Communication Technologies in Tourism pp 460-470.
- [6] Gann W-JK, Day J, Zhou S, " Twitter analytics for insider trading fraud detection system" 2014. Proceedings of the second ASE international conference on Big Data. ASE. May 27 - May 31, 2014, Stanford, CA, USA, 94305
- [7] Siddhaling Urolagin, "Text Mining of Tweet for Sentiment Classification and Association with Stock Prices," 2017 International Conference on Computer and Applications (ICCA), Doha, 2017, pp. 384-388. doi: 10.1109/COMAPP.2017.8079788
- [8] Kartik Singhal, Basant Agrawal, Namita Mittal, "Modeling Indian General Elections: Sentiment Analysis of Political Twitter Data", Advances in Intelligent Systems and Computing, vol 339. Springer, New Delhi, pp 469-477.
- [9] F. Xing and J. Zhan, "Sentiment Analysis Using Product Review Data," Journal of Big Data, vol. 2:5,2015. <https://doi.org/10.1186/s40537-015-0015-2>
- [10] M.Walaa, A. Hassan, and H. Korashy, "Sentiment Analysis Algorithms and Applications: A Survey," Ain Shams Engineering Journal, vol.5, no. 4, pp. 1093—1113, 2014
- [11] Lin Y, Zhang J, Wang X, Zhou A, " An information theoretic approach to sentiment polarity classification". 2012, Proceedings of the 2nd Joint WICOW/AIRWeb Workshop on Web Quality, ACM, New York, NY, USA, pp 35–40.
- [12] A. Mudinas, D. Zhang and M. Levene, "Combining lexicon and learning based approaches for concept-level sentiment analysis,"2012, Proceedings of the First International Workshop on Issues of Sentiment Discovery and Opinion Mining, Beijing, China ugust 12 , 2012, Article No. 5, doi:10.1145/2346676.2346681
- [13] T. H. A. Soliman, M. A. Elmasry, A. R. Hedar and M. M. Doss, "Utilizing support vector machines in mining online customer reviews," 2012 22nd International Conference on Computer Theory and Applications (ICCTA), Alexandria, 2012, pp. 192-197. doi: 10.1109/ICCTA.2012.6523568
- [14] J. Liang, P. Liu, J. Tan, and S. Bai, "Sentiment Classification Based on AS-LDA Model," Procedia Computer Science, vol. 31, pp. 511—516, 2014.
- [15] A. P. Jain and P. Dandannavar, "Application of Machine Learning Techniques to Sentiment Analysis,"2nd International Conference on Applied and Theoretical Computing and Communication Technology, Bangalore, 2016, pp. 628--632.
- [16] R. Arulmurugan, K. R. Sabarmathi, and H. Anandakumar, "Classification of sentence level sentiment analysis using cloud machine learning techniques,"Cluster Computing, <https://doi.org/10.1007/s10586-017-1200-1>
- [17] Y. Wang, Y. Zhang, and B. Liu, "Sentiment Lexicon Expansion Based on Neural PU Learning, Double Dictionary Lookup, and Polarity Association,"Conference on Empirical Methods in Natural Language Processing, Copenhagen, 2017, pp. 7–11.
- [18] C.J. Hutto and Eric Gilbert, "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text,"Eighth International AAAI Conference on Weblogs and Social Media, 2014. [available online: 08 Nov 2017].
- [19] Changbo Wang, Zhao Xiao, Yuhua Liu, Yanru Xu, Aoying Zhou, Kang Zhang, "SentiView: Sentiment Analysis and Visualization for Internet Popular Topics", in IEEE Transactions on Human-Machine Systems, Volume: 43, Issue: 6, Nov. 2013.
- [20] Megan K. Torkildson, Kate Starbird, Cecilia Aragon, "Analysis and Visualization of Sentiment and Emotion on Crisis Tweets", in International Conference on Cooperative Design, Visualization and Engineering CDVE 2014: Cooperative Design, Visualization, and Engineering pp 64-67,2014.
- [21] Farina J., Mazuran M., Quintarelli E., "Extraction, Sentiment Analysis and Visualization of Massive Public Messages",. In: Catania B. et al.

- (eds) *New Trends in Databases and Information Systems. Advances in Intelligent Systems and Computing*, vol 241. Springer, Cham, 2014.
- [22] A. Deshwal and S.K. Sharma, "Twitter sentiment analysis using various classification algorithms," 5th International Conference on Reliability, Infocom Technologies and Optimization, Noida, 2016, pp. 251-257
- [23] J. McAuley and A. Yang "Addressing Complex and Subjective Product-Related Queries with Customer Reviews," Proceedings of the 25th International Conference on World Wide Web, Montreal, 2016, pp. 625–635
- [24] L. Jure and S. Rok, "SNAP: A General-Purpose Network Analysis and Graph-Mining Library", *ACM Transactions on Intelligent Systems and Technology*, vol.8, no. 1, pp. 1–20, 2016.
- [25] Edward Loper, Steven Bird, "NLTK: the Natural Language Toolkit", *ETMTNLP '02 Proceedings of the ACL-02 Workshop on Effective tools and methodologies for teaching natural language processing and computational linguistics - Volume 1*, Pages 63-70 , 2002.

A Novel Energy Efficient Mobility Aware MAC Protocol for Wireless Sensor Networks

Zain ul Abidin Jaffri, Asif Kabir
College of Communication Engineering
Chongqing University
Chongqing 400044, China

Gohar Rehman Chughtai, S. Sabahat H. Bukhari
College of Computer Science
Chongqing University
Chongqing 400044, China

Muhammad Arshad Shehzad Hassan
College of Electrical Engineering
Chongqing University
Chongqing 400044, China

Abstract—Dealing with mobility at the link layer in an efficient and effective way is a formidable challenge in Wireless Sensor Networks due to recent boom in mobile applications and complex network scenarios. Most of the current MAC protocols proposed for WSNs generally focus on stationary network and usually provide feeble network performance in situations where mobile nodes are involved. Many MAC protocols are proposed and techniques are developed to support mobility but they undergo massive energy consumption and latency problems due to frequent connection setup and breakup. In this paper, we propose a new energy efficient mobility aware based MAC protocol (EEMA-MAC), which work efficiently in both stationary and mobile scenarios with less energy consumption. In this protocol the member nodes have sleep and awake time same like existing S-MAC protocol but it expedite the connection setup and efficiency as Cluster Head (CH) has extended wake up time and less sleep time. Simulation results show that this mechanism is effective to avoid frequent disconnection of nodes and performs well in terms of energy consumption, throughput and packet loss as compared with existing protocols, such as S-MAC and MS-MAC.

Keywords—Wireless sensor networks; energy efficiency; Media Access Control (MAC); mobility aware; cluster head

I. INTRODUCTION

Wireless Sensor Networks (WSNs) have major contribution in the recent boom of technological advancements. This technology has emerged with high potential to sense physical phenomena like environmental monitoring, medical systems, seismic events, smart spaces, etc. [1] by initially processing the collected data locally and then delivering this data over a multi-hop link [4]. Basically the network comprises of numerous distributed nodes which self-organize themselves as a wireless multi-hop network. Each node is battery operated and has one or more sensors, low power radios and embedded processors [1], [10]. Typically the nodes establish and maintain the network and coordinate to accomplish the assigned task. The network offers many attractive features due to their small size and can be easily deployed to the inaccessible places and the areas which are expensive for wired systems. Wireless Sensor Networks have large number of applications. They are used for

health monitoring, fire and smoke, temperature, vibrations etc. Similar applications include structural health monitoring using accelerometer sensors [6], observing the activities of sea birds by using light, barometric, temperature and humidity sensors [5], to monitor active volcanos using infrasonic and seismic sensors [7] and to examine large water transmission pipelines using acoustic and hydraulic sensors [8].

The main focus of most of the applications above is on the nodes which are generally static after the deployment is done. In recent years mobility has emerged as a major constituent of many WSN applications [12]. The concept of mobility in WSNs is that few network elements like base station, sensor nodes, actuators or monitored targets can be mobile to enhance the capabilities of the system so that it can react quickly in many emergency situations [11]. As gradually the quantity of mobile units monitored by sensors are increasing, the role of mobility is also becoming highly important in WSNs. Most of the medium access control protocol designed recently are adopted for stationary networks. In these stationary networks the topology is fixed and the neighboring nodes remain unmoved for long period of time. As per the survey, only few MAC protocols support mobility, which creates room for research regarding mobility in WSNs [2], [3]. The MAC protocols can be divided into three types; Time division Multiple Access based (TDMA), contention based, and hybrid. In our current research we will focus on the contention-based MAC protocols. The contention-based protocols are less complicated in terms of using scheduling algorithm for channel division and maintain equal probability to access the channel in terms of both static and mobile nodes. In addition to this they do not encounter the overhead produced by hybrid protocols.

The main objective of our novel MAC protocol for is to work effectively and efficiently in both mobile and stationery scenarios with minimum energy consumption. In order to achieve our objective, S-MAC [1], a Sensor Medium Access control protocol is our first milestone and a starting point and then extends the protocol to support mobile sensors. The mobility-aware Sensor MAC protocol (MS-MAC) [9] works almost similar to the S-MAC for energy conservation when the nodes are stationary.

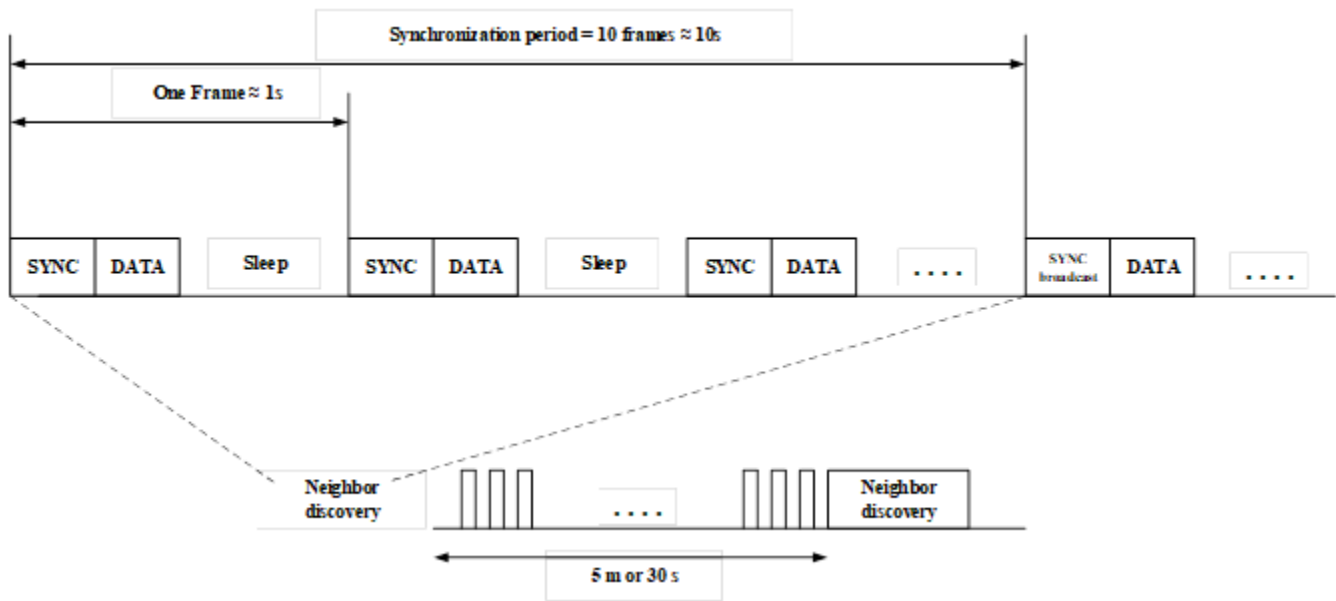


Fig. 1. Example S-MAC synchronization period description.

This medium access scheme can also switch to function like IEEE 802.11 for mobile ad-hoc scenario. S-MAC is basically a CSMA/CA MAC protocol, based on IEEE 802.11 which initiates periodic coordinated sleep/wakeup duty cycles as in Fig. 1, thus enhancing the lifetime of sensor nodes. To retain synchronization, for all predefined number of cycles, every node broadcasts its schedule in a SYNC message, so that its neighbors can update that information in their schedule tables accordingly. In addition to this, every node follow periodically the neighbor discovery scheme, so that the two neighbor nodes couldn't see each other (for example, because of SYNC packet corruption, interference, or because the medium remain busy and SYNC packets cannot be transmitted in time). The synchronization period (10 seconds) is repeated after every 2 minutes for this scheme. S-MAC does not need all the nodes in whole network, but only in each virtual cluster to synchronize. Border nodes amongst virtual clusters are required to follow more than one schedule.

The remaining paper is organized as follows: Section II elaborates the complete working of the proposed MAC protocol EEMA-MAC along with its CH election mechanism and RSSI estimation. Performance evaluation and results are described in Section III. Finally, concluding remarks are given in Section IV.

II. ENERGY EFFICIENT MOBILITY AWARE MAC (EEMA-MAC) PROTOCOL

The SMAC protocol performs very well in case of stationery nodes and when the connection setup and breakup is not rapid. The mobile node has to wait for synchronization period to detect a SYNC message, if it desires to setup a connection in another. This wait time might be long and might get disconnected from the network. MS-MAC considers only for the mobility but did not propose a solution for energy efficiency. The incidence of synchronization is varied as per the node speed which lets to drain high amount of energy for

speed moving vehicles. Hence we propose an algorithm which helps to expedite the connection setup and also provide energy efficiency.

The proposed protocol is an energy efficient mobility aware based MAC protocol. The basic structure of the proposed algorithm is derived from widely used SMAC protocol. It consists of periodic sleep and wakeup cycle as in S-MAC protocol with a listen and sleep period of 100ms. Each node periodically broadcasts a SYNC packet which consists of consists of source address, sleep time, received signal strength, speed, x-coordinate, y-coordinate and remaining energy of sender node. By periodically broadcasting the SYNC packet, nodes interchange their schedules to its close neighbors. The node on becoming a cluster head wake up for a longer wake up period and less sleep time than usual node in order to detect a new incoming node. On receiving a SYNC packet from the new incoming node not in the list, the CH quickly broadcasts its schedule in order to let the newly incoming node to synchronize with the CH. In this way the energy of all other nodes can be conserved by waking up and sleeping at the CH schedule and the role of CH is exchanged regularly in order not to let the energy of CH completely drain. We introduce a cluster based approach in which a cluster is formed and each member node then synchronizes with the schedule of Cluster Head (CH). The member nodes are scheduled to wake up at the same time as that of the cluster head. The amount of remaining energy in each node is used to choose a cluster head. The listen period is further divided into SYNC and DATA packets, respectively.

On receiving the SYNC packet each node maintains a list of immediate neighbors and is updated periodically. The cluster head election takes place after every 10 frames (10s) in order to maintain the synchronization. Clusters are formed depending on the mobility of each node. Each node calculates the relative velocity among its immediate neighbors.

ALGORITHM 1. PROPOSED EEMA-MAC

```
BEGIN
1. SET N (No. of Nodes)
2. SET radius (communication radius)
3. Initialize state
4. for I = 1 to N
    a. WHILE (1)
        i. set state = 'M'; %each node starts as a member
        ii. Each Node sends a SYNC packet
        iii. Waits for a time period T.
        iv. IF node hears a PACKET
            1. Classify Packet type
            2. IF PACKET Type == 'SYNC'
                a. Create a table for each neighbors with similar relative velocity
                b. Compare the remaining energy of each node
                c. If Erem_own > Erem_table
                    i. CH = i;
                    ii. State = 'CH';
                    iii. BROADCASTS CLUSTER_ADV_MSG to its neighboring nodes
                d. END
            3. else if PACKET type = 'CLUSTER ADV MSG'
                a. Node I schedules with the wakeup time of CH.
                b. CH = sender id;
                c. State = 'M';
            4. else if
                a. TIMEOUT
                    i. SELECTS a random sleeping time and then broadcasts schedule
                    ii. GO TO STEP 4.
                b. else
                    i. GO TO step a.
                c. END
            5. END
        b. END
    b. END
END
```

Hence the nodes with similar mobility forms a cluster and the node with highest remaining energy is elected as a cluster head, which is turn broadcasts the CLUSTER_ADV_MSG to all its immediate neighbors within its range. On receiving CLUSTER_ADV_MSG each node now schedule with the sleep time of the cluster head.

The working of EEMA- MAC protocol is explained in the form of an algorithm which can be seen in Algorithm 1 above.

A. RSSI based location estimation

Each node estimates its location depending upon the Received Signal Strength, every time it receives SYNC packet from neighboring nodes. The SYNC packet consists of the estimated x and y coordinate of the sender. The receiver node then calculates the angle between the sender node and itself depending upon these values. In our proposed algorithm we assume that a node requires at least 3 RSSI values to estimate the location. The path loss is assumed as

$$pl(d) = pl(d_0) + 10\gamma \log_{10}\left(\frac{d}{d_0}\right) \quad (1)$$

Where pl represents path loss and d_0 is the reference distance.

$$pl = 10 * \log_{10}(P_t / P_r) \quad (2)$$

Here P_t signifies transmit power and P_r indicates received signal strength.

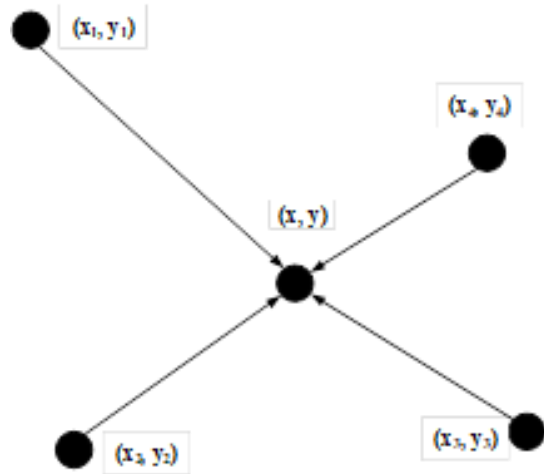


Fig. 2. Node topology.

Consider a node at location (x, y) whose location is to be estimated as shown in Fig. 2. It receives a SYNC packet from n neighboring nodes which consists of x -coordinate and y -coordinate. A target node estimates its location (x, y) and calculates angle between itself and the neighboring nodes which help to find the direction of movement of each node. Let D_i be the estimated distance between the target node and the node (x_i, y_i) which is obtained using the log-model.

$$D_i = \sqrt{(x_i - x)^2 + (y_i - y)^2} \quad (3)$$

$$D_i^2 = (x_i - x)^2 + (y_i - y)^2 \quad (4)$$

From above equation we can write,

$$D_i^2 = x_i^2 - 2x_ix + x^2 + y_i^2 - 2y_iy + y^2 \quad (5)$$

Here, $i = 1, 2, 3, 4, \dots, n$

Similarly, for n-th node the equation can be written as

$$D_n^2 = x_n^2 - 2x_nx + x^2 + y_n^2 - 2y_ny + y^2 \quad (6)$$

From (5) and (6) it can be written as

$$x_i^2 + y_i^2 - x_n^2 - y_n^2 + D_n^2 - D_i^2 = 2x(x_i - x_n) + 2y(y_i - y_n) \quad (7)$$

$$B = \begin{bmatrix} x_1^2 + y_1^2 - x_n^2 - y_n^2 - D_1^2 + D_n^2 \\ x_2^2 + y_2^2 - x_n^2 - y_n^2 - D_2^2 + D_n^2 \\ \dots \\ x_{n-1}^2 + y_{n-1}^2 - x_n^2 - y_n^2 - D_{n-1}^2 + D_n^2 \end{bmatrix} \quad (8)$$

$$A = \begin{bmatrix} 2(x_1 - x_n) + 2(y_1 - y_n) \\ 2(x_2 - x_n) + 2(y_2 - y_n) \\ \dots \\ 2(x_{n-1} - x_n) + 2(y_{n-1} - y_n) \end{bmatrix}$$

$$X = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$AX = B$$

Three nodes are enough to locate the coordinates of unknown nodes and coordinates are calculated using trilateral localization algorithm.

After the estimation of coordinates, the angle is given by

$$\theta = \tan^{-1}\left(\frac{y_i - y}{x_i - x}\right) \quad (9)$$

B. Cluster Head Election

When a node enters a sensing area it sends a SYNC packet and wakeup for a time period t. If the node does not listens to any schedule from its neighboring nodes, it randomly selects its sleeping time and broadcasts its schedule in a SYNC message signifying that it will go to sleep after certain time period.

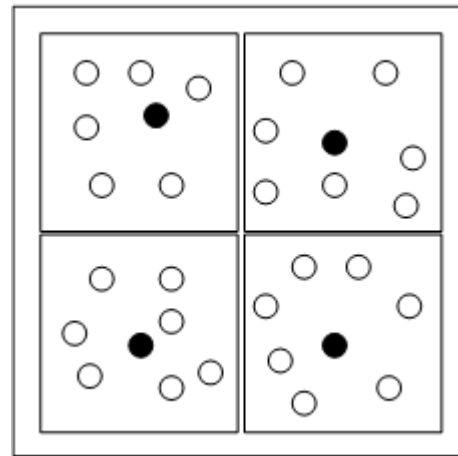


Fig. 3. Cluster formation model.

In Fig. 3, the nodes with black color are the CHs. The clustering algorithm forms a virtual cluster of nodes with similar velocity. The member nodes now synchronize with the CH. If a node receives a CLUSTER_ADV_MSG from other CH, it then synchronizes with schedule of the node from which it has received a broadcast message.

The cluster head election takes place after every 10 frames in which all nodes have equal probability to become a cluster head. Each node maintains a table of their neighboring nodes with similar mobility on receiving of SYNC packet. The relative velocity is now given by

$$\Delta v = v_A - v_B$$

Hence this helps to form a group of clusters with same mobility. Based on the information received from SYNC packet decides whether a node can become a cluster head or not. If the node has maximum energy among the group of nodes in a neighboring table, then the node immediately broadcasts CLUSTER_ADV_MSG with its own id and sleep schedule to its neighboring nodes. On hearing of this message all the neighboring nodes now synchronizes with the CH.

For the border nodes which receive a broadcast message from more than one CH chooses the node which is closer to itself and then discards the other.

C. SYNC Packet

The format of SYNC frame is shown in Fig. 4. The member nodes have an equal sleep and wake up time as that of S-MAC as can be seen in Fig. 5 whereas the cluster head has an extended wake up time and less sleep time.

Fig. 6 shows longer wakeup time which allows the newly coming node to synchronize with the CH. On receiving of SYNC message from a newly coming node immediately responds with an ACK packet by the CH. Only the cluster head acknowledges with an ACK packet.

Type	Length	Src address	Sync node	Sleep Time	Speed	x,y coordinate	Remaining Energy	State	CRC
------	--------	-------------	-----------	------------	-------	----------------	------------------	-------	-----

Fig. 4. Format of SYNC frame.

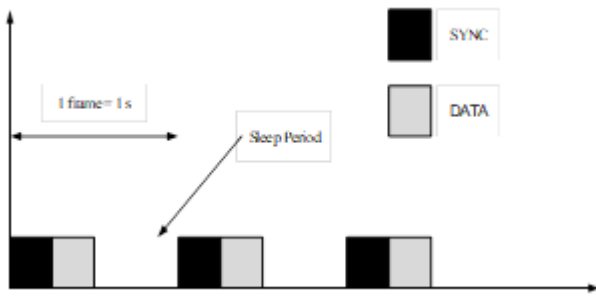


Fig. 5. Description of EEMA-MAC synchronization period.



Fig. 6. Description of EEMA-MAC synchronization period for CH.

TABLE I. SIMULATION PARAMETERS

Area	[100m x 100m]
No of Nodes	100
Speed	[1m/s – 5m/s]
Communication radius	30m
Protocols	[SMAC,MS-MAC, EEMA-MAC]
CWmax	100
Duty cycle	10%
Sleep Power	0.05mW
Idle Power	344.2mW
Transmit Power	386mW
Receive Power	368 mW
Total Energy	1000J
Packet Size	50 bytes

III. PERFORMANCE EVALUATIONS

EEMA-MAC is compared with Sensor-Medium Access Control (S-MAC) and Mobility-aware Sensor MAC (MS-MAC) protocols in order to evaluate its performance and efficiency. Simulation results conclude that EEMA-MAC enhances the overall throughput. Moreover it decreases energy consumption and packet loss compared with existing protocols.

A. Simulation Parameters

Nodes are deployed randomly over 100 m × 100 m region. Number of nodes for this simulation is fixed to 100 (n = 100). Fig. 4 shows the randomly deployed nodes for our evaluations.

The parameters given in Table I are used to simulate S-MAC [1], MS-MAC [9] and EEMA-MAC.

B. Simulation Results

To achieve the desired results, MATLAB is used as a simulation tool. Firstly, we need to analyze the energy

consumption of all the protocols. To calculate the average energy consumption we use the formula:

$$Average_Energy_Consumption = \frac{\sum_{i=1}^N E - E_{rem}}{N}$$

The speed of the nodes is varied from 1 m/sec to 5 m/sec. Fig. 7 shows that with increase of speed, the average energy consumption also increases. The energy consumption of EEMA-MAC protocol has less energy consumption as compared to other protocols. The energy consumption is highest for SMAC protocol which shows that it is highly inefficient in mobile environments since the nodes consume more energy in maintaining the schedules.

Secondly, we have to calculate the throughput which is termed as the total number of packets transferred over total simulation time. With the increase of speed, the throughput decreases. Since the increase in speed causes less packets to reach to the destination. The throughput for EEMA-MAC is the highest whereas the SMAC protocol shows lowest throughput as shown in Fig. 8.

With the increase of speed, quantity of dropped packets also increases. Thus, performance at a node cannot be just measured in terms of delay, but also in terms of probability of the packet being dropped. In order to make sure that all the data is eventually transported from source to destination, a retransmission of dropped packet can be done on end to end basis. Also, the losses between 5% and 10% of the total packet stream will disturb the network performance significantly. The higher is the speed, the higher is the number of losses. The packet drop is seen to be much less in EEMA-MAC protocol then SMAC and MS-MAC protocol as can be seen in Fig. 9. We can calculate the Packet drop by using the formula given below:

$$PD(\%) = \frac{Total.no.of.packets.sent - Total.no.of.packets.received}{Total.no.of.packets.sent} \times 100$$

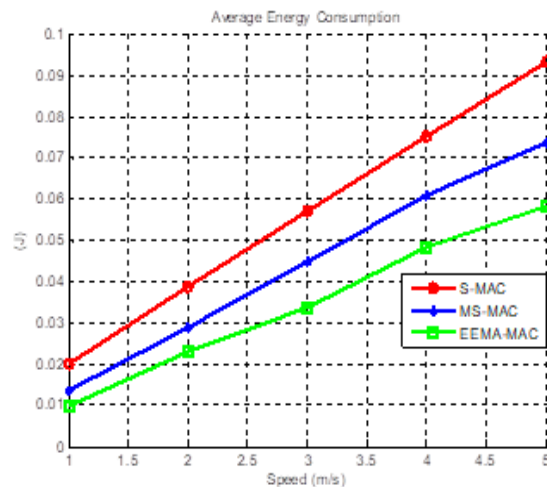


Fig. 7. Average energy consumption.

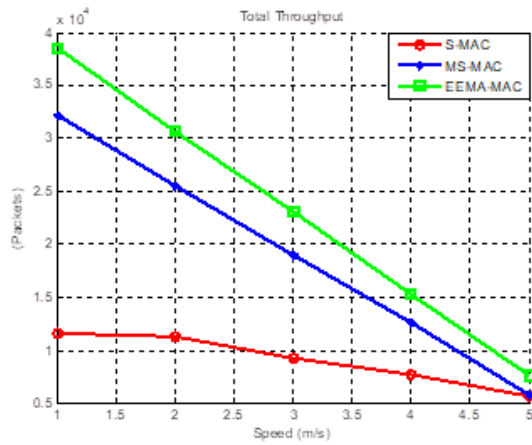


Fig. 8. Total throughput.

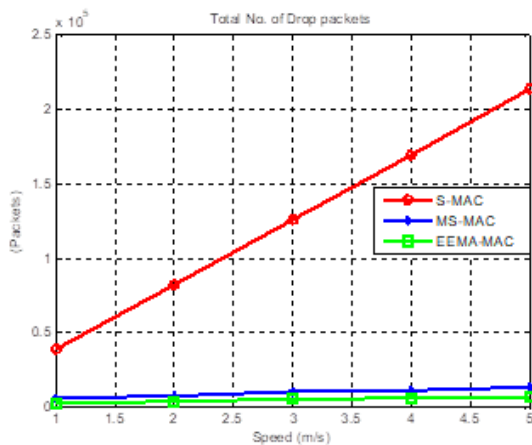


Fig. 9. Total number of drop packets.

IV. CONCLUSION

A number of potential applications (areas) are on the rise in the field of Wireless Sensor Network. The energy efficient MAC protocols turn out to be enormously significant and are indeed subject to many present research projects. In our current research, we have proposed an energy efficient MAC protocol called EEMA-MAC for WSNs. The major aim is to increase the efficiency of the network by extending the wake up time and reducing the sleep time at CH end. Performance evaluation

and simulation results conclude that EEMA-MAC has made a significant improvement in terms of energy consumption, throughput and packet loss in comparison with previously proposed MAC protocols such as S-MAC and MS-MAC. In the future, we can further evaluate the effectiveness of our proposed protocol by varying the number of nodes and can conduct more simulations involving random topologies.

REFERENCES

- [1] W. Ye, J. Heidemann, and D. Estrin, "An energy-efficient mac protocol for wireless sensor networks," IEEE Twenty-First Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings, vol. 3, pp. 1567 – 1576, 2002.
- [2] F. Peng, and M. Cui, "An Energy-Efficient Mobility-Supporting MAC protocol in Wireless Sensor Networks," Journal of Communications and Networks, vol. 17, no. 2, 2015.
- [3] F. Peng, "A novel Adaptive Mobility-Aware MAC protocol in Wireless Sensor Networks," Springer Wireless Personal Communications Journal, vol. 81, no. 2, pp. 489-50, 2015.
- [4] W. Dargie and C. Poellabauer, Fundamentals of Wireless Sensor Networks: Theory and Practice. Wiley Publishing, 2010.
- [5] A. Mainwaring, D. Culler, J. Polastre, R. Szewczyk, and J. Anderson, "Wireless sensor networks for habitat monitoring," in ACM International Workshop on Wireless Sensor Networks and Applications (WSNA 2002), 2002, pp. 88–97.
- [6] S. Kim, S. Pakzad, D. Culler, J. Demmel, G. Fenves, S. Glaser, and M. Turon, "Health monitoring of civil infrastructures using wireless sensor networks," in IPSN '07: Proc. 6th international conference on Information processing in sensor networks. New York, NY, USA: ACM, 2007, pp. 254–263.
- [7] G. Werner-Allen, K. Lorincz, M. Welsh, O. Marcillo, J. Johnson, M. Ruiz, and J. Lees, "Deploying a wireless sensor network on an active volcano," IEEE Internet Computing, vol. 10, no. 2, pp. 18–25, 2006.
- [8] I. Stoianov, L. Nachman, S. Madden, and T. Tokmouline, "PIPET: a wireless sensor network for pipeline monitoring," in IPSN '07: Proc. 6th international conference on Information processing in sensor networks. New York, NY, USA: ACM, 2007, pp. 264–273.
- [9] H. Pham and S. Jha, "An adaptive mobility-aware mac protocol for sensor networks (ms-mac)," IEEE International Conference on Mobile Ad-hoc and Sensor Systems, pp. 558 – 560, October 2004.
- [10] S. W. A. Kazmi, A. Kacso, and R. Wismuller, "Recent MAC Protocols for Mobility-Aware Wireless Sensor Networks - A Survey and Future Directions," IEEE Ninth International Conference on Ubiquitous and Future Networks (ICUFN), July 2017.
- [11] H. R. Silva, J. S. Silva, and F. Boavida, "Mobility in wireless sensor networks—survey and proposal," Computer Communications, vol. 52, pp. 1–20, 2014.
- [12] S. W. A. Kazmi, A. Kacso, and R. Wismüller, "On MAC design for mobility-aware wireless sensor networks," IEEE 2nd International Conference on Computer and Communication Systems (ICCCS), pp. 152-157, July 2017.

Routing Optimization in WBAN using Bees Algorithm for Overcrowded Hajj Environment

Ghassan Ahmed Ali*, Shah Murtaza Rashid Al Masud
College of Computer Sciences and Information systems,
Najran University, Najran, Kingdom of Saudi Arabia

Abstract—Crowded places like Hajj environment in Makkah which host from 2 to 3 million on specific area and time can pose health challenges for pilgrims who need medical care. One of the solutions to overcome such difficulties is to use Wireless Body Area Networks (WBANs). WBAN is one of the new technology using wireless sensor network to gather data about status of patient then to forward collected data to be proceeded. However, various types of challenges in WBAN should be concerned. Power consumption is critical within WBAN system. Furthermore, delay of data transfer may lead to wrong diagnosis or uncorrected report that may lead to death; therefore, the transferred data must be reliable to ensure accuracy in measurement. In this paper, we propose a framework for routing optimization in medical wireless network. The proposed framework optimize shortest path in different stages of collected data to get less energy consumption, and reduce transmission time. The proposed work is based on Bees Algorithm to overcome such challenges and find shortest path for data within shortest time during overcrowded of Hajj environment. Matlab simulation results show good performance of Bees Algorithm in terms of reducing transmission time, energy consumption, delay, and throughput.

Keywords—Wireless Body Area Network (WBAN); Bees algorithm; routing optimization; Hajj environment

I. INTRODUCTION

Bees Algorithm is inspired by the foraging behavior of honeybee in nature to find the best solution of a given optimization problem proposed by [1]. According to [2], Bees Algorithm approved its higher performance compared to many heuristic optimization in many problems. As a result of its simplicity, Bees Algorithm applied in many application including solving examination timetabling problems [3], training neural networks [4], job scheduling for machine [5], supply chain optimization [6], data clustering [7], correlation-aware service in cloud [8], and robot path in dynamic environment [9]. It is obvious that the study of Bees Algorithm and its applications in the literature increases exponentially. The Bees Algorithm shows its efficiency in terms of speed, learning, and accuracy.

In this paper, Bees Algorithm is used as optimization method to assist in deployment of WBAN and make transmission of WBAN more efficient during Hajj. The Hajj is an annual pilgrimage to Makkah in Saudi Arabia. It is performing once-in-a-lifetime obligation for all Muslims who have the ability to undertake the journey. Two to three millions pilgrims from different countries are gathering and the number of pilgrims is increasing every year. One of the

most challenges is that the complete events of Hajj must be performed in a specific locations not exceed 4 km² within specific 5 days. Therefore, places become over-crowded and difficult to be reached by medical emergency in case of injuries. However, there are some public's health care centers surrounding Makkah, but because most of pilgrims are not familiar with the place as the trip to Makkah often is first time to majority of pilgrims, thus it becomes very difficult to reach the health centers.

The review study of [10] reported that during Hajj many pilgrims' are suffering from various infectious and chronic diseases. The major health problems encountered by pilgrims are respiratory diseases (73.33%), heat effects (16.67%), diabetes (13.32%), cardiovascular (10%), gastroenteritis (10%), hypertension (6.67%), and urinary tract infection including trauma (3.33%). Whereas study of [11] observed that cardiovascular diseases is the main reason of death during Hajj.

In the last few years, new technologies have been proposed to overcome medical challenges and provide healthcare services like real-time monitoring, observing health status, managing diseases, and remote connecting to hospitals. Wireless body area network (WBAN) is one of the most promising technology that enables monitoring of health conditions, disease diagnosis, and real-time observing. The WBAN is designed as a sensor network located in patient body to collect patient medical information, and then send information to the coordinator. The coordinating or monitoring sensors search for a suitable communication network to local server to store data and then communicate with remote database server for diagnosis purpose. Until now, WBAN has been deployed for in-door based medical applications in hospitals and clinics. Hence, it is essential to ensure some of the vital requirements of quality of service of WBAN e.g. low energy consumption, higher throughput, lower delay and no collision while deploying at out-door based healthcare purposes especially for overcrowded Hajj environment. Moreover, quality of service, sensitivity of patient's information, and short time during data transmission are critical. The power consumption of WBAN keeps increased due to the growth of data rate and data transmission distance and resulting in decrement of network life.

According to [12], in crowded places like Hajj, medical data transmission must be seamless and reliable by using multi-hop based routing which provides low power consumption and consistent data routing of wireless communication. Selecting the shortest path for transferring

data plays a key role for quality of services in terms of the power consumption and delay.

II. RELATED WORKS

Recently, population-based algorithms such as Ant Colony optimization (ACO), Genetic Algorithm (GA), and Particle Swarm Optimization (PSO) have been used to improve routing algorithm in WBAN; for example, PSO algorithm presented in [12] to search for optimal location of the relay node to improve radio frequency energy in WBAN. Authors in [13] used GA to optimize extracted features in WBAN in terms of latency, classification rate and packet delivery rate. It has been showed that GA optimization algorithm is effective in sensor classification. Furthermore, ACO is proposed in [14] to find shortest route through sensor node. The distance route is calculated from patient to medical center then the shortest node is selected based on Bayesian game formulation. According to [15], Bees Algorithm has been shown to be powerful optimization methods when compare it to other population-based methods [16]. The Bees Algorithm is proposed in this paper to find best path for data to reach destination within shortest time during overcrowded Hajj environment.

III. BEES ALGORITHM OPTIMIZATION

Bees Algorithm is a population-based search algorithm inspired from nature of honeybees to find an optimal solution. Basic Bees Algorithm is divided into four components: parameter setting, initialization, local search, and global search. Bees Algorithm in its basic form uses a set of parameters need to be set for the algorithm as shown in Table I.

Fig. 1 shows pseudo code of the Bees Algorithm. The Bees Algorithm in Step 1 start generating n scout randomly as initial population. Then in Step 2, the fitness is evaluated of sites explored by scout bees. The "elite bees" are selected based on the highest fitness and neighborhood search are chosen in Step 4. In Step 5, the selected sites are recruited and more bees are employed for the elite sites as well as the fitness is evaluated. In step 6, the fittest bee is selected to produce the next bee population. The remaining bees are then assigned randomly to seek for new solutions around the domain in Step 7.

TABLE I. PARAMETERS OF BEES ALGORITHM

Ns	Scout bees
Ne	Elite sites
Nb	Best sites
Nre	Recruited bees of ne
Nrb	Recruited bees of remaining nb
Ngh	Neighborhood initial size

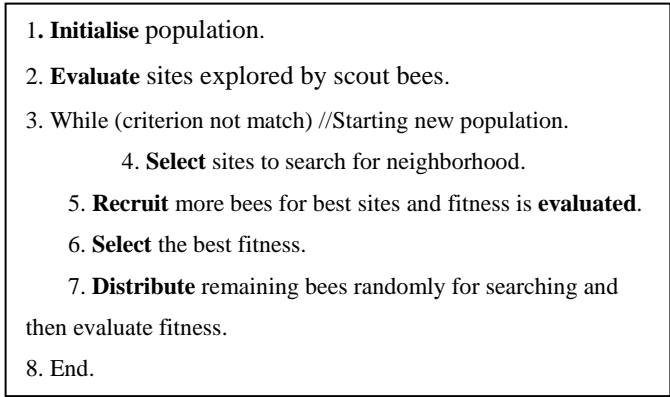


Fig. 1. Bees algorithm pseudo code.

IV. BEES ALGORITHM OPTIMIZATION IN WBAN

The proposed method has following steps in each iteration:

- Seek for a node randomly in search place.
- Test the condition of the node.
- Choose nodes for neighborhood search.
- Select best nodes with shortest time and evaluate fitness.
- Search randomly and keep evaluating fitness.
- Determine to continue searching or terminate the iteration.

Fig. 2 shows the proposed Bees Algorithm in WBAN.

WBANs lead to degradation of performance e.g. energy consumption, delay, throughput and collision, when the number of sensor density increase; interaction between the sensors with WBANs or interaction between the WBANs in the same or different environment increases; and the distance between WBANs and gateway increases. Hence, deploying WBANs for pilgrims health monitoring during Hajj requires extra attention. Pilgrims healthcare facility using WBAN as proposed in this paper is a dynamic procedure because the pilgrims may sometimes in motion walking, running, and sitting.

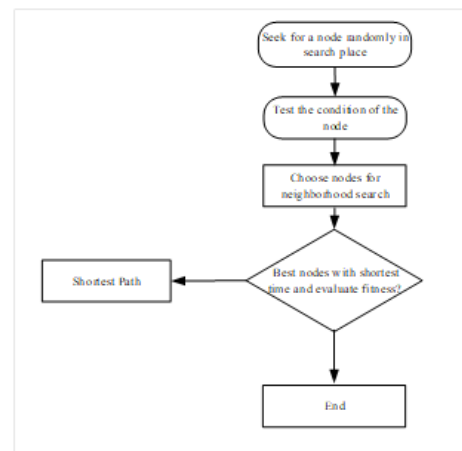


Fig. 2. Flow chart of proposed Bees Algorithm in WBAN.

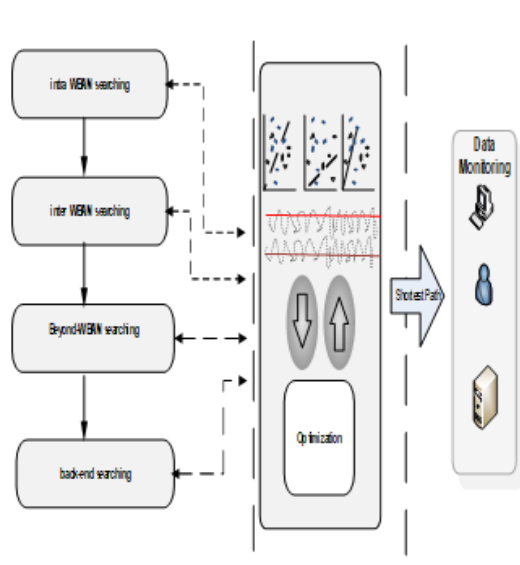


Fig. 3. Multi phases optimal search.

According to the IEEE 802.15.6 working group, nodes in WBAN are directly connected to the sink node (usually it is called as intra-WBAN and inter-WBAN) using one-hop or two-hop of star topology and thereafter the sink node which is called as coordinator is connected to access points (APs) in outer-WBAN and then to the personal server in back-end healthcare centers in a multi-hop architecture as shown in Fig. 3.

Intra-WBAN is a small network around the body can support data transmission until 1-2 meters, in some cases 2-5 meters and usually use various short-range communication infrastructure e.g. ZigBee, Bluetooth, Wi-Fi or Cellular network. And Internet is used for long-range communication to support data transmission in beyond-WBAN and back-end medical server for healthcare applications as presented in Fig. 4. Searching is performed not only in one stage but in multi phases of wireless network: intra-WBAN searching, inter-WBAN searching, beyond-WBAN searching, and back-end searching.

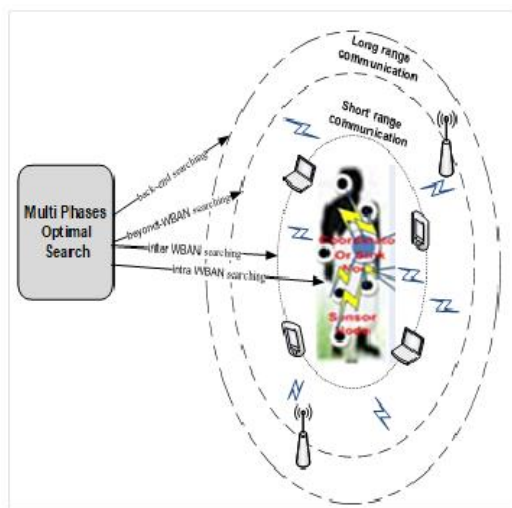


Fig. 4. Optimize long-range and short-range communications.

V. RESULTS AND DISCUSSION

The effectiveness of the proposed Bees Algorithm in WBAN is validated using Matlab simulator in the aspects of transmission time, energy consumption, throughput, and delay as following sub-section.

Define abbreviations and acronyms the first time they are used in the text, even after they have been defined in the abstract. Abbreviations such as IEEE, SI, MKS, CGS, sc, dc, and rms do not have to be defined. Do not use abbreviations in the title or heads unless they are unavoidable.

A. Analyzing the Best Effort (the Minimum Distance/ Shortest Path or Route) from Source to Destination. the Best Cost is Measured as the Minimum Time Required for Transferring Data from Source to Destination.

- **Parameters:** 1000 nodes, 30 meters distance, 20 iteration, unit of time is second.
- **Description:** Using these parameters, we will generate the following two graph. Fig. 5 describes the Iteration vs best cost using bees algorithm. Here we can see that the cost will be minimized if the number of iteration is increased. Fig. 6 describes the Iteration vs best cost without bees algorithm. From this graph, we can see that the cost is higher than the cost we got using bees algorithm.

The estimated data transmission time at the first five iteration with bees algorithm is presented in Fig. 7. It is shown that the transmission time decreased as the number of iteration increased. In Fig. 8, the same procedure has been followed but without applying the bees algorithms. From the Fig. 7 and 8 it is shown that the data transmission time at best cost by applying bees algorithm is less than that of without applying the bees algorithm.

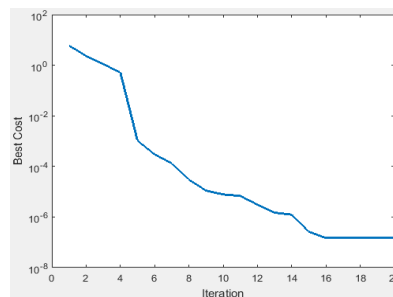


Fig. 5. Cost with applying Bees Algorithm.

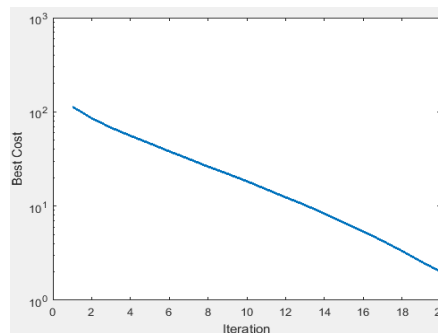


Fig. 6. Cost without applying Bees Algorithm.

```
Iteration 1: Best Cost = 6.0237
Iteration 2: Best Cost = 2.3119
Iteration 3: Best Cost = 1.1101
Iteration 4: Best Cost = 0.51486
Iteration 5: Best Cost = 0.0010406
```

Fig. 7. First 5 iteration with Bees Algorithm.

```
Iteration 1: Best Cost = 116.1093
Iteration 2: Best Cost = 86.5102
Iteration 3: Best Cost = 69.5346
Iteration 4: Best Cost = 57.4995
Iteration 5: Best Cost = 47.259
```

Fig. 8. First 5 iteration without Bees Algorithm.

B. Analyzing the Total Amount of Energy Consumption During Packet Transmission in Minimum or Shortest Distance.

- **Parameters:** 1000 nodes, 30 meters distance, 20 iteration, consider packet not bit and followed by the first analytical result.
- **Description:** Using these parameters, we will generate the following two graphs. Fig. 9 describes the total energy consumed during best cost by applying Bees Algorithm. Here we can see that the energy consumption will be minimized if the number of iteration is increased.

Fig. 10 describes the total energy consumed during best cost without applying Bees Algorithm. From this graph, we can see that the energy consumption is higher than the total energy consumption we got using bees algorithm.

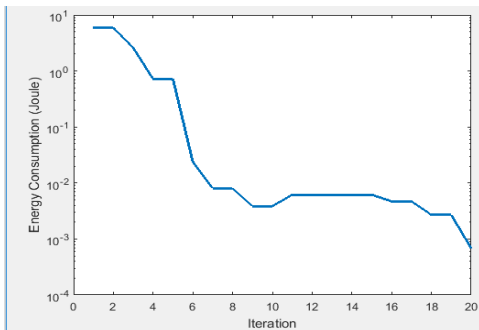


Fig. 9. Energy consumption with Applying Bees Algorithm.

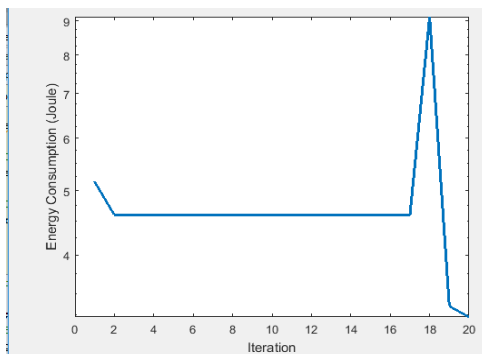


Fig. 10. Energy consumption without applying Bees Algorithm.

C. Analyzing the Delay

- **Parameters:** 1000 nodes, 30 meters distance, 20 iteration, less than 250 ms, packet size 50-300 bytes.
- **Description:** Using these parameters, we will generate the following two graphs. Both the graphs show Delay vs iteration where the unit of delay is second. Fig. 11 describes the Iteration vs delay using Bees Algorithm. Here we can see that the delay will be minimized if the number of iteration is increased. Fig. 12 describes the Iteration vs delay without Bees Algorithm. From this graph, we can see that the delay is higher than the delay we got using the Bees Algorithm.

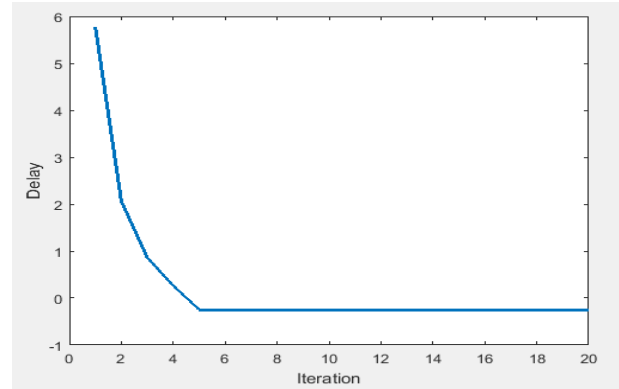


Fig. 11. Delay with applying Bees Algorithm.

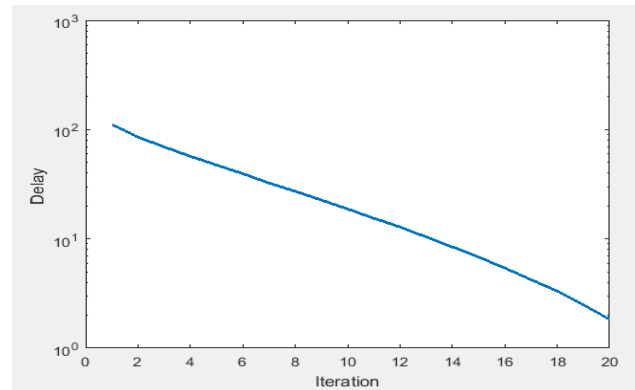


Fig. 12. Delay without applying Bees Algorithm.

D. Analyzing Throughput

- **Parameters:** 1000 nodes, 30 meters distance, 20 iteration, throughput more than 250 kbps, packet size 300 bytes.
- **Description:** Using these parameters, we will generate the following two graphs. Both the graphs show Throughput vs iteration where the unit of throughput is kilobits per second (Kbps). Figure 13 describes the Iteration vs Throughput using bees algorithm. Here we can see that the throughput will be increased if the number of iteration is increased.

Fig. 14 describes the Iteration vs Throughput without bees algorithm. From this graph, we can see that the throughput is lower than the throughput we got using bees algorithm.

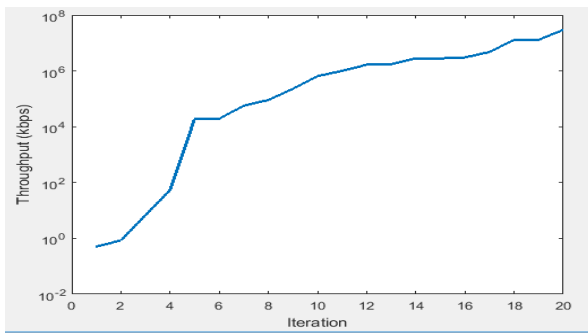


Fig. 13. Throughput with applying Bees Algorithm.

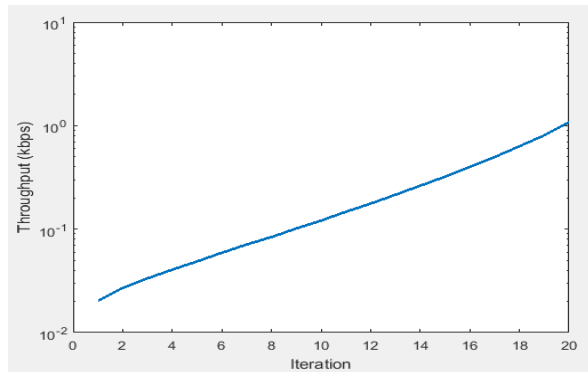


Fig. 14. Throughput without applying Bees Algorithm.

VI. CONCLUSIONS

This paper has presented Bees Algorithm in new field of WBAN. The Bees Algorithm has briefly outlined and compared to other works. Moreover, Bees Algorithm as an optimization tool used to select shortest path in multi phases that makes data reach its destination in a shortest time with low energy consumption and less delay. Simulation results show the effective of the proposed Bees Algorithm in WBAN which is promised to be very helpful for pilgrim to overcome many challenges during Hajj.

ACKNOWLEDGMENT

This research is supported by Najran University, Najran, Kingdom of Saudi Arabia, under Research Project Code: NU/ESCI/15/028.

REFERENCES

[1] Pham, D. T., Ghanbarzadeh, A., Koç, E., Otri, S., Rahim, S., & Zaidi, M. (2006). -The Bees Algorithm—A Novel Tool for Complex Optimisation Problems. In *Intelligent Production Machines and Systems*

(pp. 454-459).S.R.; and Jenkins, J.E. (1982). Evaluation of component buildup methods for missile aerodynamic prediction. *Journal of Spacecraft and Rocket*, 19(6), 481-488.

[2] Pham, D. T., & Castellani, M. (2009). The bees algorithm: modelling foraging behaviour to solve continuous optimization problems. *Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science*, 223(12), 2919-2938.

[3] Abdullah, S., & Alzaqebah, M. (2013). A hybrid self-adaptive bees algorithm for examination timetabling problems. *Applied Soft Computing*, 13(8), 3608-3620.

[4] Ali, G. A., & Jantan, A. (2011). A new approach based on honeybee to improve intrusion detection system using neural network and bees algorithm. In *International Conference on Software Engineering and Computer Systems*(pp. 777-792). Springer, Berlin, Heidelberg.

[5] Pham, D. T., Koc, E., Lee, J. Y., & Phruksanant, J. (2007). Using the bees algorithm to schedule jobs for a machine. In *Proc eighth international conference on laser metrology, CMM and machine tool performance, LAMDAMAP, Euspen, UK, Cardiff* (pp. 430-439).

[6] Yuce, B., Mastrocinque, E., Lambiase, A., Packianather, M. S., & Pham, D. T. (2014). A multi-objective supply chain optimisation using enhanced Bees Algorithm with adaptive neighbourhood search and site abandonment strategy. *Swarm and Evolutionary Computation*, 18, 71-82.

[7] Pham, D. T., Otri, S., Afify, A., Mahmuddin, M., & Al-Jabbouli, H. (2007). Data clustering using the bees algorithm. In *40th CIRP International Manufacturing Systems Seminar*.

[8] Xu, W., Tian, S., Liu, Q., Xie, Y., Zhou, Z., & Pham, D. T. (2016). An improved discrete bees algorithm for correlation-aware service aggregation optimization in cloud manufacturing. *The International Journal of Advanced Manufacturing Technology*, 84(1-4), 17-28.

[9] Haj Darwish, A., Joukhadar, A., & Kashkash, M. (2018). Using the bees algorithm for wheeled mobile robot path planning in an indoor dynamic environment. *Cogent Engineering*, (just-accepted), 1426539.

[10] Al Masud, S. M. R., Bakar, A. A., & Yussof, S. (2016). Determining the Types of Diseases and Emergency Issues in Pilgrims During Hajj: A Literature Review. *statistics and information*, 5(6), 7.

[11] Al Shimemeri, A. (2012). Cardiovascular disease in Hajj pilgrims. *Journal of the Saudi Heart Association*, 24(2), 123-127.

[12] Dhaou, I. B. (2010). Client-server network architecture for safe pilgrim journey in the Kingdom of Saudi Arabia. In *Intelligent Vehicles Symposium (IV)*, 2010 IEEE (pp. 1043-1048). IEEE.

[13] Wu, T. Y., & Lin, C. H. (2015). Low-SAR path discovery by particle swarm optimization algorithm in wireless body area networks. *IEEE Sensors Journal*, 15(2), 928-936.

[14] Kalaiselvi, K., Suresh, G. R., & Ravi, V. (2018). Genetic algorithm based sensor node classifications in wireless body area networks (WBAN). *Cluster Computing*, 1-7.

[15] Latha, R., Vetrivelan, P., & Jagannath, M. (2017). Balancing emergency message dissemination and network lifetime in wireless body area network using ant colony optimization and Bayesian game formulation. *Informatics in Medicine Unlocked*, 8, 60-65.

[16] Nasrinpour, H. R., Bavani, A. M., & Teshnehlab, M. (2017). Grouped Bees Algorithm: A Grouped Version of the Bees Algorithm. *Computers*, 6(1), 5.

An Opportunistic Dissemination Protocol for VANETs

Amina SEDJELMACI

Dept. of Telecommunications
University of Tlemcen, Algeria

Fedoua DIDI

Dept. of Computer Engineering
University of Tlemcen, Algeria

Ahmed ABDUL RAHUMAN

Dept. of Electronic Engineering
University of Tlemcen, Algeria

Abstract—In this article, we propose an Opportunistic information dissemination protocol by mixing both flooding and an enhanced DHVN (Dissemination protocol for heterogeneous Cooperative Vehicular Network) protocol, allowing them to run opportunistically in a Manhattan plan. Special additional logic is added to the existing version of DHVN protocol in order to efficiently disseminate information in two steps: 1) by adding three tags, Initial Diffusion, Standard DHVN and DHVN Near Intersection; the Initial Diffusion tag is used for the first flooding transmission only and 2) by changing the SNF (Store and Forward) period by making it adaptive depending on the region. Detailed simulation results show that our opportunistic protocol outperforms the DHVN protocol by analyzing its performances using an integrated framework VNS.

Keywords—Flooding; DHVN; SNF; opportunistic; VANET

I. INTRODUCTION

In recent years, the dissemination of data in VANET (Vehicular Ad hoc NETWORKS) has attracted a lot of attention given its imminent role in improving road safety and reducing traffic congestion. The issues weigh heavily on the economy of a country in terms of energy and time. To be able to reduce the risk of accidents, avoid dangerous situations and mitigate such problems, we have to succeed in effectively disseminating relevant information and spreading it as far as possible with a minimum bandwidth usage.

The way in which relevant information is broadcasted throughout the vehicle environment is considered as a most important aspect for the vehicles cooperation in VANETs. However, several problems can occur during this process of dissemination: (1) an excessive consumption of bandwidth in the case where we are confronted to an urban area. (2) A disconnected network problem can occur in the case of a rural area. These problems constitute a crucial challenge and the question that will arise to face them is “*what can we do to overcome the rural disconnection without excessive use of bandwidth and how can we avoid broadcast storms while keeping a high coverage ratio?*”

In this paper, we attempt to address these issues by suggesting an improved opportunistic DHVN protocol that we call oDHVN. The remainder of this paper is organized as follows. Section II describes related works about different dissemination protocols operation mode classes. Both simulation tool and scenarios are given in Section III. In Section IV, our opportunistic protocol is described. Simulation

results of the designed solution are presented and discussed in Section V. Finally, Section VI concludes the paper.

II. RELATED WORKS

An interesting aspect of VANETs is that most of their applications require an efficient and reliable multi-hop data broadcast protocol, making this task performed by the protocol one of the most difficult and indispensable network functions. For example, avoidance of vehicle collisions and post-crash warnings require efficient and robust data dissemination, especially when the distances between the sender and the intended receivers are greater than the radio transmission range [1]. We focused on approaches that focus on reducing bandwidth usage by limiting the number of re-transmissions with optimal selection of relays and transmission parameters based on network conditions.

A. Characteristics and Operation Mode Classes

Data dissemination approaches in VANETs can be classified into three main classes: Relayed Instantaneous Broadcast, broadcast with store-and-forward SNF relay, and opportunistic broadcast. In the Relayed Instantaneous Broadcast approaches, the data is broadcasted to all neighboring vehicles that will briefly store the data and then a neighbor vehicle is selected as a relay to rebroadcast it in turn. This approach works well in high density networks that avoid broadcast storms, but require efficient relay selection to ensure reliability. A good survey on this class can be found in [2].

Alternatively, in the broadcast with store-and-forward SNF relay approach, data is stored, transported and then rebroadcast on network partitions, making them more suitable for irregularly distributed vehicle density zones. In the opportunistic approach, the two previous approaches are combined to adapt according to the circumstances of the network (dense, disconnected, high occupancy rate of the channel, etc.).

Other classifications that are based on other functional aspects can be found in the literature. In [1], flooding is considered a full-fledged approach whereas it can be defined as a special case of Relayed Instantaneous Broadcast approaches since the data is disseminated to all neighboring vehicles which are all considered as potential relays for a single retransmission (there is no store and forward) to their neighboring vehicles.

In [3], the first two approaches are grouped into one approach creating two classes; stateless broadcasting and stateful broadcasting. In the first class, there is no need to

obtain information about the network topology while the second class, the protocol requires information on the local topology.

In [4], the authors distinguish two main categories: multi-hop broadcast and single-hop broadcast. The authors divide the first class into several subclasses according to the method of election of the relay (delay, probability) or according to the use of the method "network coding" [5]. In the single-hop broadcast, when receiving a message, the vehicle retains and updates the information in its embedded database and in turn broadcasts, periodically in its one-hop vicinity, its version of the information. Authors divide this class into two subclasses depending on whether the diffusion decision period is fixe or adaptive.

In a more recent work [6], the authors repeat the general multi and single hop classification but adopt another reasoning more in line with ours. In this last work, the focus is on the multi-hop class which includes the largest number of protocols where the authors consider two different categories: 1) restrictive methods, and 2) promiscuous methods. As for our classification's first class, restrictive methods combine techniques to address the problems of broadcast storms. The difference from previous work lies in the introduction of the promiscuous class subclass where VANETs can be fragmented and partitioned, hence the use of techniques such as Store-and-Forward to ensure that the information is correctly disseminated. The authors mention other approaches (not clearly classified) that combine two different techniques to improve dissemination performance. These correspond to our third class of opportunistic diffusion. Other classifications based on different points of view or spanning other higher spheres (Security, QoS, Encryption, Topology, etc.) can be found in [7], [8].

1) *Relayed Instantaneous Broadcast class*

One of the representative protocols of this class is the Distance Defer Transfer Protocol (DDT) [9]. In DTT, upon receipt of a new message, the vehicle triggers a timer that is inversely proportional to the distance to the transmitter. During this waiting time, the vehicle records the positions of all vehicles that transmit the same message and decides to abandon the retransmission if most of its own retransmission area has been sufficiently covered by its neighbors. Otherwise the vehicle retransmits the message by applying the same protocol.

Another representative protocol of the first class can be found in [10] where the original transmitter simply accesses the medium using the standard 802.11 CSMA / CA technique and broadcasts the entire emergency message. All neighboring vehicles in the transmission range calculate their corresponding SNR values as well as their Euclidean distance from the source via GPS. Subsequently, each receiver then uses these results to calculate the maximum size of the specific contention window (CW_{max}). Each node randomly chooses a time slot in the range [0, CW_{max}] and waits during this slot time. The node that chooses the shortest time interval becomes the relay and rebroadcast the emergency message. The rebroadcast message serves as an acknowledgment to the original sender.

In [11], the SIFT protocol is proposed and comprises two phases: 1) the trajectory calculation which is only executed by the source node before sending a new packet for the first time. This phase calculates the trajectories and sends the packets by triggering the multi-hop routing process. 2) The packet routing phase that is invoked by each intermediate node when receiving a packet. This phase allows the node to decide by triggering a timer according to its position with respect to the trajectory and the transmitter whether or not to transmit the packet.

In [12], the selection of the next relay is performed by the calculation of a probability by each receiver of an emergency message. The latter will determine a Backoff period (i.e. the waiting time before retransmitting the received message). The backoff duration is calculated according to the following formula: Where WT is the value of Backoff, CW is the contention window, P is the calculated probability.

$$WT = CW \times (1 - P) + \delta \quad (1)$$

In this way, the vehicle with the shortest waiting time will transmit the message first. The vehicle with the highest probability will have the shortest backoff period. This retransmission probability is a weighted sum of two parameters: the distance factor (D) and the link quality factor (LQ). It is calculated as follows:

$$P = (1 - \omega_p)D + \omega_p LQ \quad (2)$$

Where ω_p is a weight between 0.5 and 1. This is to give more importance to the quality of the link, since the security messages considered in this study are critical in nature and the reliability of the transmission is one of our main axes.

In [2], The authors take a new approach to the calculation of the waiting time and use ZigBee as a communication technology to eliminate the redundancy of broadcast messages and thus minimize the Storm Broadcast. Indeed, the waiting time is adjusted according to the distance of the vehicle from the base vehicle and the relative speed. If the vehicle speed is slow, its waiting time is increased. The vehicle that travels at the highest speed and has the largest number of nearby vehicles will have very little waiting time and will be broadcast instantly.

2) *Broadcast with store-and-forward SNF relay*

There are several protocols in the literature that belong to the second class, broadcast with SNF relay type. In (Cherif et al., 2010) [13], the ROD protocol is organized in two modules; 1) an ODDT module (Optimized Distance Defer Transfer): the same method as in the DDT is adopted where the GPS position of the vehicle is encoded in the header of the broadcast message. The ROD protocol encodes additional information in addition to the GPS_pos which represents the position of the sender, OI_pos represents the outbound intersection position and II_pos represents the incoming intersection position that will be used by the ODDT module to optimize the data dissemination, in sections of road (between two intersections) and in intersections. 2) The Store and Forward SNF module where if no relay vehicle is found due to temporary network

fragmentation, the vehicle in charge of the message uses the Store and Forward module to keep the data until a better retransmitter is found.

Another protocol representative of the broadcast class with store and forward relay is the DHVN protocol [3] which gives particular attention to the network connectivity, the road structure and the heterogeneity of the vehicles. In this protocol, the source car broadcasts the packet in both directions where each receiver on the same route triggers a timer based on the distance to the transmitter. It retrieves the sender's position information from the packet header and calculates the backoff delay as follows:

$$T = \frac{1}{dist+car_height*MD} \quad (3)$$

Where *dist* is the distance between sender and receiver, *Car_height* is the height of the vehicle, *MD* is the maximum extra distance when the vehicle is 1 meter higher than a standard vehicle. A relay is chosen for each route and each direction to propagate the message. Since vehicle networks are also highly partitioned networks, continuous connectivity is not guaranteed. To allow longer-term dissemination of information in the case of highly partitioned networks, the DHVN protocol uses the SNF approach where nodes carry information with their movement and transmit it periodically. In DHVN, the choice of the retransmission period is crucial. Indeed, a small period causes a loss of bandwidth and a high period implies a significant delay. The algorithm of DHVN protocol is summarized in Fig. 1.

```
While (Position is in Dissemination Area)
{
  function Receive (msg){
    if (same_road) {
      if (first reception) Trigger timer`
      elseif (duplicate && sender is before) Cancel
      timer`
    } //end same road
    if (Intersection Zone){
      if (first reception) Trigger timer`
      else { //if the message is already received
        if (duplicate && sender is not in the same
        road)
          Ignore the reception and continue to
          disseminate`
      }
    } //end IO_Zone
  } //end event receive
} //end while
Function Timerfired()
  Trigger timer with SNF period
```

Fig. 1. Original DHVN protocol algorithm.

The TrAD protocol [14] requires beaconing to maintain a list of vehicles and their status in a single-hop neighborhood in order to work seamlessly in both urban and rural scenarios. It is composed of two main components: 1) The flooding

suppression technique for a well-connected network which makes it possible to constrain the broadcast storm problem and improve the reliability of the transmission; and 2) the SCF mechanism (Store-Carry-Forward) that selects the appropriate vehicles to act as relay which save the message and rebroadcast it in a disconnected network. Several concepts of TrAD are defined as (a) Directional cluster: It is a group of vehicles in the neighborhood of a sender *S*, which are in a similar direction with respect to the sender *S*. (b) Coordinator: the coordinator is the vehicle that is located at an intersection. (c) Breaker: In a well-connected network, the breaker is not only the farthest vehicle but also the one that moves out of the network.

The first component contains two mechanisms: 1) Classification of vector-angle clusters which consists of designing several clusters according to the vector angle with the sender *S* to identify if the vehicles belong to a directional Cluster. 2) Traffic-adaptive sorting technique that takes account both for road traffic and network traffic status and for assigning the transmission task to the neighbor who has both a dense neighborhood, a greater distance from the sender and the lowest rate of transmission (occupation of the canal).

The second component also includes two mechanisms: 1) Selection of the SCF agent: this technique makes it possible to identify the breaker; when the vehicle receives a data message, the protocol checks and eliminates the possibility of being a coordinator. After that, the vehicle checks whether its direction of travel is the same as the direction of the data. If so, the vehicle will look for another neighbor even further, which also moves in the direction of data transmission. Otherwise, the vehicle is defined as a breaker. This procedure will go to the limit of the connected network. 2) SCF redundant redistribution technique: This technique is intended to trigger the re-broadcast of the SCF or restrict it if more than one SCF receives the same request. Thus, a different broadcast delay for each SCF is calculated based on the shortest distance of the sender and the lowest channel occupancy.

3) Opportunistic broadcast.

Rare are the protocols that belong to this class [15], [1]. Such protocols take advantage of the strengths of the two first classes and mitigate their weak points by combining them to propagate data via a vehicular ad hoc network (VANET). The authors in [15] propose a hybrid protocol that allows to merge the two classes by classifying the field relative to each vehicle as a so-called multi-hop broadcast MHB area or a so-called store and forward (SF) area. The MHB and SF regions are partitioned via a broadcast area radius (*R*) around the sender where the data is broadcasted via a multi-hop broadcast in the broadcast area. Outside the MHB area, the message spreads via the store and forward protocol.

III. METHODOLOGY

A. Simulation Tool

Two aspects of the VANET simulation exist, the first aspect concerns the simulation of the vehicles mobility and the second concerns the simulation of the network components. Mobility simulators have the ability to generate realistic vehicular mobility traces reflecting the movements of vehicles. These will be introduced as input for the network simulator.

There are several mobility simulation software environments and their main characteristics are: a) Supported trace formats, b) Roadmap type, c) Supported mobility model, d) Implemented traffic model.

Network simulators can be used to simulate network components in a detailed way such as source, destination, channel and data traffic transmission. The main features of the network simulator are a) large-scale simulation capability, b) ease of installation and use, c) MAC protocols and supported networks. For the network components simulation, we use NS3 which is a discrete-event network simulator for Internet systems targeted primarily for research and educational use [16]. Nowadays, NS-3 is still evolving where some models are still under development. However, the NS-3 simulator remains one of the smartest choices like OMNeT ++, JISt are. We chose NS-3 because it is a free and open source software, for its potential to expansion that materializes in the large and growing number of source code contributed by the NS-3 computing community and for its extensibility and stability. To allow easy integration of new models in NS-3, its network architecture is inspired by the real world in terms of hardware and software.

Finally, the VANET or Framework simulators allow integration and coupling between the network simulator and the traffic simulator. The degree of integration level will or will not allow the vehicle to change direction or leave the road in response to a network event.

Some popular platforms of vehicular networks simulation can be cited such as VEINS [17], [18], iTETRIS [19]-[21] and VNS [22] which are considered as very powerful platforms to simulate and evaluate VANETs protocols.

The interface module of the first two platforms which is used to interconnect both traffic and network simulators, introduces a communication and synchronization overhead, consequently reducing the efficiency of the global simulator. On the other hand, VNS proposes a unique framework which makes it possible to avoid additional calculation overhead by completely integrating the aforementioned NS-3 network simulator and the traffic simulator DIVERT 2.0.

Various realistic roadmap formats can be imported by DIVERT 2.0. It also provides vehicle models with different driver behaviors as well as a realistic traffic generation model.

B. Simulation Scenario

In VNS, the author proposes a modified version of the Manhattan model where traffic lights are present at each intersection. Several lanes are present at each street. The Divert 2.0 mobility model defined in VNS introduces elements relating to the characteristics of a micro-mobility model such as the driver model in each vehicles that offers more realism. The rate of creation of the vehicles is given by the user.

This modified version of Manhattan in VNS is a result of a manifestation of an initial transitional period that depends on the size of the map (number of streets). This period is followed by a stationary vehicle density and in which the simulation results are taken into account.

Our simulation was structured as follows:

- The transformation into parameterized variables of several simulation parameters considered important and which are likely to change during the simulation.
- The ability to change the simulation parameters from the command line.
- The choice to be able to control the simulation by means of a python script.
- The evaluation and the tracing of the results are done by the use of the library Scipython.

We performed a VNS simulation using the input parameters represented in Table I. The first dissemination message is triggered randomly after the fulfillment of the following condition:

$$nbr_V \geq Trunc\left(\frac{nbr_L * len_L * dens_V}{100}\right) \quad (4)$$

To avoid the initial transition period of the simulation, we established this formula empirically after several attempts. To avoid re-executing the simulation several times, the sending of a new emergency message is retriggered each time period defined by the “time_to_send_new_message” parameter. For each transmission, the triggering vehicle must be within 30 meters of the center of the map.

TABLE I. INPUT PARAMETERS OF A VNS SIMULATION

Input parameters	Definition
gui	controls the visualization of the result as an animated graph
finishtime	control the duration of the simulation
Nbr_L	sets the roads number
len_L	defines the length of the slices of roads and thus the size of the blocks (block of house)
dens_V	defines the rate of creation of vehicles on each road
time_to_send_new_message	sets the minimum time between two sendings of a new emergency message
maxspeed	Sets the maximum speed for each route
normal_range	defines the standard height vehicle reception range
high_range	sets the highest vehicle receipt range (Bus, Truck)
high_vehicle_ratio	the percentage of high vehicle creation
dhvn_snf	sets the period of transmission of the message by the relay DHVN
dhvn_message_TTL	The lifetime of the message

C. Evaluation Parameters

We performed simulations for the DHVN protocol in an urban and rural environments. The different simulation parameters are represented in Table II.

TABLE II. SIMULATION PARAMETERS OF DHVN PROTOCOL (URBAN AND RURAL)

Parameters	Values
Simulation zone	1200m*1200m in urban and 3000m*3000m in rural
Number of Road	5
Length of the Roads	500m in rural and 200 m in urban (5x5 roads)
Number of nodes	700 a 1750 in urban and 2000 à 3000 in rural
The propagation model	Three Log Distance Propagation Loss Model and Nakagami propagation loss model with two variations of emission power for each type of vehicle height
Vehicle height	1m (ordinary vehicle) -- 2m (truck or bus)
High vehicle rate	20%
Coverage area (Radio Range)	250m (ordinary vehicle) --375 m (truck or bus)

To evaluate the performance of the DHVN protocol, we will focus on the following metrics:

- Coverage which equals the total number of successful first receptions divided by the estimated number of vehicles in circulation.
- The average number of duplicate reception which is equal to the total number of duplicates receptions for a message divided by the estimated number of vehicles in circulation.
- Performance index, which is calculated according to the formula $\frac{coverage^2}{Duplicate_reception_average/3+100}$ where we have emphasized the importance of coverage and reduced that of duplicate reception. We will rely on this index, whose values vary between 0 and 100, for the choice of the optimal result.

IV. OUR PROPOSITION

A. Exploring Some Enhancement Approaches

We performed a simulation of the DHVN protocol for different SNF periods and obtained the results of both urban and rural environments as illustrated in Fig. 2 and 3, respectively. We can clearly see that the performance of the DHVN protocol has great sensitivity to the SNF period with the performance index varying between 26 and 51 for the urban environment and between 20 and 54 for the rural environment. The Best results are obtained for an SNF period of 45s and 20s for the urban and rural environment respectively. This is apparently tied to the environment (lane length, intersection density and vehicle density) and the DHVN protocol user must choose the SNF period very carefully. In our modified algorithm, we will try to attenuate this sensitivity and find a logic to select an adaptive SNF period without the intervention of the user.

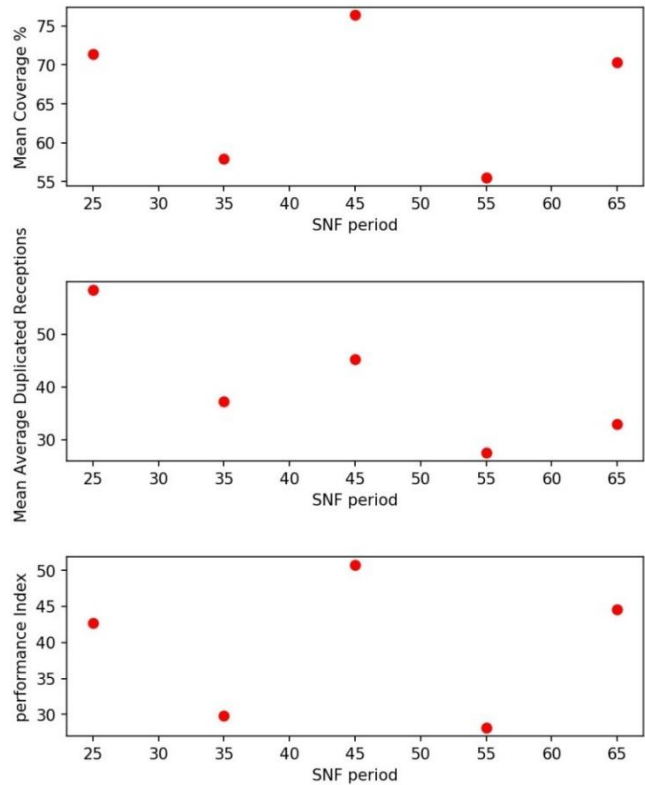


Fig. 2. DHVN protocol simulation results for different SNF periods in urban environment.

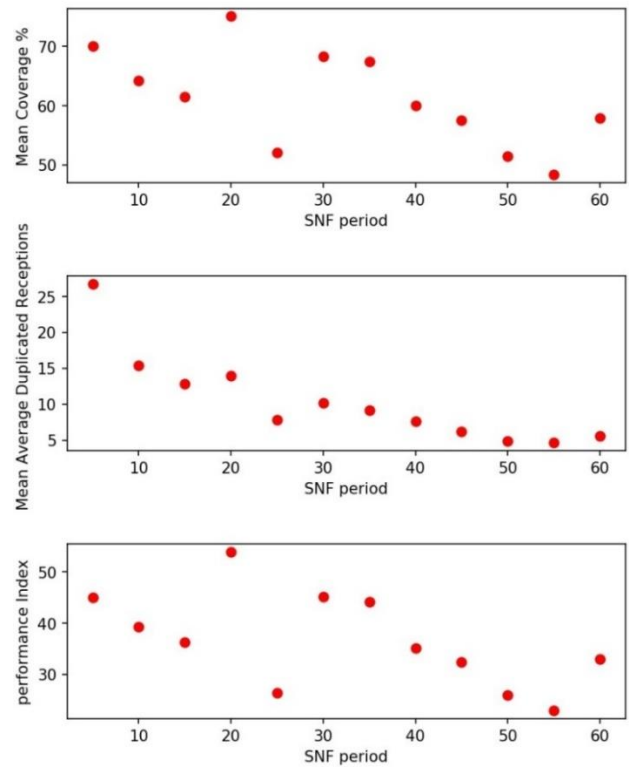


Fig. 3. DHVN protocol simulation results for different SNF periods in rural environment.

Additionally, as we can see in Fig. 3, in some special cases where the original transmitter is in very low vehicle density zone, the DHVN protocol misfires and resulted performance is very poor. This is very likely the case for the SNF periods where the coverage is lower than 50%.

Another cause for these poor results would be tied to SNF period termination being triggered in a very low vehicle density zone. This means that the relay vehicle would be in the waiting state while passing high density zones like road intersections. We also consider this case to be a DHVN protocol misfire. On the other hand, in urban environment where mostly all the zones are high vehicle density ones, the DHVN protocol produces too many duplicate messages which means a high wastage of the transmission bandwidth. We will be taking these observations into account in our enhanced algorithm by opportunistically synchronizing the transmission of the DHVN message with the passage of the relay near an intersection zone.

To further explore the DHVN protocol performance variation, we performed a simulation for different TTL durations and an SNF period equal to 45 seconds in the urban environment and equal to 20 seconds in the rural environment. The obtained results of both urban and rural environments are illustrated in Fig. 4 and 5, respectively.

We can observe again a certain sensitivity of the resulted performance to the TTL duration with the performance index varying between 36 and 54 for the urban environment and between 11 and 61 for the rural environment. The best results are obtained for the TTL duration of 150s and 400s for the urban and rural environment, respectively. This is easily explained due to the SNF mechanism; a short TTL duration doesn't allow for the message to propagate far enough, while a large TTL duration results in excessive bandwidth usage. So logically, an ideal TTL duration would be a function of the size of the dissemination area and the mobility model (speed, acceleration, stop sign waiting time, etc.). We will not be introducing a selection mechanism in our enhanced algorithm but we will be using these values for further simulation so we can have the best performance results.

Another aspect of the original DHVN protocol that can be criticized is the use of the distance between the receiver and the transmitter to calculate the backoff delay. While in an ideal scenario, the DHVN relays would be moving away from the position where the first message originated; this is especially not true in a sparsely distributed road map. So we had the idea of calculating the backoff delay using the distance between the receiver and the position of the original transmitter. This later position is preserved in the transmitted DHVN message.

B. Final Detailed Algorithm

The resulted enhanced algorithm for the Opportunistic DHVN is illustrated in Fig. 6 with D_{to_inter} is the distance of the vehicle to the next intersection, T_Range is the transmission range of the vehicle and SNF_P is the adaptive SNF period.

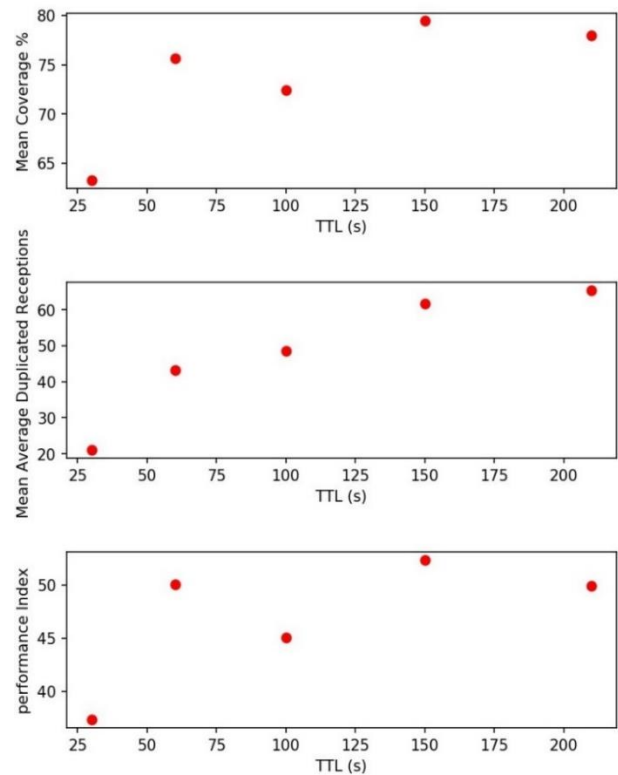


Fig. 4. DHVN protocol simulation results for different TTL duration and SNF periods of 45s in urban environment.

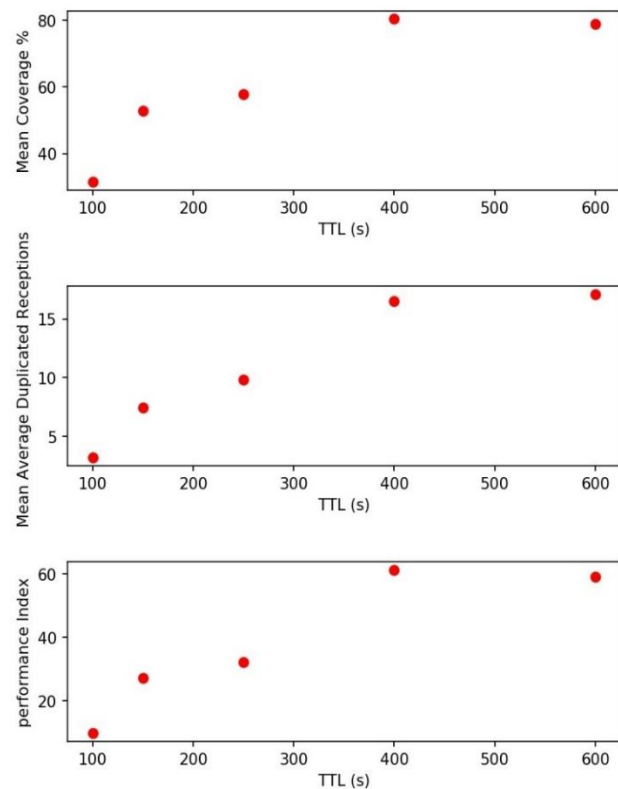


Fig. 5. DHVN protocol simulation results for different TTL duration and SNF periods of 20s in rural environment.

```
While (Position is in Dissemination Area){
  function Receive (msg){
    if (msg.tag is ID){
      Wait (backoff_time+SNF_P)
      Set msg.tag equal to SD
      Retransmit(msg)
    }
    elseif (msg.tag is SD)|| (msg.tag is DNI){
      if (msg.tag is DNI)
        Set timer_delai equal to
        backoff_time+SNF_P
      else Set timer_delai equal to
        backoff_time
      if (same_road) {
        if (first reception)
          Trigger timer`
        elseif (duplicate &&
          sender is in same_road)
          Cancel timer`
      }
      if (Intersection Zone){
        if (first reception)
          Trigger timer`
        elseif (duplicate &&
          sender is not in the same road)
          Continue to disseminate
      }
    }
  }
}
} //end while

Function Timerfired(){
  if ((T_Range>=D_to_inter) &&
    (D_to_inter >= T_Range/4)){
    Wait Until (D_to_inter <= T_Range/4)
  }
  if (Intersection Zone) Set Tag_Diffusion
  equal to DNI
  Retransmit(msg)
  Trigger timer with SNF_P
}
```

Fig. 6. Opportunistic DHVN protocol algorithm.

The algorithm introduces the following changes:

- Three tags are introduced: Initial Diffusion (ID), Stranded DHVN (SD) and DHVN Near Intersection (DNI). The Initial Diffusion tag is used for the first transmission only while DNI tag is used if the relay is inside an intersection. The SD tag is used elsewhere.
- The condition “Same road” cancels the DHVN Timer instead of “Same road before”.
- An adaptive SNF period is introduced and is calculated as $(T_{transm_Range}/V_{max_speed})$.
- If the vehicle relay is close to an intersection ($T_{Range} \leq D_{to_inter} \leq T_{Range}/4$) and the SNF period termination is triggered, the transmission is delayed until the relay is about to enter the intersection ($D_{to_inter} < T_{Range}/4$).

- In the first DHVN retransmission the vehicles wait for $Backoff_Time+SNF_period$ instead of just $Backoff_Time$ to counter balance the Initial Diffusion flood message.

V. RESULTS AND DISCUSSION

We performed a simulation of the new oDHVN protocol for different TTL durations taken around the previously outlined 150s and 400s TTL durations for the urban and rural environments, respectively. The obtained results are illustrated in Fig. 7 and 8, respectively.

In the case of the urban environment and except for the very low TTL duration of 50s, we can observe a stable performance index around 60 with the best result being 65 for the TTL duration of 200s.

In the case of the rural environment we can observe a slightly larger variation of the performance index between 56 and 77 with the best result being for the TTL duration of 600s.

We can clearly observe the enhancement in the performance index value and stability. Additionally we didn't have to choose an SNF period which gives the introduced oDHVN protocol a great advantage in autonomy and ease of use.

VI. CONCLUSION

In this work we studied a palette of the VANET information dissemination protocols that we classified into three major classes, Relayed Instantaneous Broadcast, broadcast with store-and-forward SNF relay, and opportunistic broadcast. We underlined the performance of a recently introduced algorithm in [3], which is the DHVN protocol. We consider the latter to be a very good representative of the second class of dissemination protocols and we chose it as a basis for a new enhanced opportunistic dissemination protocol, the oDHVN.

Our contribution is represented by the enhanced algorithm which we have taken as the basis for our new oDHVN protocol. This algorithm was the fruit of well thought analysis of the behavior and performances of the standard DHVN protocol in both urban and rural environments. We have presented the results of the simulations carried out on the VNS framework for both the standard DHVN and the introduced oDHVN.

These results state clearly the enhancements brought by the oDHVN protocol. These enhancements are observed, on one hand, in the maximum value of the performance index which translates an equilibrium between the coverage percentage and the bandwidth usage. They are also observed, on the other hand, in the overall stability of the performance indicators.

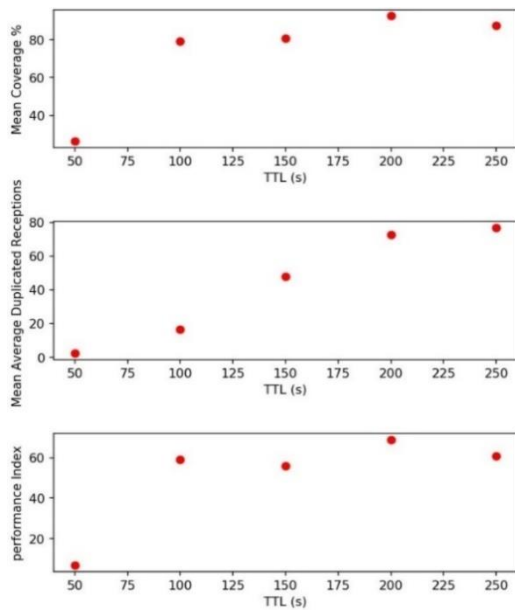


Fig. 7. Opportunistic DHVN protocol simulation results for different TTL duration in urban environment.

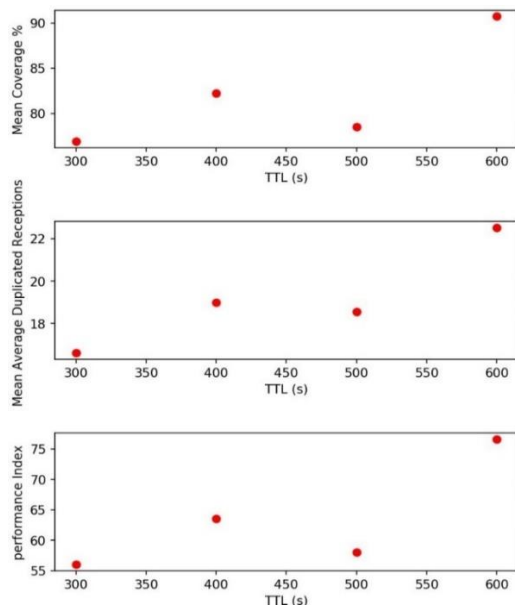


Fig. 8. Opportunistic DHVN protocol simulation results for different TTL duration in rural environment.

REFERENCES

[1] L. Aparecido, "Data dissemination in vehicular networks: Challenges, solutions, and future perspectives", in 2015 7th International Conference on New Technologies, Mobility and Security (NTMS), 2015, p. 1-5.
 [2] U. Hayat, R. Iqbal, et J. Diab, "Eliminating Broadcast Storming in Vehicular Ad-Hoc Networks", *Int. J. Adv. Comput. Sci. Appl.*, vol. 7, no 1, 2016.
 [3] S. Mehar, S. M. Senouci, et G. Rémy, "Dissemination protocol for Heterogeneous Cooperative Vehicular Networks", in 2012 IFIP Wireless Days, 2012, p. 1-6.
 [4] S. Panichpapiboon et W. Pattara-Atikom, "A Review of Information

Dissemination Protocols for Vehicular Ad Hoc Networks", *Commun. Surv. Tutor. IEEE*, vol. 14, no 3, p. 784-798, Quarter.
 [5] I. Achour, T. Bejaoui, et S. Tabbane, "Network coding approach for vehicle-to-vehicle communication: Principles, protocols and benefits", in 2014 22nd International Conference on Software, Telecommunications and Computer Networks (SoftCOM), 2014, p. 154-159.
 [6] J. A. Sanguesa, M. Fogue, P. Garrido, F. J. Martinez, J.-C. Cano, et C. T. Calafate, "A Survey and Comparative Study of Broadcast Warning Message Dissemination Schemes for VANETs", *Mobile Information Systems*, 2016. [En ligne]. Disponible sur: <https://www.hindawi.com/journals/misy/2016/8714142/>. [Consulté le: 27-déc-2017].
 [7] A. Mchergui, T. Moulahi, B. Alaya, et S. Nasri, "A Survey and Comparative Study of QoS Aware Broadcasting Techniques in VANET", *Telecommun Syst*, vol. 66, no 2, p. 253-281, oct. 2017.
 [8] R. Ghebleh, "A comparative classification of information dissemination approaches in vehicular ad hoc networks from distinctive viewpoints: A survey", *Comput. Netw.*, vol. 131, p. 15-37, févr. 2018.
 [9] M.-T. Sun, W.-C. Feng, T.-H. Lai, K. Yamada, H. Okada, et K. Fujimura, "GPS-based message broadcast for adaptive inter-vehicle communications", in *Vehicular Technology Conference, 2000. IEEE-VTS Fall VTC 2000. 52nd, 2000*, vol. 6, p. 2685-2692.
 [10] R. C. Voicu, H. I. Abbasi, H. Fang, B. Kihei, J. A. Copeland, et Y. Chang, "Fast and reliable broadcasting in VANETs using SNR with ACK decoupling", in 2014 IEEE International Conference on Communications (ICC), 2014, p. 574-579.
 [11] N. Ababneh et H. Labiod, "Safety message dissemination in VANETs: Flooding or trajectory-based?", in *Ad Hoc Networking Workshop (Med-Hoc-Net), 2010 The 9th IFIP Annual Mediterranean, 2010*, p. 1-8.
 [12] S. Zemouri, S. Djahel, et J. Murphy, "Short paper: Road-Casting: A new distributed dissemination protocol for safety messages in urban areas", in 2013 IEEE Vehicular Networking Conference (VNC), 2013, p. 234-237.
 [13] M. O. Cherif, S.-M. Secouci, et B. Ducourthial, "How to disseminate vehicular data efficiently in both highway and urban environments?", in 2010 IEEE 6th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob), 2010, p. 165-171.
 [14] B. Tian, K. M. Hou, et J. Li, "TrAD: Traffic Adaptive Data Dissemination Protocol for Both Urban and Highway VANETs", in 2016 IEEE 30th International Conference on Advanced Information Networking and Applications (AINA), 2016, p. 724-731.
 [15] M. Rathod, I. Mahgoub, et M. Slavik, "A hybrid data dissemination scheme for VANETs", in *Wireless Days (WD), 2011 IFIP, 2011*, p. 1-7.
 [16] J. Balen, J. Matijaš, et G. Martinović, "Simulation and testing of VANET protocols", *Tridesettreći Skup O Prometu. Sustavima Med Junarnodnim Sudjelov. Autom. U Prometu KoREMA*, p. 5-8, 2013.
 [17] C. Sommer, R. German, et F. Dressler, "Bidirectionally Coupled Network and Road Traffic Simulation for Improved IVC Analysis", *IEEE Trans. Mob. Comput.*, vol. 10, no 1, p. 3-15, janv. 2011.
 [18] C. Sommer, Z. Yao, R. German, et F. Dressler, "Simulating the influence of IVC on road traffic using bidirectionally coupled simulators", in *INFOCOM Workshops 2008, IEEE, 2008*, p. 1-6.
 [19] V. Kumar et al., "itetr: Adaptation of its technologies for large scale integrated simulation", in *Vehicular Technology Conference (VTC 2010-Spring), 2010 IEEE 71st, 2010*, p. 1-5.
 [20] J. Harri, F. Filali, et C. Bonnet, "A framework for mobility models generation and its application to inter-vehicular networks", in 2005 International Conference on Wireless Networks, Communications and Mobile Computing, 2005, vol. 1, p. 42-47 vol.1.
 [21] J. Harri, F. Filali, et C. Bonnet, "Mobility models for vehicular ad hoc networks: a survey and taxonomy", *IEEE Commun. Surv. Tutor.*, vol. 11, no 4, p. 19-41, 2009.
 [22] R. Fernandes, F. Vieira, et M. Ferreira, "VNS: An integrated framework for vehicular networks simulation", in *Vehicular Networking Conference (VNC), 2012 IEEE, 2012*, p. 195-202

Student Facial Authentication Model based on OpenCV's Object Detection Method and QR Code for Zambian Higher Institutions of Learning

Lubasi Kakwete Musambo
School of Engineering
Dept. of Electrical & Electronics Engineering
The University of Zambia
Lusaka, Zambia

Jackson Phiri
School of Natural Sciences
Dept. of Computer Science
The University of Zambia
Lusaka, Zambia

Abstract—Facial biometrics captures human facial physiological data, converts it into a data item variable so that this stored variable may be used to provide information security services, such as authentication, integrity management or identification that grants privileged access or control to the owner of that data variable. In this paper, we propose a model for student authentication based on facial biometrics. We recommend a secure model that can be used in the authentication and management of student information in the registration and access of resources, such as bursaries, student accommodation and library facilities at the University of Zambia. Since the model is based on biometrics, a baseline study was carried out to collect data from the general public, government entities, commercial banks, students, ICT regulators and schools on their understanding, use and acceptance of biometrics as an authentication tool. Factor analysis has been used to analyze the findings. The study establishes that performance expectancy, effort expectancy, social influence and user privacy are key determinants for application of a biometric multimode authentication. The study further demonstrates that education and work experience are regulating factors on acceptance and expectancy of a biometric authentication system. Based on these results, we then developed a biometric model that can be used to perform authentication for students in higher learning institutions in Zambia. The results of our proposed model show 66% acceptance rate using OpenCV.

Keywords—Biometrics; authentication; model; integrity

I. INTRODUCTION

Applying a secure biometric infrastructure is a key in ensuring that organisational and or private data is well managed and accessed only by the intended party. It is important that a possibility to authenticate only those individuals that are registered as students of a high institution exists [1]-[3]. This study focuses on the University of Zambia (UNZA), which is Zambia's biggest institution of higher learning. The findings can be generalised to cover the rest of Zambia's higher institutions of learning. The current authentication processes for UNZA are paper-based systems installed by the management of the university. Though these

authentication processes which are prone to data redundancy are implemented; issues of over payment to ghost students on bursaries always arise, issues of illegal residents at the university arise from time to time, issues of non-availability of student records or non-available student records which were earlier created and filed with the office of Dean of Students arise. To overcome such problems generated by a lack of a secure student authentication system, we present a biometric model based on two-factor authentication that can be used.

II. LITERATURE REVIEW

A. Erroneous Payments

During a student registration process especially for the postgraduate level of studies at UNZA, a student is required to make payment at 4 different points; the Dean of Student office, the School of study, the Office of the Directorate of Research & Graduate Studies and the Library. In this process, students have been known to make erroneous payments to different pay points and correcting this error has been problematic for the student. Students have had to make extra payments in-order to meet the payment plan set by the University.

B. Squatting

The UNZA student handbook prescribes rules and regulations that a student must abide by. The rules indicate the number of students that can share a room for the purposes of being accommodated by the University. In 2018 first quarter, a student living conditions audit necessitated by an outbreak of cholera, a water-borne disease showed that more than 2 students shared a room. This trend referred to as squatting is an illegal activity. A lack of authentication (binding a student to a room) is the problem here as it appears difficult to determine which student is in which room.

Library Access: Access to the UNZA Library resources and access to an examination hall is guaranteed via a student ID card as shown in Fig. 1. No other form of authentication is available. It is possible that one may create a bogus ID card and gain access.

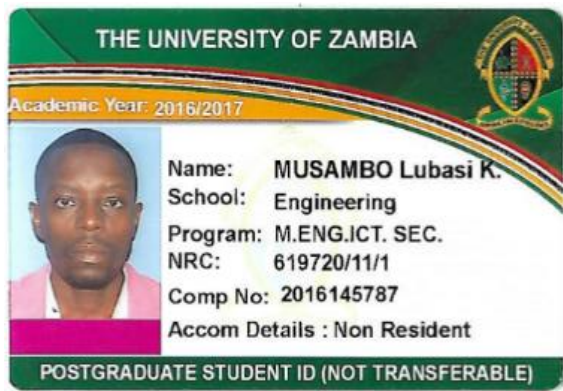


Fig. 1. UNZA Student ID.

C. Ghost Students

In Zambia, all government run institutions are audited annually by the Auditor General's office. The audits are conducted to determine usage of public funds. The audit concludes in what is called The Auditor General's Report. UNZA is a government institution and as such is audited. The Auditor General's report for the period ended 2016 shows that the Zambian Ministry of Higher Education paid tuition and meal allowances to 543 individuals who were not registered by Zambia's two most prominent learning institutions; The University of Zambia and The Copperbelt University (CBU). These funds are termed 'meal and project allowances'. The report further states that over 8,000 UNZA students had been left out of these payments resulting into the students rioting. These funds are meant to secure accommodation, meal and project allowances for the students. The payments in question are within the region of K8,521,629 (or approximately \$852,162.9). These are payments to ghost students. The reason this is possible is because a lack of authentication (entity binding exists). The Higher Education Loans and Scholarships Board Management is a Zambian government institution that manages student bursaries on behalf of the Zambian government. The board has no mechanisms to authenticate student status due to a lack of university presence and reliance on paper based systems that are easily manipulated.

D. Loans and Bursaries

The Government of the Republic of Zambia awards loans and Educational bursaries to deserving students at Zambia's high institutions of learning. These bursaries are meant to ease the pressure of meeting school fees by a student. Payment of school fees is a requirement before one is admitted into school. Verification of who is entitled and who has been registered into this loan scheme is problematic as a clear authentication procedure is weak because it is paper-based. The current authentication solution is that the loans and bursaries board is forced to physically setup office at the institution of learning when learning institutions open. These temporal offices are used as screening facilities to screen and activate accounts of the students that have qualified for the bursaries. This is a labour intense and time consuming activity. A solution provides an automated authentication mechanism that can plug-into the student database system to perform the authentication when students register.

E. Student Registration

Registration at UNZA is a process of being enrolled into the school's student database system. Usually passport photos and physical copies of the student's credentials are needed for filing purposes. The credentials are authenticated by any third party referred to in Zambia as a commissioner of oath. The registration process can be enhanced if a centralised civil registration biometric database bound to an education database existed. This would authenticate a potential student's credentials to a higher institution of learning.

III. SUMMARY OF REVIEWED BIOMETRIC AUTHENTICATION SYSTEMS

Biometrics can be collected from either a physiological characteristic or a behavioral characteristic. A physiological characteristic is a relatively stable human physical feature. An example of a physiological characteristic is a fingerprint, retina iris pattern, or a hand-geometry pattern. Physiological measurements are static and non-alterable. This type of measurement is unchanging and irreversible or permanent apart for deformity caused by external significant duress such as ailment or physical injury [4]. A behavioral characteristic on the other hand attempts to resemble a person's psychological makeup. This is affected by a person's build stature and gender among others. Behavioral characteristics can be identified in activities such as speech, hand-writing speed and pressure exerted on paper when writing among others. Four methods of biometric authentication systems were reviewed employing both physiological and behavioral characteristics. These have been reviewed in terms of basic operation, advantage and disadvantage of implementation.

A. Fingerprint Authentication

Fingerprints are made up of ridge patterns on a person's fingers. These ridge patterns have capacity to uniquely distinguish and identify individuals. Fingerprint features are made up of arches, loops, and whorls. An individual fingerprint will exhibit at least one of these major features. The minor details that are collected from these fingerprint features are referred to as minutiae. Fig. 2 and 3 show a finger print sample and finger print features. The authentication processes is an automated method of verifying a match among different human fingerprints [5].

1) Advantages

- a) Individualistic features guarantee authentication of subject [4].
- b) Systems are relatively inexpensive to purchase and install.
- c) Longevity of life of the fingerprint pattern's individualistic feature composition guarantees long term usage [4].
- d) Once in use a subject does not have to rely on memory for passwords as fingerprint authentication will guarantee access.
- e) A fingerprint identity point cannot be spoofed [6].



Fig. 2. Fingerprint image sample [17].

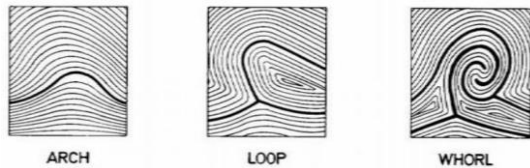


Fig. 3. Fingerprint features [9].

2) Disadvantages

a) Limitation of capture is reduced to an individual finger with further limitation of capture reduced to a section or part of that finger only and not the entire finger.

b) Susceptible to FAR (false acceptance error) whereas a wrong subject is enrolled and access is allowed.

c) Hand injury (fingers included), chemical prone jobs and labour prone activities such as brick-laying or metal fabricating present a within-person variation that makes the reading and capture of finger prints difficult.

d) Washing with a soap detergent or submerging a finger in water for period of time (approximately 30 minutes) works as a contraceptive to finger-print scanners and this may impede the scanners from capturing or enrolling the finger prints until the finger reverts to its original form it was in during capture or enrolment [7].

B. Retina Authentication

This is one of the two forms of eye biometrics; the other being iris recognition. This form of biometrics is one of the most secure authentication systems in place today. The installed technology requires that an impression of a retina pattern must be taken and stored. The authentication process involves evaluating a subject's retina with a stored version (impression enrolled) of that subject's retina. Retina recognition has a low FAR (false acceptance error) as well as low rejection rates [8]. An image sample of an eye is shown in Fig. 4.

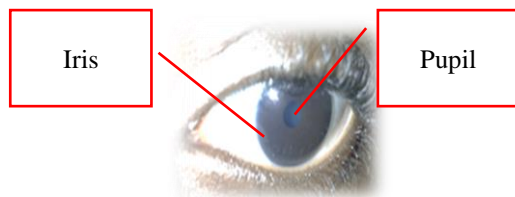


Fig. 4. Eye image sample – for iris Recognition.

1) Advantages

a) Different even in identical twins.

b) Highly specific with unique structure shape and limits the possibility of fake retina presentation.

c) Longevity of structure throughout life time of subject.

d) Wearing of glasses or contact lenses does NOT work as a contraceptive to technological accuracy.

e) High accuracy and High recognition process speed.

2) Disadvantages

a) Eye injury or sickness may render this biometric system ineffective.

b) Intrusive technology and may not be welcomed by many individuals.

c) Lighting may affect the accuracy of the reader.

d) Fairly expensive to acquire when compared to other systems of biometrics.

C. Voice Authentication

This technology allows the conversion of voice or sounds from human voice into an electrical signal that can be coded. Voice recognition software is designed to identify an individual via their unique voiceprint. Voiceprints are generated from physical characteristics of an individual's throat in conjunction with their mouth. Research indicates that no two voices are the same and therefore voice biometrics provides a rare opportunity to use one's voice to authenticate or identify individuals [3]. A sample of a voice pattern is shown in Fig. 5 below.

1) Advantages

a) No need for user training as users can simply speak into the voice biometric reader.

b) Voice communications is a natural activity for human beings.

c) Voice communications eliminates the need to learn keyboard operations (and in this way helps to bridge the gap between the able-bodied and individuals who experience restricted capabilities in hand based motion activities such as writing). By eliminating the learning aspect, voice overcomes the need to learn how to operate some complex biometric technology's operations.

d) It eliminates the need to be accurate in written statements as is for password based authentication.

e) Because one uses voice, the speed of operation is enhanced. People generally speak faster than they are able to write.

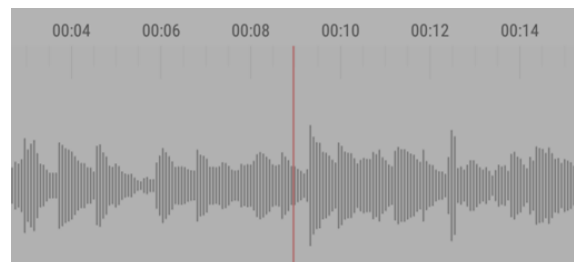


Fig. 5. Voice print. Adapted from [5].

2) Disadvantages

- a) Impulse noise may affect the accuracy of the voice signal and render the system ineffective.
- b) Microphone proximity must be precise for the system to work well.
- c) A pre-recorded audio may by-pass this system.
- d) A person may speak different languages and this may affect the accuracy of the device should that individual use a different language or dialect.
- e) Certain words have a homonym characteristic this may affect the accuracy of the device.
- f) The learning curve for the system may be long as it is trained per voice.
- g) Most voice controlled biometrics is expensive.

D. Face

Facial biometrics divides into two aspects namely the face detection and face recognition programs. Face recognition extracts a face from a given image while face recognition compares a captured face against saved faces in order to match the face. The entire process is run by a series of complex algorithms. One of the options of face recognition is to select features of a face and match those features to a face. Fig. 6 below shows a facial image sample with facial image mapping that is used to collect facial features. The facial features or dataset is normally stored in a database. In ideal situations this database must be encrypted to achieve sufficient security [9].

1) Advantages

- a) Non-intrusive technology and can be performed stealthily without the subject knowing, therefore, proves ideal for investigation purposes.
- b) Certain algorithms can be adjusted to scan a large scale of a population and thus this technology proves ideal in crowded environments.
- c) Ideal for person tracking and incident reporting.
- d) User friendly as far as users are concerned as no need of complex training for the subjects to be captured.
- e) Can be developed and run from a basic computer camera without buying any other tools. This proves to be one of the strongest advantage and reduces the cost of this technology exponentially.

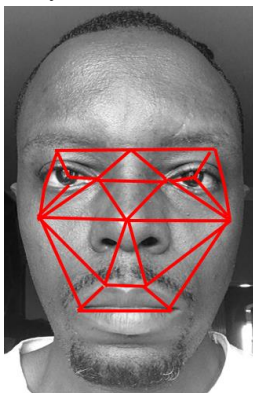


Fig. 6. Facial Image Sample with facial map for Facial Biometrics. Adapted from [5].

- f) Some easy to install ready to use pre-trained facial calibration tools are available. This again reduces cost of setup.
- g) Facial biometric algorithms have a within-person variation calculation that can detect aging and basic facial deformity and reduce a face to a known variable [10].

2) Disadvantages

- a) Certain algorithms may NOT work well on black faces.
- b) Light conditions and camera capabilities may affect the accuracy of the technology.
- c) Within-person variations may affect the accuracy levels of the technology [11].
- d) When used for security purposes, extra equipment to provide lighting can increase cost of setup.

IV. QR (QUICK RESPONSE) CODE

A typical QR code will have the shape as shown below:



Fig. 7. Sample QR Code [12].

As shown in Fig. 7 above, a QR code is a machine readable imprint made out of an array of black and white squares that normally embed certain information within the print. QR codes were developed by a Japanese company called Denso Wave for purposes of tracking manufacturing processes. QR codes however, provide an opportunity to authenticate as well as identify an entity. In this way QR codes may be used as an added security feature especially in logging into networks. Networks may be designed to read QR codes, verify the data and offer or deny access to an entity. Because QR Code information is non-human readable, this provides a basic form of information hiding in plain sight (encryption). This hidden information can then be transmitted. When used with geo-tagging, QR Codes can be used to determine a location status of an entity [12].

V. STUDENT AUTHENTICATION AND PREFERRED BIOMETRIC MODEL

A. Problem Statement

Student authentication is equivalent to entity authentication. 'Entity authentication is the assurance that a given entity is involved and currently active in a communication session'. A need to bind a student registered with a learning institution to a

current resource access of that institution exists [1]. There is need to grant student privileges such as accommodation, bursaries, allowances and loans to a deserving student automatically; a need to allow a student writes an exam without the need for unnecessary paper work is eminent; a need to ensure that a student’s location status within UNZA facility exists. To achieve these functions, we recommend a facial biometric solution with a mobile QR Code reader.

B. Understanding the Haar based Frontal Face Biometric Algorithm

Based on a rapid object detection scheme based on boosted cascade of simple feature classifiers introduced by Paul Viola and Michael Jones, a facial biometric model can be developed based on Haar-like features and implemented to detect and recognise a student’s face. This recognition facility allows for authentication. Facial features to form a Haar classifier are collected after a facial mapping as shown in Fig. 8. The biometric model utilises Haar basis features as used by Papageorgiou et al. [13].

An adaption of the algorithm based on an OpenCV Open Source technology which is readily available from OpenCV has been used. This algorithm uses Haar like features and OpenCV pre-trained classifiers for face detection. A classifier is a program that can decide whether an image is positive or not. A positive image is an image face (image having a face) while a negative image is a non-face image. Classifiers are trained from a huge volume of faces (both positive and negative images) to learn how to classify a new image correctly. This is a machine learning concept. The classifiers used for this student authentication is the HaarClassifier which is earlier developed by Viola et al. [14]. Haar Classifiers process data in grey scale (non-colour). Colour is inconsequential in determining whether an image has a face or not.

1) Haar Classifier function logic

Viola et al states each object has features that are unique and can be used to identify and recognize that object. Haar features can be picked out from edge, line, center and diagonal features of an object as shown in Fig. 9.

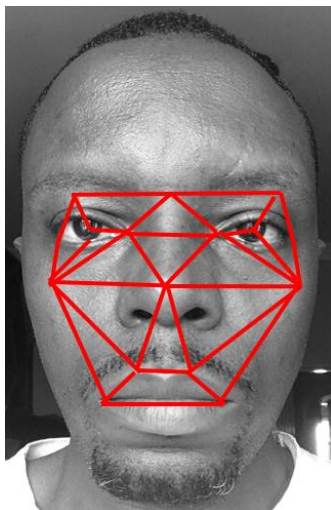


Fig. 8. Identifying features by a biometric reader. Adapted from [5].

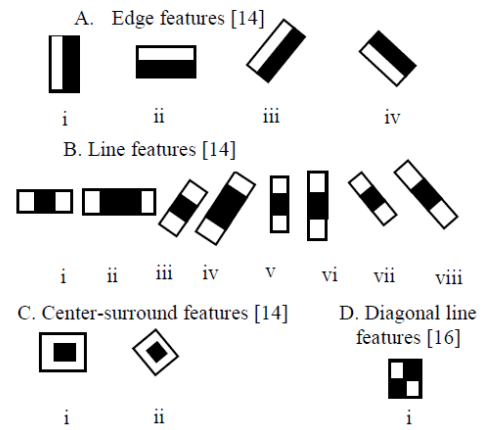


Fig. 9. Example feature determination for extraction [15].

Edge features are characteristics of an image that are unique and at unique distances from each other. No two people share the same features. The features can be mapped by placing an object identifying feature. A biometric model developed to pick up the readings from the facial recognizer can pick up the features and collectively store them to perform identification and recognition. The features can be collected into small elements referred to as a weak classifier which when collectively used identify and recognize an object [15]. Feature collection is done via rectangles. Haar like features consist of two or more rectangular regions enclosed in a template. Each of the rectangles is a window that is placed on an image as shown in Fig. 10 that is to be captured and recognized. A feature is extracted from subtracting the sum of pixels under the white part from the black part of that window (rectangle).

In determining the haar like features an understanding that the area around the eyes have a darker area then the nose bridge is used. This view is also held for the cheeks (brighter than other areas), though the data from the cheeks is not necessarily used.

Rectangles are placed on an image so as to pick the features using a weak classifier. The features of a rectangle are computed using an integral function of the form:

$$ii(x, y) = \sum_{x' \leq x, y' \leq y} i(x', y'), \quad (1)$$

In this function an object or image at location x, y contains the sum of pixels above and to the left of x, y inclusive.

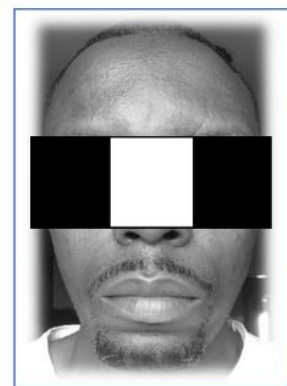


Fig. 10. Feature Determination. Adapted from [18] [20].

Where, $ii(x, y)$ and $i(x, y)$ is the original image. Using the following pair of recurrences:

$$s(x, y) = s(x, y - i) + i(x, y)$$

$$ii(x, y) = ii(x-i, y) + s(x, y)$$

(Where $s(x, y)$ is the cumulative row sum, $s(x-a) = 0$, and $ii(-i, y) = 0$). Using the integral image any rectangular sum can be computed in four array references [14]-[16].

The rectangle itself can be understood to have an object of pixels $W \times H$ (i.e. to say width x Height) [14]. Fig. 11 shows the determination of a rectangular region of an integral image.

To determine the sum of pixels, the logic can be deduced as follows:

$$a = \text{sumRec}(\text{pixels}) \tag{2}$$

$$b = 1 + 2,$$

$$c = 1 + 3$$

$$d = 1 + 2 + 3 + 4$$

The sum is then derived as $d + a - (b + c)$.

Using the OpenCV library of face detectors and recognizers a function can be developed into a web based application that can perform an online web authentication at UNZA where a student is interacting with the institution such as a library service. Between the web and the OpenCV recognisers a batch file mechanism as shown in Fig. 12 below can be incorporated to pass control to the OpenCV recognizers. OpenCv recognizers are developed in python [14]. A means of communication with a web application developed in a programming language python is achieved via the batch files as shown in Fig. 13 below. The algorithm has been set to capture 100 faces per student.

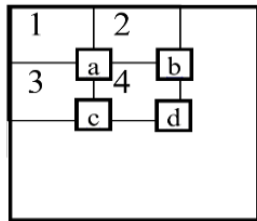


Fig. 11. Rectangular regions of an integral image [19].

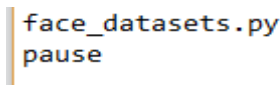


Fig. 12. Batch call function from PhP to python.

```
<title>VR|Registration </title>
<base />
<?php include_once 'userheadfiles.inc'; ?>
<script type="text/javascript" src="js/qrcode_gen.js"></script>
<script type="text/javascript">
function myFunction(){
    WshShell = new ActiveXObject("Wscript.Shell"); //Create WScript Object
    WshShell.run("c://xampp/htdocs/face/run.bat"); // Please change the path and file name
    execute .exe file as well
}
</script>
<script type="text/javascript">
```

Fig. 13. Batch call function from PhP to python to call for OpenCv face recognizer.

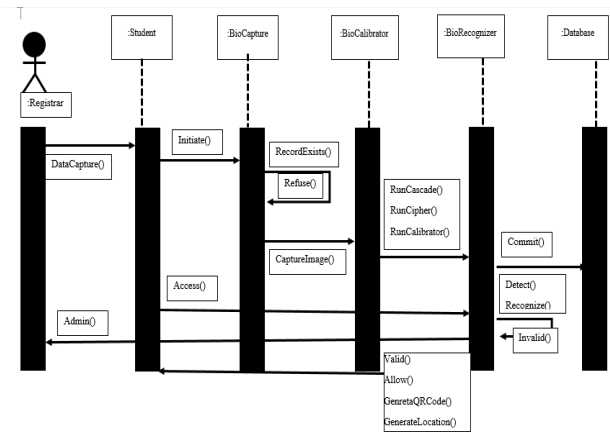


Fig. 14. UML interaction sequence for Student authentication.

A student can be registered only once after that the recognizer would perform the authentication for every other function. The interaction sequence for the student facial model is shown in Fig. 14 above.

VI. MODEL TEST RESULTS

The classifier described in the paper was implemented on authenticating students at different times of the day. This image set collected used 3000 student image faces. The system achieves a person detection rate of 66% with a 33% false acceptance error.

VII. DISCUSSIONS

The biometric model is able to yield a positive result of 66%, the false acceptance rate of 33% has been determined to be due to lighting conditions when the images are captured and the dark faces enrolled. Performance of the model has been observed to be higher or accurate when lighter faces are used. The researchers hold the view that that the darker regions around the eyes become fairly complex for the algorithm to determine on black faces. Improving lighting conditions has been observed to correct the recognition and detection process.

We believe that biometric and QR code authentication is the right approach to the management of student authentication. Frontal face biometrics appears easier to confront as it is not expensive to develop.

A web camera mounted in a laptop or computer is sufficient for this task. It must however be understood that sufficient research is needed into ensuring that false positives are dealt with as frontal face biometrics presents false positive errors. It is also necessary to understand that trying a new technology requires ownership sense in the users and the subjects in question. An understanding of where a student's biometric data is kept is critical as most students interviewed showed little understanding of where their biometric data must be stored. This survey finding is shown in Fig. 15. It is then important that institutions of higher learning that will implement biometric technology explain to the students where their personal biometric data will be stored. It is recommended that ISO 24745 is used to guide higher institutions of learning in the secure management and usage of biometric data.

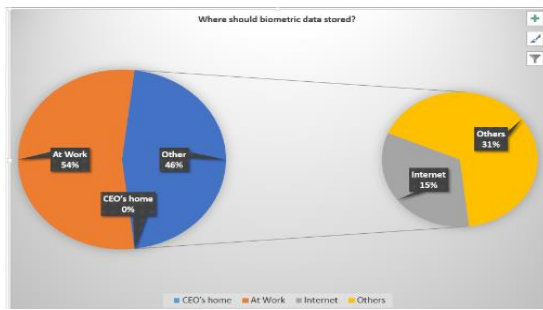


Fig. 15. Public understanding of where biometric data can be stored.

VIII. CONCLUSION

In this paper, we give the results of the implementation for a student authentication system based on our local university called UNZA. The results can however be generalized to cover other higher institutions of learning. The test results show the proposed system was able to give up to 66% accuracy level. For a developing country like Zambia with no form of automation in student's identification, this would be a good starting point.

Zambia currently does not have technological advancements that cover biometrics in detail let alone a biometric standard to determine suitable security that can be implemented in the use of biometrics. This paper recommends a frontal facial biometric model that can be used to perform authentication at various points within the university but can be generalized to any higher learning institution. The frontal facial biometrics uses OpenCV's boost algorithms which are open source and readily available for adaptation.

IX. SUMMARY

In this paper, we began by a review of the various forms of biometrics that can be used in authentication systems. We then presented the general security challenges in developing countries especially higher institutions of learning. One of the solutions to these challenges is the integration of biometrics features in the authentication systems. A cheaper solution for most developing countries is the use of open source tools and cheaper devices. Our study was proposing the use of OpenCV for Biometric Facial recognition and simple cheaper Web Camera such as one that comes integrated in most mobile computing devices. For future works, we recommend a large dataset testing comprising of a majority of black people for a full proof authentication system based on facial biometrics is required.

REFERENCES

[1] J. Phiri and J. I. Agbinya, "Modelling and Information Fusion in Digital Identity Management Systems," in International Conference on Systems and International Conference on Mobile Communications and Learning Technologies, 2006., Mome, Mauritius, 2006.

[2] I. J. Agbinya, N. Mastali, R. Islam and J. Phiri, "Design and Implementation of a Multimodal Digital Identity Management system using fingerprint matching and face recognition," *Broadband and Biomedical Communications (IB2Com)*, pp. 272-278, 21-24 Nov 2011.

[3] V. a. Tripathi, "A Comparative Study of Biometric Technologies with Reference to Human Interface," *International Journal of Computer Applications*, vol. 14, no. 5, pp. 1-6, 2011.

[4] J. Phiri, T.-J. Zhao, H. C. Zhu and J. Mbale, "Using Artificial Intelligence Techniques to Implement a Multifactor Authentication System," *International Journal of Computational Intelligence Systems*, vol. 4, no. 4, pp. 420-430, 2011.

[5] R. Saini and N. Rana, "COMPARISON OF VARIOUS BIOMETRIC METHODS," *International Journal of Advances in Science and Technology*, vol. Vol 2, no. I, pp. 1-7, 2014.

[6] N. Ferguson, B. Schneier and T. Kohno, *Cryptography Engineering: Design Principles and Practical Applications*, Indianapolis: Wiley, 2010.

[7] K. Martin, *Everyday Cryptography: Fundamental Principles & Applications*, New York: Oxford University Press, 2012.

[8] J. M. Stewart, E. Tittel and M. Chapple, *Certified Information System Security Professional*, Canada: Wiley, 2008.

[9] F. Alonso-Fernandez, J. Fierrez and J. Ortega-Garcia, "Quality Measures in Biometric Systems," in *IEEE*, 2011.

[10] A. Lanitis, "Facial Biometric Templates and Aging: Problems and Challenges for Artificial Problems and Challenges for Artificial," in *IAAI-2009 Workshops Proceedings*, 2014.

[11] E. Bilgin and B. Sankur, "Effects of Aging over Facial Feature Analysis and Face Recognition," *Bogaziçi Un. Electronics Eng. Dept.*, pp. 1-4, 2010.

[12] A. Mehta, "QR Code Recognition from Image," *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 5, no. 12, pp. 781-785, 2015.

[13] A. Mohan, C. Papageorgiou and T. Poggio, "Example Based Object detection.," *IEEE Transactions on pattern Analysis and Machine Intelligence*, vol. 23, no. 4, pp. 349-361, 2001.

[14] P. Viola and M. Jones, "Rapid Object Detection using a Boosted Cascade of Simple Features," in *CONFERENCE ON COMPUTER VISION AND PATTERN RECOGNITION 2001*, Cambridge, 2001.

[15] S.-K. Pavani, D. D. Delgado and A. F. Frangi, "Haar - like features with optimally weighted rectangles for rapid object detection," *Elsevier*, vol. 43, no. 160-172, pp. 160-172, 2010.

[16] R. Lienhart, A. Kuranov and V. Pisarevsky, "Empirical Analysis of Detection Cascades of Boosted Classifiers for Rapid Object Detection," *MRL Technical Report*, pp. 1-7, 2002.

[17] F. Alonso-Fernandez, J. Fierrez, J. Ortega-Garcia, J. Gonzalez-Rodriguez, H. Fronthaler, K. Kollreider and J. Josef, "A Comparative Study of Fingerprint Image-Quality Estimation Methods," *IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY*, VOL. 2, NO. 4, DECEMBER 2007, vol. 2, no. 4, pp. 734-743, 2007.

[18] D. Yadav, R. Singh, M. Vatsa and A. Noore, "Recognizing Age-Separated Face Images: Humans and Machines," *Pone*, 2014.

[19] M. S. Uddin and A. Y. Akhi, "Horse Detection Using Haar Like Features," *International Journal of Computer Theory and Engineering*, vol. 8, no. 5, pp. 1-4, October 2016.

[20] R. Rezaei, H. Z. Nafchi and S. Morales, "Global Haar-Like Features: A New Extension of Classic Haar Features for Efficient Face Detection in Noisy Images," R. Klette, M. Rivera, and S. Satoh (Eds.), pp. 302-313, 2014.

BLOT: A Novel Phase Privacy Preserving Framework for Location-Based Services

Abdullah Albelaihy, Jonathan Cazalas, Vijey Thayanathan
King Abdulaziz University
Jeddah, Saudi Arabia

Abstract—The inherent challenge within the domain of location-based services is finding a delicate balance between user privacy and the efficiency of answering queries. Inevitably, security issues can and will arise as the server must be informed about the query location in order to provide accurate responses. Despite the many security advancements in wireless communication, servers may become jeopardized or become infected with malicious software. That said, it is possible to ensure queries do not generate fake responses that appear real; in fact, if a fake response is used, mechanisms can be employed for the user to identify the query's authenticity. Towards this end, the paper propose Bloom Filter Oblivious Transfer (BLOT), a novel phase privacy preserving framework for LBS that combines a Bloom filter hash function and the oblivious transfer protocol. These methods are shown to be useful in securing a user's private information. An analysis of the results revealed that BLOT performed markedly better and enhanced entropy when compared to referenced approaches.

Keywords—Privacy; location-based services (LBS); oblivious transfer; Bloom Filter Oblivious Transfer (BLOT); bloom filter

I. INTRODUCTION

Owing to the proliferation of smartphones and similar mobile devices, location-based services (LBS) have found widespread use in recent years. In particular, the increasing use of smartphones has contributed to the growing popularity of LBS [1]. Consequently, new means of communication have been developed, and the information obtained from LBS has provided users with greater awareness of their surroundings [2]. Various iOS and Android applications enable users to download and use LBS and submit queries to LBS servers [3].

Users can receive LBS data through points of interest (POIs). For example, users can search for nearby stores or restaurants and check whether the data regarding the prices at these locations are accurate. Thus, LBS find increasing use in various applications. They provide simple solutions for location awareness and location sharing. Therefore, they are regarded as both helpful and advantageous [4].

While proven and demonstrated to be profoundly helpful, the inherent challenge of location-based services is one of juggling between two seemingly diametrically opposed requirements: user privacy and security and the efficiency of server response times when answering user queries. On one hand, the users demand accurate answers to their location-based queries, and, on the other hand, these same users expect their personal information to remain private and secure. The quagmire is that a server simply cannot answer a location-based query without having the location itself. More clearly,

one cannot expect an answer to, "what are the nearest gas stations to XYZ address?", unless they are willing to, somehow, someday, reveal the said address. In short, users want to have their cake and eat it too! The challenge, therefore, lies in finding a delicate balance between providing accurate and efficient responses to queries while maintaining an acceptable level of privacy for the user. LBS queries submitted by users must not only be efficient but must also ensure privacy; unfortunately, such queries often lead to security issues. The first issue is one of implementation: is it possible to make queries to a data store (service, database, or website) in a way that reveals the minimum amount of information about the data store while providing maximum utility to the user making the queries? For many years, these conflicting requirements were thought to be difficult, if not impossible, to satisfy. However, recent studies have shown that a variety of query protocols can, in fact, be used to achieve both goals, at least in a probabilistic sense [5].

Protection can be provided by a security protocol even in the presence of a malicious server, malicious eavesdropper, or malicious man-in-the-middle attacker. The development of "1-in-q" algorithms (which are described in detail below) has facilitated the creation of security protocols that can provide probabilistic guarantees about the validity of response even in the presence of a compromised or actively malicious data store. Recently, a major cryptographic breakthrough was made, whereby the possibility of fully homomorphic encryption (FHE) was demonstrated [6]-[8]. In this scenario, the data store operates on opaque binary blobs of encrypted data and has no access to any unencrypted information. FHE makes it possible to assert that responses are unforgeable, not merely probabilistically unforgeable.

The use of advanced cryptographic algorithms, such as FHE, has led to the third issue associated with privacy-preserving queries: the question of efficiency. It is known that "1-in-q" algorithms exhibit runtime performance of the order of $O(n)$ or even $O(1)$; thus, the strength of the probabilistic guarantee is only limited by the user's privacy goals. The accuracy of the response can be increased by incurring a linear or even constant time computational overhead. Partially homomorphic encryption schemes, such as routing and spectrum allocation (RSA), may exhibit exponential computational growth. However, extensive experience with such algorithms has led to the development of some implementation heuristics that allow the computations to avoid the exponential neighborhood. Thus, the computations will take polynomial time, not exponential time, except with negligible probability [9], [10]. Efficient computation has not

been realized for FHE, thus far. At present, the implementation technology is immature and even a straightforward FHE computation may take several seconds [11], [12]. Obviously, current FHE performance is not compatible with end-user requirements [13], [14].

Finally, one must consider the requirements specifically imposed by location-based data. The response to an LBS query is not a single data item but an ensemble of data items [15], [16]. Thus, malicious actors, such as a man-in-the-middle eavesdropper, may find it easier to make correlations over multiple queries that lead to an unacceptable loss of privacy. Any privacy preserving solution to the query problem for location-based servers must take such behavior into account [17], [18].

Many techniques have been proposed to provide an optimal solution for privacy-preserving queries in LBS, and there are many ways to create solutions for privacy issues. The ultimate objective is to ensure the efficiency and privacy of LBS and their queries. However, this requirement raises questions regarding the security of servers in general and that of LBS in particular. One of these issues is that servers could be compromised and may become malicious. Moreover, servers may produce queries that cannot have forged responses. If for some reason, there is a forged response, the recipient of that response will be able to recognize that it is a fake response. In [19], the authors proposed mobile online social networks (mOSNs), where a location sharing service was presented to the mOSNs. They examined the current problems of location sharing and proposed BMobishare as an enhanced security mechanism that ensures location privacy. Also, they employed the Bloom filter to provide greater security for sensitive data compared to other existing methods.

To overcome the problems mentioned above, the paper proposes the BLOT approach, which combines oblivious transfer (OT) with Bloom filters (BF). Here, the sender is the LBS, while the receiver is the user. The user initiates a transaction by sending a set of q queries to the LBS. These queries cause a benign LBS to generate q responses, i.e., the messages $m_0 \dots m_{(q-1)}$. The protocol proceeds as above. The benign LBS cannot guess which of the q queries were the relevant queries, except with probability $1/q$; therefore, the expectation is that BLOT will enhance communication security between the client and the server against attacks.

The contributions can be summarized as follows:

- The proposed method involves the creation of an entirely new type of cryptographic communication protocol in a problem domain and the development of a secure implementation that can be used in resource-constrained systems, such as smartphones.
- The amount of information exposed to attack is minimized. Hence, the BLOT transformation slows down information leakage from the BF from 0.5 bits per query to $0.5/N$ bits per query, on average, where N is arbitrarily large. Both OT (transformation 2) and BF can be used to reduce information leakage.

- More entropy means less information leakage. Thus, it is more difficult for an attacker to learn anything by studying the encrypted responses.

The remainder of this paper describes the approach to overcoming the challenges described above using OT and BF. A persuasive argument can be made that these diverse protocols can be successfully combined to yield solutions that are stronger than any individual approach.

II. RELATED WORK

Recently, many solutions have been proposed to protect the users of LBS. Information access control [20] is a technique that provides location privacy to LBS users. Specifically, it equips the LBS provider with a mechanism that enables the LBS users to control access to their location data. Toward this end, LBS providers enforce the access policy's usage to control access to the users' location data. The drawback of this technique is that the LBS providers could be potential adversaries who misuse the location data of the users.

One method is based on a mix zone, i.e., it depends on an intermediate server to hide the user's location [21]. The intermediate server assigns a pseudonym to the user when he/she enters a mix zone. The pseudonym is used by the user to send queries to the LBS server via the intermediate server. A new pseudonym is assigned to the user when his/her mix zone is changed. An important application area of the mix zone technique is road networks. The drawback of this technique is that the intermediate server is vulnerable to adversaries.

The k -anonymity technique [22], [23] is based on the concept that a user cannot be distinguished from other $k-1$ users. Toward this end, mobile users are grouped in clusters of k members. For each group, a bounding region is defined. Each user uses the bounding region of the group as his/her location and attaches this region for all the queries sent to the LBS provider. An intermediate server is required to construct the bounding region. An adversary can identify the location of the user with a maximum probability of $1/k$. The drawback of this technique is that adversaries can compromise the intermediate server.

Using dummy locations is another technique for achieving location privacy [24], [25]. In this technique, the mobile user confuses the LBS server by sending his/her query many times, where one of them contains his/her real location whereas the others contain fake locations. The drawback of this technique is that adversaries can use the side information to analyze the user's sent locations and identify the dummy locations.

In [26], the authors proposed an optimal approach to tackling current location-based query issues regarding the leakage of location data from a query database, generally known as points of interest (POIs). The authors focused on the sharing of location data by the owner with every user. Toward this end, a two-stage approach was proposed; the first step involves oblivious transfer, and the second step involves private information retrieval. The proposed method was implemented on a desktop system and a mobile device with the aim of examining its efficiency.

In [27], the authors proposed the use of oblivious transfer, k -anonymous oblivious transfer, deterministic encryption, and Bloom filters. The k -anonymous oblivious transfer protocol was combined with symmetric key encryption (block cipher) to prevent the server from determining which client data are being used. This approach provides anonymity and is beneficial when privacy is preferred over speed.

In [28], the authors showed that oblivious transfer (OT) is essential regarding cryptography, and it is used in many protocols as a security measure for multi-party computation. There is a need for such protocols in the case of a hands-on application, and OT extensions allow the base OTs to be used to calculate a large amount of OTs at a low cost. Consequently, counterparts that use more efficient protocols are created. This model is safe when used in both the random oracle model and the standard model. If the efficacy of the OT extensions is enhanced, they become even safer.

On the other hand, the paper [29], explored bloom filters and their setbacks. Hash functions provided issues for users, even though bloom filters seemed straightforward at first. This paper proposed a method to create independent and orthogonal hash functions while limiting information leaks.

While in [30], presented a survey of the current trends of privacy methods used to protect users that had used LBSs. This paper showed the effectiveness, or lack thereof, of each technique showcased in this paper. It should be recommended that in the future, researchers should conduct a security analysis, using a simulation setup, which would evaluate their proposed algorithm using a modified Hilbert Curve. This would assure the effectiveness of the privacy in the proposed scheme.

The flexibility of sharing location in privacy protection is key to social networking services. Opaque identities developed by Smoke Screen in 2007 [31] help share presence among trusted friends and untrusted strangers. Although the previous work [32] has resolved the problems of flexibility in location sharing in a privacy protecting way, the strictness of the method prevents a direct use.

Recently a Mobishare mechanism [33], proposed by Wei et al. provides a flexible privacy-preserving location sharing to both trusted social relations and untrusted strangers in mOSNs. However, a Mobishare mechanism is actually an extension of Smoke Screen but use a dummy of technique like a real fake identity that prevents complete identities locations between LBS providers and social network providers. A later improvement by Liu et al. broadcast encryption to preserve users' location privacy allowing users to add or removal of friends [34].

Even with such improvements, multiple attacks on a location-based provider can essentially obtain all the friends' relation of user plus the locations in real time. The Paillier Cryptosystem developed by Li et al. could, however, prevent this problem [35]. Li et al. proposed a mechanism that employs a private set intersection protocol to prevent named Mobishare+, that prevent the social network server from copying individual information. While not much is improved in terms of security, this scheme should achieve a lot more

complex encryption and decryption operations that sustain computation overhead appreciably. There are still other encryption-based methods particularly in a cloud environment that used for privacy protection [36]-[38].

III. PRELIMINARIES

A. Deficiency of BMobishare

A dummy query is used, which allows services to protect their users' real and fake entities.

A Bloom filter serves as a bridge between LBS and social network servers. Although this approach safeguards private data, it is flawed, as there is a risk of leakage between users and servers, which may increase the overhead and decrease the entropy.

B. Oblivious Transfer Scenario in LBS

Create a data transfer protocol between an LBS client and an LBS server with the following properties:

- 1) *The protocol has 1-in- q capability.*
- 2) *The LBS server cannot distinguish information-bearing queries.*
- 3) *Even if an eavesdropper obtains all queries and responses, no information is leaked.*
- 4) *A malicious LBS server can be detected with probability q .*
- 5) *The probability q can be increased exponentially, close to 1, by increasing the queries.*

Oblivious transfer (OT) is a 1-in- q type of algorithm in which the sender transmits a series of messages to the receiver without knowing (except with negligible probability) which, if any, of these messages, contain information used by the receiver, will describe OT for the case of two messages, but it is obviously extendable to the case of q messages. Note that OT is typically built on top of a public key cryptosystem, such as RSA [39], [40], although there is no compelling reason why a system based on symmetric cryptography could not also be used [41].

The steps of the protocol are as follows. The sender has two messages, m_0 and m_1 , and wishes to provide both messages to the receiver without being able to determine which message is actually used by the receiver. The sender generates an RSA key pair. Let M denote the modulus of this key pair; p_e , the public exponent; and p_d , the private exponent. The sender generates two random nonces, x_0 and x_1 , and sends them to the receiver along with the RSA public key. The receiver generates a random bit b and uses it to set $x = x_b$. The receiver then generates a large random large number L and computes the value $V = (x + Lp_e) \bmod M$ [42]. Then, the value V is transmitted to the sender.

The sender cannot know (except with random probability) the value of b . Therefore, the sender computes two values, $v_0 = (V - x_0)p_d \bmod M$ and $v_1 = (V - x_1)p_d \bmod N$. The sender can deduce that the value of L is either v_0 or v_1 [43]. The RSA cryptosystem is based on the presumed intractability of the discrete logarithm problem. Thus, the sender cannot know (except with random probability) which of these two values is

correct. The sender computes $m_0' = m_0 + v_0$ and $m_1 = m_1 + v_1$, and sends both to the receiver. The receiver knows the value of b and is thus able to determine which of m_0' and m_1' is valid. Note that no eavesdropper can obtain any information from the value pair (m_0', m_1') except with negligible probability. In the 1-in- q generalization, the receiver has a secret index, i , and can receive a set of q messages from the sender such that the sender has no way of determining the value of i .

In the OT implementation of an LBS, the sender is the LBS, and the receiver is the user. The user initiates a transaction by sending a set of q queries to the LBS. These queries cause a benign LBS to generate q responses, i.e., the messages $m_0 \dots m_{(q-1)}$. The protocol proceeds as above. The benign LBS cannot guess which of the q queries are the relevant queries, except with probability $1/q$. If N exchanges are performed, rather than one, the benign LBS cannot guess which is the relevant query, except with probability $1/qN$ [44]. Thus, by adjusting q and N , the user can achieve the desired degree of certainty. Note that a malicious LBS server will fail to produce a meaningful response, except with random probability; hence, in the worst case, a malicious LBS server can use the OT protocol to launch a denial-of-service attack. Such an attack can be readily detected by the user with probability $(1 - 1/qN)$, with as much certainty as the user desires.

C. Bloom Filter Scenario in LBS

A protocol between a client, an LBS server, and a second server must have the following properties:

- 1) The client and second server can send a query if there is something in the range of the LBS server.
- 2) The original server can only determine one bit of information from a list of queries, and the second server cannot determine any of these data.
- 3) No set of objects in the range of the server needs to be itemized.
- 4) Any malicious LBS server can be detected with probability q .
- 5) The probability q can be exponentially increased close to 1 by increasing the queries.

The Bloom filter (BF) allows testing of a fixed membership and never reveals more than one bit of information regarding this set. If A is a reputed element of set S , then the BF creates a probabilistic algorithm to determine whether A is an element of S without going through the elements of S [45].

The LBS includes a Bloom filter that uses location information encoded as bit vectors. The user creates q pairs (query, assertion). The queries will be sent to the LBS, which will create a bit vector of q values. If the i -th position of this vector is 0, then the pair (Q_i, A_i) is absolutely inconsistent; if the i -th position of this vector is 1, then the pair (Q_i, A_i) is probabilistically consistent [46]. If the server is benign, the evaluations will be accurate. However, if the server is malignant, the evaluations will be inaccurate. Because users can create pairs (Q, A) with a known state, the server will perform tests against known values. A false positive means that

the server is malicious. False positives can be reduced with N iterations of q queries. Thus, a malicious server can be found with probability $(1 - p^N)$, where p is the legitimate false positive. Malicious servers cannot learn anything from these exchanges, except with a negligible probability due to the inverse H_k .

D. BLOT Overview in LBS

The ultimate objective of both protocols is to reduce information leakage. Both BF and OT can be used. We achieve this by applying a two-stage approach shown in Fig. 1. The first stage is based on Bloom filter and the second stage is based on an oblivious transfer.

E. Security for LBS: The BLOT protocol

This paper describes a new approach for LBS applications that provides a user with three security guarantees: (a) it will take at least Q_1 queries for an eavesdropper to learn a single bit of information from the protocol exchange; (b) it will take at least Q_2 queries for a malicious LBS provider to learn a single bit of information from the protocol exchange; and (c) the information leaked by the database handling the LBS information is at most Q_3 bits for each query. In these statements Q_1 , Q_2 , and Q_3 are adjustable parameters based on the specific implementation of the protocol. If the total number of queries that the user must make is A , and if the protocol can be adjusted such that $A < Q_1$, $A < Q_2$, and $(A * Q_3) < 1$, then the security guarantees become absolute: no eavesdropper or malicious LBS provider can learn any information from the protocol exchange, except with negligible probability.

First, describe the component technologies. BLOT approach is in the form of a framework. The LBS provider will always use one particular storage mechanism: Bloom filters. The protocol will be a pluggable protocol that can take any form as long as it satisfies the security guarantees stated above. In this paper, will consider the OT transfer protocol. Then, will show how it can be integrated to create the BLOT protocol. The construction of the Bloom filter hash functions enables us to demonstrate security guarantee (c). Next, the precise steps in the BLOT exchanges are described in detail, with full mathematical justification, allowing us to demonstrate the first two security guarantees.

In the remainder of this paper, the entity requesting LBS information will be known as the user or the receiver, whereas the entity giving out LBS information will be known as the provider or the sender. The terms client and server will not be used because the BLOT protocol does not require client/server architecture to be deployed. Toward the end of this paper, will explain how it is possible to deploy the BLOT protocol entirely on a mobile device, entirely within the LBS provider, or in some combination of both.

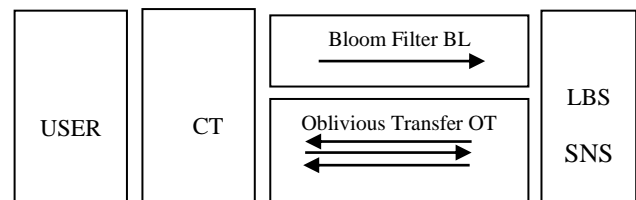


Fig. 1. BLOT overview in LBS.

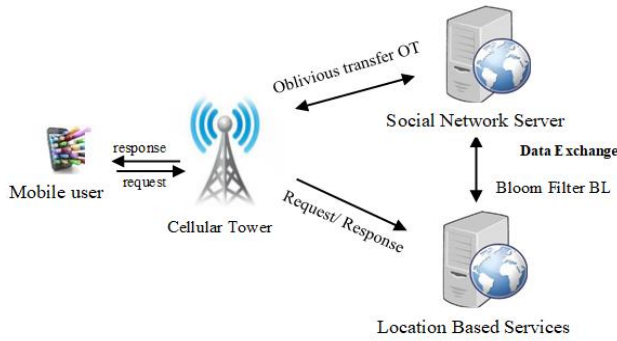


Fig. 2. BLOT architecture in LBS.

F. BLOT System Architecture

As shown in Fig. 2, the scenario of LBS consists of four entities in the mobile online network.

- The mobile device user can access the Internet with wireless technology (such as 3G/4G) and then share his or location and inquire about the location of friends who are nearby.
- SNS manages the user’s identity-related information (user profiles, friend lists, etc.).
- LBS stores the user’s anonymized location information and provides LBS, as per the user’s request, with the real-time locations of people nearby.
- CT aids the user’s communication with SNS and LBS.

BLOT assumes that SNS, LBS, and CT are connected with links that are high-speed and secure.

G. OT Standard

The purpose of OT is to securely transfer information from a sender to the receiver such that both sides learn the minimum amount of information about the requested transfer. Specifically, suppose that the sender has N pieces of information, $(I_0 \dots I_{N-1})$, and the receiver has an index n in the range $(0, N)$. The receiver wishes to receive I_n without the sender knowing the value of n . The sender also wishes to ensure that the receiver will only be able to decrypt one of the received N encrypted messages. Both sides do not want an eavesdropper to learn any information. OT is a type of minimal knowledge protocol [47].

OT is typically built on top of some form of public key cryptography. The RSA protocol is typically used. The RSA’s security is based on the inversion of the discrete logarithm in polynomial time with negligible probability [48]. The algorithm below [49] shows the standard form for OT (note that this exchange shows the protocol for a single sender and a single receiver, as is typical for LBS transactions; however, OT can be extended to multiparty communications).

- 1) The Sender generates RSA key pair, with the public modulus MO and the public exponent e .
- 2) The Sender generates N random numbers r_i .
- 3) The Sender sends MO , e , and all the r_i to the Receiver.

- 4) The Receiver selects the random number r_n where n is the index of the desired message.
- 5) The Receiver generates a random number z and computes: $v = (r_n + ze) \bmod MO$.
- 6) The Receiver sends v , known as the RSA blind value of r_n , to the Sender.
- 7) The Sender computes all N values $v_i = (v - r_i) d \bmod MO$, where d is the private exponent.
- 8) The Sender computes all $I'_i = I_i + v_i$ and sends all I'_i to the Receiver.
- 9) The Receiver computes the correctly decrypted value $I_n = I'_n - z$.

The Sender has a 1-in- N chance of guessing the value n . Thus, from a security standpoint, both parties would like to make N as large as possible. However, from a performance standpoint, the transfer time increases linearly with the value of N . LBS often involve a large number of data transfers. Hence, the performance issue can be partially offset by modification of the messaging protocol; nevertheless, linear performance degradation will occur as the security of the system enhances.

Although we can assume that the receiver (the user) is honest, the same may not be true for the sender (the LBS provider). A dishonest sender could deliberately send weak RSA values, as well as weak random numbers, to improve the odds of guessing the value n . The selected system architecture is summarized in Table I.

The selected notations are summarized in Table II.

In the next section, will describe the BLOT protocol, which significantly improves the security of the system by combining BF with OT.

TABLE I. SYSTEM ARCHITECTURE

Symbol	Description
User (receiver)	Mobile user
LBS	Location-based server (backing storage db)
SNS (sender)	Social network server
CT	Cellular tower
ID	An authorized user's unique ID
qf	Distance threshold in friends' locations query
df	User's friend-case distance
qd1	Friends' locations query 1
qd2	Friends' locations query 2

TABLE II. NOTATIONS

Symbol	Description
ID_A	User A's unique social network identifier
$PubMO_A$	Public modulus
$Pub e_A$	Public exponent
ID_{CT}	Cellular tower's identifier
$PubKey$	Public key encryption-decryption
$Skey (AES)$	Symmetric key encryption-decryption
ds	User's strange-case distance
BF_f	Bloom filter with inserted friend-list

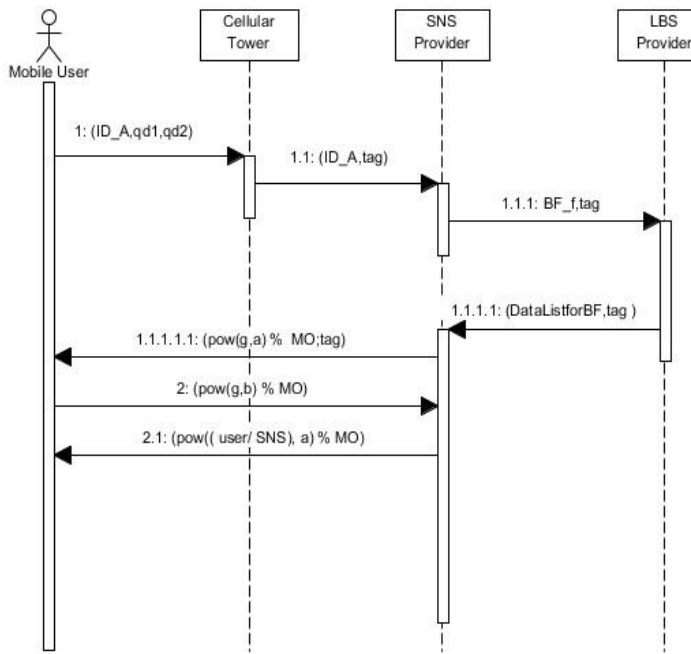


Fig. 3. BLOT Querying friends' locations.

IV. BLOT MECHANISM

A. BLOT Querying Friends' Locations

OT builds a security mechanism between users and the SNS, which ensures privacy while avoiding leaks. BF creates a security mechanism between LBS and SNS, and these protect each other's privacy. Fig. 3 shows methods for discovering friends' locations:

Mobile User: First, the user must submit a friend's location query (IDA, qd1, qd2) to CT.

CT: The appropriate entry $\langle IDA, dfA \rangle$ for IDA will be read in its table. The sequence number and identifier are appended, and the query is forwarded as (IDA, tag) to SNS, where tag = (IDCT, qd1, qd2, seq) and seq is a sequence number.

B. Bloom Filter Stage BL

1) First, BF initializes itself and updates itself using the file (DataListforBF.txt).

2) This file contains the list of words that will be set in BF for fast searches in the data structure.

a) In the filter, Add Element(line) method in the initializing function will operate adding the BF element.

b) It uses two hash functions, "FNV" and "MURMUR," to generate the two indices for an element and sets the bit to 1 at the selected indices [50].

3) Now, with another file, "DataToQuery.txt.". This file has the list of words that the receiver (USER) is going to query from the SNS.

4) Then, the SNS sends the query to the LBS and asks about the two words (QueryData1"qd1" and QueryData2"qd2).

5) The LBS takes these 2 words and searches in the BF. It again generates the indices using the same FNV and MURMUR hash functions and checks whether the bit is set to 1.

6) If it does not find the data in the BF, then it sends a message to the receiver (USER): "Element (QueryData) is not here."

Otherwise, if it finds the data in the filter, it confirms the same by searching for the data in the LBS having backing storage. The backing storage, in this case, is the same file DataListforBF.txt. It has been used here as a database.

7) If both sets of data are found in the backing storage in the LBS, they are forwarded to the SNS, and then the process of OT is initiated for transferring the desired encrypted data to the mobile user.

C. OT Stage

1) For OT, consider the following parameters:

```
int g = 7; // public exponent
int MO = 101; // public modulus
int N = 2; // number of encrypted packets to send
int C = 0; // indicates that message to be selected by sender SNS (it could be 0 or 1)
```

2) The Sender (SNS) generates a random number 'a' between 1 and 5.

3) Based on the generated random number, SNS calculates the following and sends it to the user:

Double Sender = (long long int) pow(g,a) % MO;

4) After receiving the SNS input, the Receiver (user) performs the following mathematical operations and sends it to the SNS:

```
int b = rand() % 5 + 5; // generate random number between 5 and 10
```

```
if (C == 0) Receiver(user) = (long int) pow(g,b) % MO;
```

```
else if (C == 1) Receiver(user) = Sender SNS * ((long int) pow(g,b) % MO);
```

5) The Sender (SNS) does not know which message the Receiver (user) is going to decrypt. Hence, it will encrypt both pieces of data and send it to the user.

a) SNS generates two different keys using the "sha256" key generator.

b) Using these keys, it encrypts both data1 and data2 and sends it to the user. For encryption, it uses "AES256-ECB"

6) After receiving the data, the Receiver (user) generates a key based on the value of "SNS" received from the user. The key is generated using the same "sha256" key generator.

7) The Receiver (USER) then uses this key to decrypt data. The user will only be able to decrypt one of the two datasets correctly.

Fig. 3 shows how friends' locations are queried for BLOT sequences in LBS.

V. BLOT PERFORMANCE MEASURES

We have already discussed the three performance measures (Q1, Q2, and Q3) that will use as measures for the security of the system. Also introduce a functional performance measure, namely the performance entropy (PE). Both of the protocols have packets that flow between the user and the LBS server as well as in the reverse direction.

Entropy is related to the length (L) of the messages [51], we can divide the contents of each packet into message-carrying information and non-message-carrying information. If the total length of all packets in a single exchange is L, then we can write $L = C + N$, where C is the sum of the lengths of the message-carrying subset, and N is the sum of the lengths of the non-message-carrying subset. Then, define the performance entropy of the protocol exchange as

$$PE = N/L \tag{1}$$

The larger the value of PE, the higher is the entropy. Thus, for a practical system, want this measure to be as small as possible. For example, if all data was transmitted in the clear, then $N = 0$ and $PE = 0$; there is no entropy. Of course, this situation is completely insecure [51], so it is important to minimize this performance measure for the two encrypted protocols described in this paper.

VI. PERFORMANCE ENTROPY CALCULATION OF BMOBISHARE

To calculate the performance entropy of BMOBishare, considered the same scenario as that considered in BLOT, i.e., querying a friend's location can be divided into two steps: service registration and authentication, and querying a friend's location. The selected sizes of information (bits) are summarized in Table III.

A. Service Registration and Authentication

This scenario involves some message exchanges between the user, SNS, and CT, as shown in the message sequence chart (MSC). Based on MSC and Table III, the number of bits exchanged is calculated first, and then PE is calculated. Fig. 4 shows the MSC.

Service registration and authentication involve 6192 bits of information exchanged, which is considered to be non-message-carrying information.

B. Query a Friend's Location

Querying a friend's location involves 7936 bits of information exchanged between User, CT, SNS, and LBS, as shown in Fig. 5, shows how friends' locations are queried.

The sizes of message-carrying information (C) and non-message-carrying information (N) are listed in Table IV.

TABLE III. SIZE OF INFORMATION (BITS)

S. No.	Information	Size of information (bits)
1	PubKey	2048
2	Skey (AES)	128
3	df _A	64
4	ds _A	64
5	ID _A	64
6	Time stamp (ts)	64

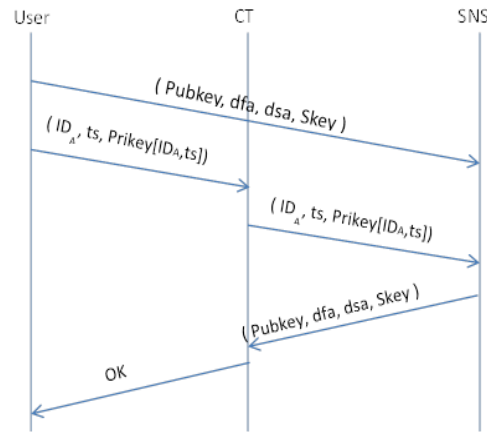


Fig. 4. Message Sequence Chart (MSC).

$$L = C + N = 14128 \text{ bits} \tag{2}$$

$$P_1 = \frac{C}{L} = \frac{1600}{14128} = 0.11325 \tag{3}$$

$$P_2 = \frac{N}{L} = \frac{12528}{14128} = 0.886 \tag{4}$$

By using (13), (shown in the Appendix at the end of the paper), we can calculate entropy as [51]:

$$H_{BMOBishare} = (0.11325) \log_2 \left(\frac{1}{0.11325} \right) + (0.886) \log_2 \left(\frac{1}{0.886} \right) = 0.5105 \tag{5}$$

While,

$$H_{BLOT} (Avg) \text{ has been calculated from simulation founds to } be = 0.9756 \tag{6}$$

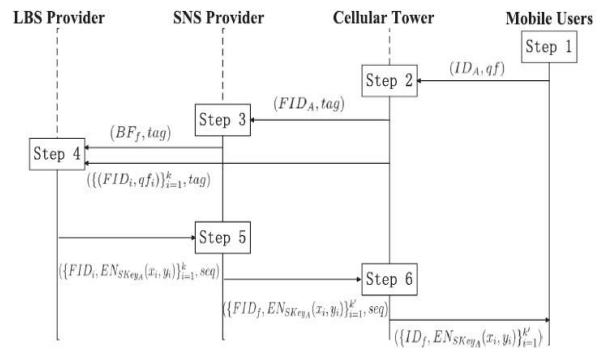


Fig. 5. Querying friends' locations.

TABLE IV. THE TOTAL SIZE OF NON-/MESSAGE-CARRYING INFORMATION

S. No.	Function	C	N
1	Service registration and authentication	0	6192
2	Query friends' locations	1600	6336
	Total	1600	12528

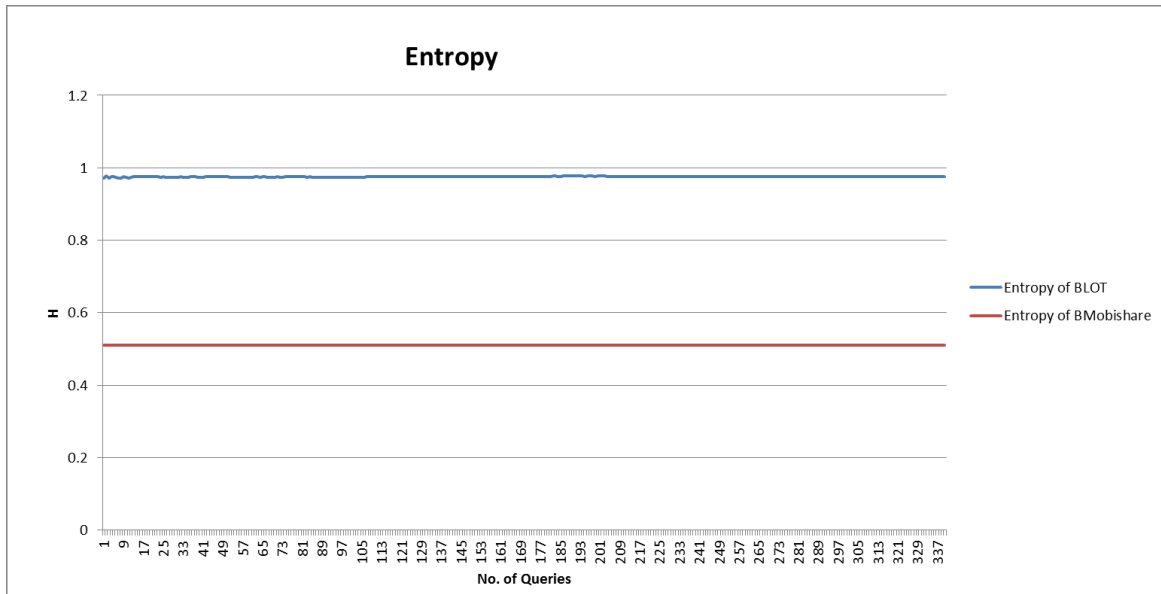


Fig. 6. Performance entropy comparison.

VII. ANALYSIS OF RESULTS AND DISCUSSION

The total number of information exchanges in BLOT is around 960 bits as per simulation, which is much smaller than that of BMobishare. Fig. 6 shows the performance entropy comparison. The average performance entropy of BLOT is 0.9756, which is greater than that of BMobishare 0.5105, because BMobishare uses a large number of anonymous messages for increasing security, whereas BLOT ensures security by using an encryption algorithm (AES256).

Based on the result, we can realize a significant entropy value, which means less information leakage. Further, can confirm the effectiveness of the BLOT mechanism as well as low overhead and high performance owing to the avoidance of massive numbers of fake IDs, in contrast to BMobishare.

VIII. PERFORMANCE RESULTS

Note: In the given paper on BMobishare, the author did not mention what he means by “connection performance” or “computing performance”. The conclusion below is based on the assumption that “connection performance” and “computing performance” have the following definitions:

A. Connection Performance

The time a sender takes to establish a connection to the receiver to start further secure communication using any protocol. Connection performance of BLOT would be better than that of BMobishare because BMobishare establishes the connection before querying the friend’s location using the service registration and authentication mechanism. In BLOT there is no need for service registration and authentication. Hence connection time is 0s. In BMobishare, average connection time is 1.5s.

Therefore, we can conclude that the connection performance of BLOT is better while computing performances are comparable.

B. Computing Performance

The time a computer takes to execute a code/algorithm/protocol. Computing performance can be calculated based on the number of packets communicated in each protocol because a higher number of packets involved means more computation required to process them.

a) In BLOT, a total of 7 messages need to be communicated to query a friend’s location.

b) In BMobishare, 7 messages are communicated to query a friend’s location.

C. Computing Performance Measurement

In BLOT average computation time measured from the simulation is 10 ms. while the BMobishare shows that the average computing time is 1.75s.

But the author of BMobishare measured the time with a processor of 1.2 GHz and 1GB RAM.

While in BLOT have measured time over 2.5 GHz and 4GB RAM. As the protocol code is not very large, RAM cannot be a factor in execution time. 1GB RAM is also enough to accommodate the code data as well. Here we can see that the BLOT processor is twice the speed of the BMobishare processor. Therefore to make the comparison fair, we can assume that if run BLOT on a 1.2 GHz machine, it will double the execution time. (That is, it would multiply the time by a factor of 2. 10ms x2 = 20ms can give us the time if the code runs on a 1.2 GHz machine). The following table shows the performances result in Table V.

TABLE V. PERFORMANCE RESULTS IN SUMMARY

S. No.	Protocol	Computation time (ms)	Connection time (s)
1	BLOT	20	0
2	BMobishare	1.75	1.5

IX. CONCLUSION

In this paper, is described the BLOT mechanism, which combines two protocols that attempt to address the security deficiencies commonly found in current LBS application environments. In particular, BLOT addresses the security of backing storage information, information that is visible to an eavesdropper, or information that is vulnerable to a malicious presence on the LBS provider machine. Also, it enables a user to perform consistency checks on the LBS provider.

Each component individually provides high security, while the deployment of both components significantly enhances the security using either transfer approach. Analyzed the performance entropy of BLOT and found it greater than the solution by BMobishare; additionally, the performance was more efficient.

X. APPENDIX

A. Entropy Derivation

By the definition of Shannon's entropy [51]:

$$H(x) = -\sum_{i=1}^n P_i(x) \log_2(P_i(x)) \quad (7)$$

$$H(x) = \sum_{i=1}^n P_i(x) [\log_2(1) - \log_2(P_i(x))] \quad (8)$$

$$H(x) = \sum_{i=1}^n P_i(x) \log_2\left(\frac{1}{P_i(x)}\right) \quad (9)$$

In BLOT, case $n=2$

$$i = 1, \text{ for message acrrying information} \quad (10)$$

$$i = 2, \text{ for non - message acrrying information} \quad (11)$$

Therefore the above equation becomes:

$$H(x) = \sum_{i=1}^2 P_i(x) \log_2\left(\frac{1}{P_i(x)}\right) \quad (12)$$

$$H(x) = P_1(x) \log_2\left(\frac{1}{P_1(x)}\right) + P_2(x) \log_2\left(\frac{1}{P_2(x)}\right) \quad (13)$$

$$P_1 = \text{probability message acrrying information} = \frac{C}{L} \quad (14)$$

$$P_2 = \text{probability nonmessage acrrying information} = \frac{N}{L} \quad (15)$$

Equation (13) has been used to calculate the entropy.

REFERENCES

- [1] B. Niu, Z. Zhang, X. Li, and H. Li, "Privacy-area aware dummy generation algorithms for location-based services," in Proc. Of IEEE ICC 2014.
- [2] J. Krumm, "A survey of computational location privacy," Pers. Ubiquit. Comput., vol. 13, no. 6, pp. 391-399, Aug. 2009.
- [3] Through spatial and temporal cloaking," in Proc. Of ACM MobiSys 2003.
- [4] A. Albelaihy, J. Cazalas, "A Fine-Grained Spatial Cloaking With Query Probability Levels for Privacy in LBS," International Journal of Computer Networks and Applications (IJCNA), vol. 2, no. 5, pp. 212-221, 2015.
- [5] C. Cachin, C. Crepeau, J. Marcil and G. Savvides, "Information-Theoretic Interactive Hashing and Oblivious Transfer to a Storage-Bounded Receiver," IEEE Trans. Inform. Theory, vol. 61, no. 10, pp. 5623-5635, 2015.
- [6] N. Aggarwal, C. Gupta, and I. Sharma, "Fully Homomorphic symmetric scheme without bootstrapping," in Cloud Computing and Internet of Things (CCIOT), 2014 International Conference on, 2014, pp. 14-17.
- [7] J. Alwen, M. Barbosa, P. Farshim, R. Gennaro, S. D. Gordon, S. Tessaro, et al., "On the relationship between functional encryption, obfuscation, and fully homomorphic encryption," in Cryptography and Coding, ed: Springer, 2013, pp. 65-84.
- [8] V. Garg and M. Jhamb, "A Review of Wireless Sensor Network on Localization Techniques," International Journal of Engineering Trends and Technology (IJETT)-Volume4Issue4-April, 2013.
- [9] C. Gentry, J. Groth, Y. Ishai, C. Peikert, A. Sahai, and A. Smith, "Using fully homomorphic hybrid encryption to minimize non-interactive zero-knowledge proofs," Journal of Cryptology, pp. 1-24, 2014.
- [10] N. Smart, "Investigations of Fully Homomorphic Encryption (IFHE)," DTIC Document2015.
- [11] Z. Zhang, T. Plantard, and W. Susilo, "Reaction attack on outsourced computing with fully homomorphic encryption schemes," in Information Security and Cryptology-ICISC 2011, ed: Springer, 2012, pp. 419-436.
- [12] A. Leiva, N. Pavez, A. Beghelli, and R. Olivares, "A joint RSA algorithm for dynamic flexible optical networking," in Communications (LATINCOM), 2014 IEEE Latin-America Conference on, 2014, pp. 1-6.
- [13] G. D. Sutter, J.-P. Deschamps, and J. L. Imaña, "Modular multiplication and exponentiation architectures for fast RSA cryptosystem based on digit serial computation," Industrial Electronics, IEEE Transactions on, vol. 58, pp. 3101-3109, 2011.
- [14] G. Asharov, A. Jain, A. López-Alt, E. Tromer, V. Vaikuntanathan, and D. Wichs, "Multipart computation with low communication, computation and interaction via threshold FHE," in Advances in Cryptology-EUROCRYPT 2012, ed: Springer, 2012, pp. 483-501.
- [15] Z. Brakerski and V. Vaikuntanathan, "Lattice-based FHE as secure as PKE," in Proceedings of the 5th conference on Innovations in theoretical computer science, 2014, pp. 1-12.
- [16] S. Devadas, M. van Dijk, C. W. Fletcher, and L. Ren, "Onion ORAM: A Constant Bandwidth and Constant Client Storage ORAM (without FHE or SWHE)," 2015.
- [17] M. Naehrig, K. Lauter, and V. Vaikuntanathan, "Can homomorphic encryption be practical?," in Proceedings of the 3rd ACM workshop on Cloud computing security workshop, 2011, pp. 113-124.
- [18] S. Mudda and S. Giordano, "Mobile P2P queries over temporal data," in Pervasive Computing and Communications Workshops (PERCOM Workshops), 2014 IEEE International Conference on, 2014, pp. 278-283.
- [19] N. Shen, J. Yang, K. Yuan, C. Fu and C. Jia, "An efficient and privacy preserving location sharing mechanism", Computer Standards & Interfaces, vol. 44, pp. 102-109, 2015.
- [20] M. Youssef, V. Atluri, and N. R. Adam, Preserving mobile customer privacy: An access control system for moving objects and custom proles. s.l. : In Proc. MDM 2005.
- [21] Mobimix, B. Palanisamy and L. Liu, Protecting location privacy with mix-zones over road networks. s.l. : In Proc. ICDE 2011.
- [22] B. Bamba, L. Liu, P. Pesti, and T. Wang, Supporting anonymous location queries in mobile environments with PrivacyGrid. s.l. : In Proc. WWW 2008.
- [23] C.-Y. Chow, M. F. Mokbel, and W. G. Aref, "Casper*: Query processing for location services without compromising privacy," ACM Trans. Database Syst., vol. 34, no. 4, 2009.
- [24] P. Shankar, V. Ganapathy and L. Iftode, Privately querying location-based services with SybilQuery. s.l. : In Proc. Ubicomp 2009.
- [25] O. Han, H. Zhao, Z. Ma, K. Zhang, and H. Pan, Protecting Location Privacy Based on Historical Users over Road Networks .. s.l. : In Proc. WASA 2014.
- [26] R. Paulet, M. G. Kaosar, X. Yi, and E. Bertino, "Privacy-preserving and content-protecting location based queries," Knowledge and Data Engineering, IEEE Transactions on, vol. 26, pp. 1200-1210, 2014.

- [27] V. Gupta, T. S. Vineeth, and V. Aggarwal, "Make Your Query Anonymous With Oblivious Transfer," in Proceedings of the Sixth International Conference on Computer and Communication Technology 2015, 2015, pp. 345-349.
- [28] G. Asharov, Y. Lindell, T. Schneider, M. Zohner: More Efficient Oblivious Transfer Extensions with Security for Malicious Adversaries. EUROCRYPT (1) 2015: 673-701.
- [29] A. Albelaihy and J. Cazalas, "Privacy preserving queries for LBS: Hash function secured (HFS)," 2017 2nd International Conference on Anti-Cyber Crimes (ICACC), Abha, 2017, pp. 7-12.
- [30] A. Albelaihy and J. Cazalas, "A survey of the current trends of privacy techniques employed in protecting the Location privacy of users in LBSs," 2017 2nd International Conference on Anti-Cyber Crimes (ICACC), Abha, 2017, pp. 19-24.
- [31] L.P. Cox, A. Dalton, V. Marupadi, Smokescreen: flexible privacy controls for presence-sharing, Proceedings of the 5th International Conference on Mobile Systems, Applications and Services, ACM 2007, pp. 233-245.
- [32] K.P. Puttaswamy, B.Y. Zhao, Preserving privacy in location-based mobile social applications, Proceedings of the Eleventh Workshop on Mobile Computing Systems & Applications, ACM 2010, pp. 1-6.
- [33] W. Wei, F. Xu, Q. Li, Mobishare: flexible privacy-preserving location sharing in mobile online social networks, INFOCOM, 2012 Proceedings IEEE, IEEE 2012, pp. 2616-2620.
- [34] Z. Liu, J. Li, X. Chen, J. Li, C. Jia, New privacy-preserving location sharing system for mobile online social networks, P2P, Parallel, Grid, Cloud and Internet Computing (3PGCIC), 2013 Eighth International Conference on, IEEE 2013, pp. 214-218.
- [35] J. Li, J. Li, X. Chen, Z. Liu, C. Jia, Mobishare+: security improved system for location sharing in mobile online social networks, J. Internet Serv. Inf. Secur. (JISIS) 4 (2014) 25-36.
- [36] Z. Liu, J. Li, X. Chen, J. Yang, C. Jia, Tmds: thin-model data sharing scheme supporting keyword search in cloud storage, Information Security and Privacy, Springer 2014, pp. 115-130.
- [37] J. Li, Z. Liu, X. Chen, F. Xhafa, X. Tan, D.S. Wong, L-encdb: a lightweight framework for privacy-preserving data queries in cloud computing, Knowl.-Based Syst. 79 (2015) 18-26.
- [38] Z. Liu, X. Chen, J. Yang, C. Jia, I. You, New order preserving encryption model for outsourced databases in cloud environments, J. Netw. Comput. Appl. (2014),
- [39] D. Boneh, A. Sahai, and B. Waters, "Functional encryption: a new vision for public-key cryptography," Communications of the ACM, vol. 55, pp. 56-64, 2012.
- [40] A. J. Menezes, Elliptic curve public key cryptosystems vol. 234: Springer Science & Business Media, 2012.
- [41] K. B. Swamy and K. K. Raju, "Multi-keyword Ranked Search over Encrypted Cloud Data Using RSA Algorithm," IJSEAT, vol. 3, pp. 307-311, 2015.
- [42] S. Bai, R. Brent, and E. Thomé, "Root optimization of polynomials in the number field sieve," Mathematics of Computation, 2015.
- [43] V. Grolmusz, "Separating the communication complexities of MOD m and MOD p circuits," in Foundations of Computer Science, 1992. Proceedings., 33rd Annual Symposium on, 1992, pp. 278-287.
- [44] O. Chowdhury, D. Garg, L. Jia, and A. Datta, "Equivalence-based Security for Querying Encrypted Databases: Theory and Application to Privacy Policy Audits," in Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, 2015, pp. 1130-1143.
- [45] A. Broder and M. Mitzenmacher, "Network applications of bloom filters: A survey," Internet mathematics, vol. 1, pp. 485-509, 2004.
- [46] W. Zhang, S. Liu, and Y. Xiaoyuan, "RLWE-Based Homomorphic Encryption and Private Information Retrieval," in Intelligent Networking and Collaborative Systems (INCoS), 2013 5th International Conference on, 2013, pp. 535-540.
- [47] R. Ahlswede, "Founding Cryptography on Oblivious Transfer," in Hiding Data-Selected Topics, ed: Springer, 2016, pp. 337-344.
- [48] D. R. Stinson, Cryptography: theory and practice: CRC press, 2005.
- [49] M. O. Rabin. How to exchange secrets by oblivious transfer. Technical Report TR-81, Harvard Aiken Computation Laboratory, 1981.
- [50] "Oblivious transfer", Asecuritysite.com, 2017. [Online]. Available: <https://asecuritysite.com/encryption/ot>. [Accessed: 26- Jul- 2017].
- [51] M. Eigen, From strange simplicity to complex familiarity. Oxford: Oxford Univ. Press, 2013.

Development of Mobile-Interfaced Machine Learning-Based Predictive Models for Improving Students' Performance in Programming Courses

Fagbola Temitayo Matthew
Department of Computer Science
Federal University, Oye-Ekiti, Nigeria

Obe Olumide
Department of Computer Science
Federal University of Technology, Akure

Adeyanju Ibrahim Adepoju
Department of Computer Engineering
Federal University, Oye-Ekiti, Nigeria

Olaniyan Olatayo, Esan Adebimpe, Omodunbi Bolaji
Department of Computer Engineering
Federal University, Oye-Ekiti, Nigeria

Oloyede Ayodele
Department of Computer Science
Caleb University, Imota, Lagos, Nigeria

Egbetola Funmilola
Department of Computer Science & Engineering
LAUTECH, Ogbomosho, Nigeria

Abstract—Student performance modelling (SPM) is a critical step to assessing and improving students' performances in their learning discourse. However, most existing SPM are based on statistical approaches, which on one hand are based on probability, depicting that results are based on estimation; and on the other hand, actual influences of hidden factors that are peculiar to students, lecturers, learning environment and the family, together with their overall effect on student performance have not been exhaustively investigated. In this paper, Student Performance Models (SPM) for improving students' performance in programming courses were developed using MSP Decision Tree (MDT) and Linear Regression Classifier (LRC). The data used was gathered using a structured questionnaire from 295 students in 200 and 300 levels of study who offered Web programming, C or JAVA at Federal University, Oye-Ekiti, Nigeria between 2012 and 2016. Hidden factors that are significant to students' performance in programming were identified. The relevant data gathered, normalized, coded and prepared as variable and factor datasets, and fed into the MDT algorithm and LRC to develop the predictive models. The developed models were obtained, validated and afterwards implemented in an Android 1.0.1 Studio environment. Extended Markup Language (XML) and Java were used for the design of the Graphical User Interface (GUI) and the logical implementation of the developed models as a mobile calculator, respectively. However, Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), Relative Absolute Error (RAE) and the Root Relative Squared Error (RRSE) were the metrics used to evaluate the robustness of MDT and LRC models. The evaluation results obtained indicate that the variable-based LRC produced the best model in terms of MAE, RMSE, RAE and the RRSE having yielded the least values in all the evaluations conducted. Further results obtained established the strong significance of attitude of students and lecturers, fearful perception of students, erratic power supply, university facilities, student health and students' attendance to the performance of students in

programming courses. The variable-based LRC model presented in this paper could provide baseline information about students' performance thereby offering better decision making towards improving teaching/learning outcomes in programming courses.

Keywords—*Student-performance; predictive-modeling; MSP-Decision-Tree; mobile-interface; linear-regression-classifier; programming-courses*

I. INTRODUCTION

Computer programming courses are a fundamental part of many Universities' curricula and among the most important subjects for computer science and information technology students. This requires the knowledge of programming tools and languages, problem-solving skills and effective strategies for program design and implementation [1]. Furthermore, students are being exposed to various programming specifications and techniques which normally entails an overview of algorithms, concept of programming, basic data structure, problem analysis and illustrations describing the application of various techniques to problems which are quite difficult to understand [2]. Furthermore, the high level of abstraction and very complex language syntax and semantic structures induced in programming makes it a much dreaded task in which most students fail [2]. This is evidenced by the notion that the same set of students who failed programming courses performed better in other non-programming courses [3]. As a matter of fact, the failure rate in programming courses at the University level suggests that learning to program is a difficult task [3]. The perception of the complexity ascribed to programming courses can be described as one of the main reasons that may have attributed to the decline in number of undergraduates who offer or intend to offer computer science in various institutions [4].

Chermahini [5] noted that students are different based on their ability to learn, how they respond to instructional practices, their motivational differences from one individual to another and that the more students understand the differences in their abilities, the better are the chances they have to meet their different learning needs in order to achieve good scores in examinations. Students' performance is majorly affected by several social, economic, institutional, environmental, psychological and personal factors which vary across individuals and regions [6]-[8]. Unfortunately, poor performances have ravaged the academic institutions due to indices of those factors which influence students' performance including poor funding, lack of frequent curricular review, overpopulation, students' unrest, staff strikes, poor facilities, coarse relations between the university and government, inadequate teaching and research facilities needed to enhance students' learning and performance. More specifically, Ogbogu [6] and Irfan and Shabana [9] emphasized that challenges such as poorly equipped departmental and central libraries, overcrowded lecture rooms, method of collating and accessing semester results, interruption of electricity supply, poor access to internet facilities, incessant strike and closure of school and poor accommodation facilities which are pertinent to developing countries affect student performance.

Students' performance assessment has become a pressing issue that requires fair attention from all regardless of differences in interest and intentions [9], [10]. However, different methods have been used to evaluate students' performance, and more than ever before, information generated by evaluation can be helpful for students and tutors to take timely, meaningful and effective decisions. Most existing student performance models have adopted statistical techniques for prediction which are probability-induced, depicting that results may not be scientifically correct but rather are based on estimation. To this end, several authors have adopted data mining and soft computing techniques in educational domain and/or to evaluate students' performance [11]-[17].

Ashish, Saeed, Maizatul, and Hamidreza [14] focused on consolidating the different types of clustering algorithms been applied within the context of Educational Data Mining (EDM) to harnessing the power of the massive didactic data recently being generated in institutions. EDM was employed to analyze data generated in an educational setup by the various intra-connected systems in a bid to develop a model for improving learning and institutional effectiveness. Among the slightly numerous clustering algorithm consolidated by the authors are Expectation Maximization, Hierarchical Clustering, Simple k -Means and x -Means, Apriori Algorithm (as applied to academic records of students in a guise to obtain the best association rules which helps in student profiling), C -Means clustering, Ward's clustering, Markov Clustering (MCL) algorithm, Unique Clustering with Affinity Measure (UCAM), Fuzzy sets, Transitive Closure and a hierarchical cluster analysis which was performed on the questionnaire data. As concluded by these authors, data mining methods in the educational sector sets to uncover the previously hidden data

to meaningful information that can be used for strategic and learning gains.

Kolo, Adepoju and Alhassan [18] aimed at predicting the performance of students with the decision tree approach. Gurmeet and Williamjit [13] employed data-mining approach for an effective prediction of student performance based on personal, social, psychological and environmental variables. This was to ensure a high accuracy in the prediction of student performance, thereby assisting to identify students with low academic achievements. The parameters employed in the study include gender, hometown, family income, previous semester grade, attendance, communication language (medium), seminar performance and participation in sports. Analysis of these parameters was conducted by implementing the algorithms in WEKA tool. Naïve Bayes and J48 algorithms were used for classification and the result showed that the Naive Bayes algorithm provided an accuracy of 63.59% while the J48 algorithm provided an accuracy of 61.53%.

Generally, the educational sector in developing countries is being faced by a series of multi-factored challenges that contribute to the rapid decline in the performance of students located within such contemporary environments. Teachers and students alike have for so long been unable to estimate the impact that certain factors have on academic performances but rather anticipate good performances in the long run. This way, it becomes impossible for student to quickly re-adjust and retune performance demeaning challenges surrounding them or probably their responses to such surrounding factors. More often than not, the actual influences of hidden factors that are peculiar to students, lecturers, learning environment and the family, together with their overall effect on student performance have not been exhaustively investigated in existing studies.

In this paper, M5P decision tree and linear regression classifier, which are among the most widely adopted machine learning techniques, are employed to develop the student performance predictive models. Metrics used to evaluate the performance of the machine learning techniques employed include mean absolute error, root mean squared error, relative absolute error and the root relative squared error, correlation coefficient, time taken to build the model and the time taken to test the model.

The major contributions of this paper are as follows:

a) Exhaustively investigated, examined, identified and established new hidden factors and associated variables on which students' performance in programming courses is dependent and that are particularly peculiar to a prototype University in a developing economy. These are significant and technical extensions beyond most student performance models that currently exist;

b) Beyond the spheres of statistical approaches commonly used for student performance modeling which are based on probability and estimation in most existing works, this study applied machine learning techniques (M5P Decision Tree and Linear Regression Classifier) to predicting

student performance in programming courses to guarantee precision and accuracy of the resultant predictive models;

c) Towards facilitating the accessibility, availability and ubiquity of the developed predictive models, a mobile application, that visually interfaces the stakeholders and all student performance indices with the models, was developed. This is to realize real-time use in predicting students' performance and for promoting effective and efficient decision making on education planning by all stakeholders.

The rest of this paper is organized as follows: Section 2 discusses the materials and method including the M5P decision tree and linear regression classifier, data acquisition, the development and validation of the machine learning-based predictive models and the performance evaluation metrics for the machine-learning based approaches. In Section 3, the design and implementation of the mobile-frontend application for the developed predictive models are presented and discussed. The results of performance evaluation of the machine learning approaches are presented and discussed in Section 4 while the conclusion and future works are presented in Section 5.

II. MATERIALS AND METHOD

In this research, models for predicting students' performance in programming courses were developed based on M5P and linear regression classification algorithms in three basic steps. These include data acquisition, development of the predictive models and finally model validation. Furthermore, the performance evaluation of the machine learning approaches employed and the mobile implementation of the predictive models developed were conducted.

A. The Classification Algorithms

1) M5P Decision Tree: This is a decision tree model that learns regression tasks. The M5P learns efficiently and can cope with highly-dimensional data with up to several hundreds of distinct attributes. According to Quinlan [19], M5P decision tree is the most accurate among the family of regression tree learners with much smaller model trees than regression trees. It uses mean squared error as the impurity function. A M5P tree is constructed by recursive partitioning of a data into a collection of set T which can either be associated with a leaf or a split function that segregates T into some subsets based on some split function criteria [20]. The subsets that emerge are further partitioned following the same process repeatedly. However, the quality of split (goodness of fit) is evaluated using a function $\emptyset(S, T)$ where S is the split candidate in node T such that the split candidate that maximizes the value of quality of fit is selected as the next node of tree [21]. That is, $\emptyset(S, T) = \Delta I(S, T) = I(T) - p_L I(T_L) - p_R I(T_R)$ (1) where $I(T)$ is the impurity function at node T for k classes in a dataset defined as:

$$I(T) = I(p(1|T), p(2|T), \dots, p(k|T)) \quad (2)$$

p_L and p_R are the probabilities that an instance is going to the left branch and right branch of T according to split S , $p(j|T)$ is the estimated posterior probability of class j given a

point in node T , $\Delta I(S, T)$ is the difference between the impurity measure of node T and two child nodes T_L, T_R according to split S . The information gain in M5P is determined by the difference in the values of standard deviation obtained before and after the split function test. Simply put, given data T , where T_i denotes the subsets of T corresponding to the i_{th} outcome of a split function test, then the expected error reduction value is determined by Hieu [22]:

$$\Delta error = sd(T) - \sum_i \frac{|T_i|}{|T|} sd(T_i) \quad (3)$$

The split function test criterion that maximizes this expected error reduction is then selected. To avoid overfitting, subtrees that do not improve the performance of the tree are pruned via an error-based estimation procedure, from the leaves to the root node [23]. This is determined by the difference in the estimated error of a node and estimated error of the subtree below at each internal node.

2) Linear Regression Classifier: The linear regression classifier is a mathematical measure depicting the mean relationship among two or more variables based on the original units of the data [24]. This often involves the estimation and prediction of an unknown value of one variable from the known value of another variable [25]. This implies that there exists a linear regression between the variables should the regression curve be a straight line. With linear regression, the values of the dependent variable increase by a constant absolute amount for a unit change in the value of the independent variable. However, the general form of linear regression measure is given as [26]:

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n = \theta^T(x) \quad (4)$$

where $h_{\theta}(x)$ is a straight line with variable intercept

and parameters $\theta_0, \theta_1, \dots, \theta_n$, if $x_0 = 0$ is assumed.

Algorithm: Linear Regression Classification [27]

Inputs: Class models $X_i \in R^{q \times p_i}$, $i = 1, 2, \dots, N$ and a test input student performance factors' vector $y \in R^{q \times 1}$.

Output: Class of y

- i. For each class model, $\hat{\beta}_i \in R^{p_i \times 1}$ is evaluated such that $\hat{\beta}_i = (X_i^T X_i)^{-1} X_i^T y$, $i = 1, 2, \dots, N$
- ii. \hat{y}_i is computed for each $\hat{\beta}_i$, $\hat{y}_i = X_i \hat{\beta}_i$, $i = 1, 2, \dots, N$
- iii. Distance between original and predicted response variables is determined by $d_i(y) = \|y - \hat{y}_i\|_2$, $i = 1, 2, \dots, N$
- iv. Decision is made with regard to the class that has the minimum distance $d_i(y)$

B. Data Acquisition

Hidden factors that are significant to student performance were identified via a thorough literature review, interview and field observations. Questionnaire was developed for the University under study with respect to information on programming courses and associated scores as presented at the Appendix section. In Table I, the contextual definition of the variables is presented. Copies of the questionnaires were disseminated to students that had offered programming courses and their respective lecturers in the University.

Relevant data were gathered, normalized and coded. The coded data was utilized by the machine learning techniques to develop the student performance models and were further validated for prediction purpose.

TABLE I. VARIABLE EXPRESSION FROM DESIGNED QUESTIONNAIRE

S/N	Expressions
x_1	I had enough time to study programming
x_2	Studying before attending a class aided my assimilation during programming classes.
x_3	Studying programming was never a wasted effort
x_4	Programming sounded very scary
x_5	I was always nervous during programming classes
x_6	I was always nervous during programming examinations
x_7	I attended programming classes regularly
x_8	Blending in after missing a class was very easy
x_9	I was very serious with programming classes
x_{10}	I believed I could understand the programming course
x_{11}	I had interest in programming beyond class level
x_{12}	Programming was not confusing and did not cause headache
x_{13}	Programming is relevant to my pursuit
x_{14}	Group discussions helped me to understand programming
x_{15}	Attending programming tutorials was very helpful
x_{16}	Programming courses tutorials helped me so much
x_{17}	Motivation of programming lecturers encouraged my commitment towards learning programming
x_{18}	Programming language lecturers helped me develop interest in programming
x_{19}	Programming languages lecturers were never partial in their dealings with students
x_{20}	Programming lecturers were friendly during lectures
x_{21}	Programming language lecturers enforced discipline during their lectures
x_{22}	Programming languages lecturers were too serious during lectures
x_{23}	Teaching methods and styles of programming lecturers inhibited lecture clarity
x_{24}	Programming language lecturers wasted time on matters with less relevance in class
x_{25}	Programming language lecturers were always clear, precise and communicates understandably
x_{26}	Programming language lecturers made use of enough relevant instructional materials
x_{27}	Programming language lecturers delivered course contents well and to my understanding
x_{28}	Programming language lecturers were very clear and explicit
x_{29}	Programming language lecturers didn't miss classes
x_{30}	Programming language lecturers attended to me whenever I had difficulties with their course(s)
x_{31}	Programming lecturers were always available
x_{32}	Programming course lecturers allowed students to ask questions and take time to explain
x_{33}	Programming course lecturers came to class fully prepared
x_{34}	Programming languages lecturers spent extra time to explain things during class
x_{35}	Programming language lecturers usually came early to class
x_{36}	I fell sick quite often
x_{37}	Prolong usage of computer caused me headache
x_{38}	I took a few compulsory medications frequently
x_{39}	It was difficult to charge my computer even within the campus
x_{40}	Erratic power supply reduced the effectiveness of my practice
x_{41}	Consistent power supply helped me in programming courses
x_{42}	I had a good background in physics
x_{43}	I had a good background in mathematics
x_{44}	I had a good background in English
x_{45}	Strong background in Physics and Mathematics helped me in programming

x_{46}	Absence of accessible ICT facilities inhibited my programming performance
x_{47}	The environment where we had programming lectures was not conducive
x_{48}	Lack of computer programming facilities disrupted clear understanding of programming lessons
x_{49}	The school library was not equipped with materials relevant to programming
x_{50}	Large class population disrupted my concentration during programming lectures
x_{51}	Population of students offering programming courses debarred my commitment to learning
x_{52}	Effectiveness of the programming lecturers' teaching was reduced by huge programming class population.
x_{53}	Programming lectures were scheduled after an equally tiring lecture
x_{54}	Programming courses were scheduled to non-conductive times
x_{55}	We had programming classes at unfavorable times
x_{56}	Programming lecture theatres were equipped with audio-visuals and learning aids
x_{57}	Programming courses were analyzed clearly to sight
x_{58}	I had a visual understanding of what the programming lecturer was implying
x_{59}	Expensive cost of living did not affect my performance in programming classes
x_{60}	My family could afford to buy enough programming textbooks
x_{61}	My family sponsored my academic pursuit
x_{62}	Quarrel between family members is normal
x_{63}	I had to travel to settle quarrels within my family
x_{64}	Quarrel between my family members escalates a times
x_{65}	My father is familiar with computers
x_{66}	My mother is familiar with computers
x_{67}	My parents are well educated
x_{68}	My parent would want me to offer programming courses
x_{69}	I received educational advices from family members often
x_{70}	My family believed that a proper study will help me in programming courses

However, twenty-one (21) factors were investigated via this study with a total of 81 variables. Each factor was coded based on the cumulative of the variables designated to investigate it as conducted by Fagbola *et al.* [11]:

a) *Student Study Habit (SSH)*: This is the amount of the student's effective study in programming courses offered relative to the frequency of revision and practice and hours spent on revising the lecture notes. It was investigated by three variables x_1, x_2, x_3 .

b) *Student Fear and Perception (SF)*: This is the students' fearful perception of programming courses where a positive perception implies a reduction in fear factor of the student. This was investigated by the variables x_4, x_5, x_6 .

c) *Student Attendance (SATD)*: This is the level of effort, seriousness and devotion of students towards learning to program, investigated by the variables x_7, x_8, x_9 .

d) *Student Attitude (SAT)*: This is the level of responsiveness of a student relative to their interest, behavior and seriousness to programming courses, and characterized by student's participation in class activities, assignment, willingness to learn, and motivation from friends, colleagues and lecturer(s). This was represented by the variables $x_{10}, x_{11}, x_{12}, x_{13}$.

e) *Tutorials and Extra Classes (ST)*: These are the extra effort put in place by students in other to have a clear

understanding of the subject matter(s) discussed programming classes. This includes extra-classes attended, assistance from friends and use of online forums and materials. This factor was investigated by the variables x_{14} , x_{15} , x_{16} .

f) *Lecturer Attitude (LAT)*: This is defined as the lecturers' assertiveness, interest to explicitly expatiate on the subject matter, ability to motivate the student and relate with the student in a means to improve their interest in the course. This was investigated by variables x_{17} , x_{18} , x_{19} , x_{20}

g) *Teaching Style (LTS)*: This is defined as the pattern of teaching of the lecturer in charge (probably dishes out voluminous handouts or excessive assignments). Whether he carries the class along and helps the student conceptualize the concept of that particular programming course. This was investigated by variables x_{21} , x_{22} , x_{23} , x_{24} .

h) *Communication Skills (LCS)*: This is the ability of the lecturer to deliver the course content in a less ambiguous manner and to the understanding of the students. This entails the clarity and explicitness of the lecturer. This was investigated by variables x_{25} , x_{26} , x_{27} , x_{28} .

i) *Lecturer Availability (LA)*: This is the presence and accessibility of the lecturers' when they are needed by the student(s). This factor was investigated by the variables x_{29} , x_{30} , x_{31} .

j) *Lecturer Dedication (LD)*: This is the devotion of the lectures to the programming courses they tutor. This includes the assertiveness of the lecturers to their duty and extra effort put in place to ensure an excellent student performance. This factor was coded as presented in Table III and was investigated by the variables x_{32} , x_{33} , x_{34} , x_{35} .

k) *Health (OH)*: This is the influence of medical condition on students' performance in programming courses. This factor was coded and was investigated by the variables x_{36} , x_{37} , x_{38} .

l) *Electricity (OE)*: This is defined as the erraticism of power supply as it affects the students' practice using computers and also other laboratory works. This factor was coded and was investigated by the variables x_{39} , x_{40} , x_{41} .

m) *Background knowledge (OB)*: This is the academic strength of the student in other courses that are elementarily related to computer programming (mathematics and physics). This factor was investigated by the variables x_{42} , x_{43} , x_{44} , x_{45} .

n) *Facilities (UF)*: This is the availability of appropriate programming learning facilities (computer laboratory) within the university environment. This factor was investigated by the variables x_{46} , x_{47} , x_{48} , x_{49} .

o) *Class population (UCP)*: This is the student to tutor population ratio during the programming course class. This factor was investigated by the variables x_{50} , x_{51} , x_{52} .

p) *Lecture time (ULT)*: This is the conduciveness of the lecture schedule. This factor was investigated by the variables x_{53} , x_{54} , x_{55} .

q) *Teaching aids (UTA)*: This is the availability of teaching aids (audio visuals) for the demonstration of the

concept of programming courses. This factor was investigated by the variables x_{56} , x_{57} , x_{58} .

r) *Family income (FI)*: This is the robustness of the family income of the student. As it influence the ability of the student to afford textbook materials, print handout or even own a personal computer for effective study. This factor was investigated by the variables x_{59} , x_{60} , x_{61} .

s) *Family stress (FS)*: This is the degree of disturbance from home. An unsettled home creates a paranoid atmosphere which seemly affects student performance. This factor was investigated by the variables x_{62} , x_{63} , x_{64} .

t) *Parent education (FPE)*: This is the degree of education of the students' parent. A poor motivation from home might destabilize the student cognitive sense, hence influencing the students' performance in programming. This factor was investigated by the variables x_{65} , x_{66} , x_{67} .

u) *Proper guidance (FPG)*: This is the student's family guidance and support level for programming courses. A student from a family of computer scientist is prone to having huge support and guidance from home. This factor was investigated by the variables x_{68} , x_{69} , x_{70} .

After final normalization and cleaning process were completed, the entire data acquired was divided into variable and factor datasets and each data split was used to train the machine learning classifiers.

C. Development of the Machine learning-based Student Performance Predictive Models

M5P decision tree and the linear regression classifier, having industrially-packaged working implementations in WEKA environment, were trained using the variable and factor datasets and further applied to generate predictive models which are of exclusive significance to the determination of students' performance. The variable-based student performance model generated by the linear regression classifier is presented in (5).

$$\begin{aligned} x_{80} = & 0.0444 * x_1 + 0.3166 * x_2 + 0.0746 * x_3 - 0.0415 * \\ & x_4 - 0.239 * x_5 + 0.3153 * x_6 - 0.1467 * x_7 + 0.3464 * \\ & x_8 + 0.6227 * x_9 - 0.1404 * x_{11} - 0.3228 * x_{12} + \\ & 0.1179 * x_{13} - 0.4613 * x_{14} - 0.3948 * x_{15} + 0.4249 * \\ & x_{16} - 0.2241 * x_{17} - 0.1389 * x_{18} + 0.2025 * x_{19} + \\ & 0.0664 * x_{20} + 0.133 * x_{21} + 0.1745 * x_{22} - 0.3222 * \\ & x_{23} - 0.3334 * x_{24} - 0.2479 * x_{25} - 0.1623 * x_{26} + \\ & 0.0665 * x_{28} - 0.2556 * x_{29} + 0.2829 * x_{30} - 0.2215 * \\ & x_{31} - 0.4575 * x_{33} + 0.135 * x_{34} + 0.3312 * x_{35} - \\ & 0.2152 * x_{36} + 0.2407 * x_{37} + 0.1757 * x_{38} - 0.2986 * \\ & x_{39} + 0.1768 * x_{40} - 0.2375 * x_{41} - 0.1969 * x_{42} + \\ & 0.2352 * x_{43} - 0.098 * x_{44} + 0.4561 * x_{45} - 0.136 * \\ & x_{46} - 0.387 * x_{47} + 0.1525 * x_{48} - 0.2215 * x_{49} + \\ & 0.0481 * x_{50} + 0.1292 * x_{51} + 0.1508 * x_{52} + 0.4368 * \\ & x_{53} - 0.3313 * x_{54} - 0.1794 * x_{55} - 0.0523 * x_{56} - \\ & 0.3505 * x_{57} + 0.4718 * x_{58} + 0.269 * x_{59} + 0.086 * \\ & x_{60} - 0.3004 * x_{61} - 0.444 * x_{62} + 0.3544 * x_{63} - \\ & 0.2301 * x_{64} - 0.538 * x_{65} + 0.0899 * x_{66} + 0.2394 * \\ & x_{67} - 0.0681 * x_{68} - 0.1007 * x_{69} - 0.3858 * x_{70} + \\ & 9.8865 \end{aligned} \quad (5)$$

The learned models developed are further used to generate predictions on new instances. The factor-based Student

Performance Model obtained using linear regression classifier is expressed in (6).

$$grade = -0.074 * sf + 0.0942 * satd + 0.065 * sat + 0.0449 * lat - 0.0448 * lcs - 0.0407 * la + 0.0493 * oh + 0.0814 * oe - 0.0792 * uf + 0.0621 * fi - 0.0663 * fs - 0.0533 * fpe - 0.1233 * fpg + 5.6703 \quad (6)$$

The M5 pruned model tree for the variable dataset is presented in Fig. 1. However, the variable-based M5P decision tree classifier generated smoothed Linear Models (LM) through 22 refinement processes. The first and the last generated models are presented in (7) and (8), respectively

although the latest refinement was used to predict student performance.

$$x80 = 0.0297 * x1 + 0.0187 * x4 - 0.0376 * x5 + 0.1263 * x9 - 0.017 * x12 - 0.0826 * x14 + 0.021 * x15 + 0.0316 * x19 - 0.0209 * x22 + 0.0389 * x57 + 0.0211 * x59 - 0.0343 * x65 - 0.0217 * x69 + 3.9539 \quad (7)$$

$$x80 = 0.0155 * x1 + 0.0098 * x4 - 0.0652 * x5 + 0.0552 * x7 + 0.1046 * x9 - 0.0089 * x12 - 0.0143 * x14 + 0.011 * x15 - 0.0503 * x19 - 0.1112 * x22 + 0.032 * x29 - 0.0288 * x33 + 0.0324 * x59 - 0.06 * x65 - 0.188 * x69 + 5.9906 \quad (8)$$

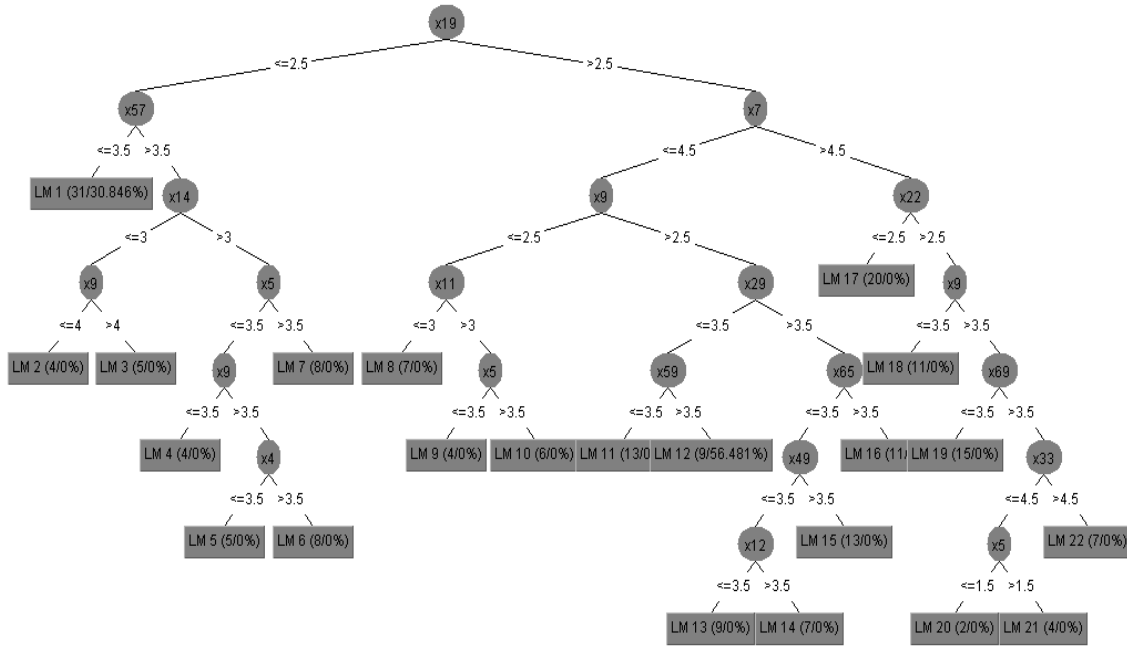


Fig. 1. The M5 pruned model tree for the variable dataset.

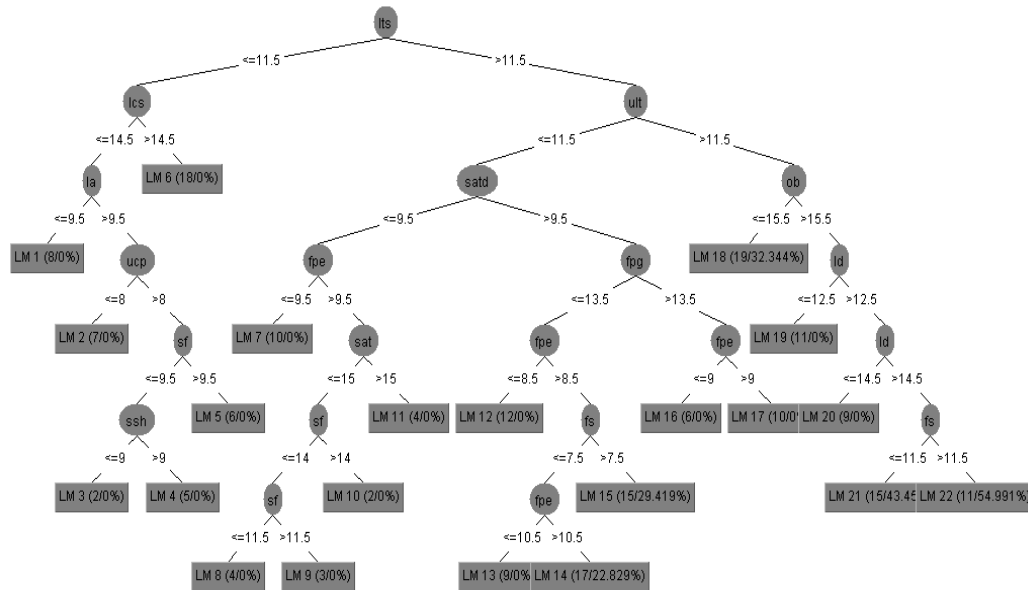


Fig. 2. The M5 pruned model tree for the factor dataset.

The M5 Pruned model tree for the factor dataset is presented in Fig. 2. However, the factor-based M5P classifier generated smoothed Linear Models (LM) through 22 refinement processes. The first and the last models generated are presented in (9) and (10), respectively.

$$grade = 0.0481 * ssh + 0.1057 * sf + 0.0343 * satd + 0.0084 * sat - 0.0083 * st + 0.0127 * lat + 0.0475 * lcs - 0.1963 * la + 0.0141 * oe + 0.0232 * ucp - 0.0248 * fs - 0.0097 * fpe - 0.0293 * fpg + 5.0805 \quad (9)$$

$$grade = -0.0144 * sf + 0.0307 * satd + 0.0106 * sat - 0.0101 * st + 0.0045 * lat - 0.0548 * lcs + 0.0202 * ld + 0.0273 * oe + 0.0675 * ob - 0.0159 * ult - 0.0466 * fs - 0.0034 * fpe - 0.0321 * fpg + 4.0297 \quad (10)$$

D. Validation of the Developed Machine Learning-based Student Performance Predictive Models

The variable and factor datasets were employed in the development of the students’ performance predictive models, which were then validated using the test dataset. Some instances of the validation results of the predictive models generated by the machine learning classifiers are presented in Table II. It is important to note that with limited data used for validation, the results of validation test cannot be exclusively used to justify the correctness of the developed models but rather by some standard evaluation measures. Based on some validation results obtained, the best performing model is the factor dataset-based SPM generated by the linear regression classifier. This is followed by variable dataset-based SPM generate by M5P decision tree classifier, factor dataset-based M5P decision tree and the variable dataset-based SPM based on linear regression classifier in decreasing order of performance. Note that the best prediction values are marked in “bold”.

TABLE II. MODEL VALIDATION INSTANCES FOR LINEAR REGRESSION AND M5P DECISION TREE CLASSIFIERS

Actual Grade	Linear Regression Classifier Algorithm (variable dataset)-based SPM	Linear Regression Classifier Algorithm (factor dataset)-based SPM	M5P Decision Tree Classifier (variable dataset)-based SPM	M5P Decision Tree Classifier (factor dataset)-based SPM
4	4.3618	4.0124	3.865004	4.0347
6	6.2135	5.9675	5.96883	5.9505
4	4.2946	4.1055	4.036288	4.1578
6	5.0878	5.9583	5.375602	5.2558
5	4.6443	5.0572	4.774742	4.4751
5	5.1855	4.9381	4.881071	5.2582
6	5.3058	5.8879	6.184321	5.8878
5	4.8282	4.8246	4.146855	4.8255
6	5.7855	6.5766	5.697118	5.9423
6	4.3962	6.039	5.271766	5.3175

E. Performance Evaluation Metrics for the Machine Learning-based Approaches Used

The mean absolute error, root mean square error, relative absolute error, root relative squared error, time taken to build and test the models are the standard metrics used to evaluate the performance of the learning techniques.

a) *Root Relative Squared Error (RRSE)* is determined using the relation:

$$RRSE = \sqrt{\frac{\sum_{j=1}^n (P_{(ij)} - T_j)^2}{\sum_{j=1}^n (T_j - \bar{T})^2}} \quad (11)$$

where $P_{(ij)}$ represents the predicted value by each individual program i for any sample case j which is a subset of n sample cases, T_j is the target value for sample case j ; and \bar{T} is given by [28]:

$$\bar{T} = \frac{1}{n} \sum_{j=1}^n T_j \quad (12)$$

b) *The Relative Absolute Error, RAE*, accepts the total absolute error and divides it with the actual absolute error of the model predictor. Relative Absolute Error is determined using the relation [24]:

$$RAE = \frac{\sum_{j=1}^n |P_{(ij)} - T_j|}{\sum_{j=1}^n |T_j - \bar{T}|} \quad (13)$$

c) *Mean Absolute Error, MAE*, is determined by adding the absolute values of the error, e_i , and then dividing the total error by n [24]:

$$MAE = \frac{1}{n} \sum_{i=1}^n |e_i| \quad (14)$$

d) *Root Mean Square Error*: This is a measure of the differences between the sample values predicted by a model and those which are actually observed from the system that is being modelled [28]. That is, the change between the model performance of a predictive model and another. Analytically,

$$RMSE = \sqrt{MSE} \quad (15)$$

where $MSE = \sum_{i=1}^n \frac{(\mu_i - \hat{z}_i)^2}{n}$ such that \hat{z}_i is the model-predicted response for input x_i .

e) *Time taken to build the model*: This is the total time required to learn the discriminating features and to develop a model

f) *Time taken to test the model*: This is the time taken to validate and ascertain the correctness of the developed model.

III. THE DESIGN AND IMPLEMENTATION OF A MOBILE FRONT-END APPLICATION FOR THE DEVELOPED PREDICTIVE MODELS

The developed student performance models were implemented within an Android 1.0.1 Studio environment, using XML for the design of the Graphical User Interface (GUI) and Java for the logic that unifies the GUI and the implementation of the developed models. The flowchart representation for the implementation of the developed student performance models is presented in Fig. 3. The code and design interface is presented in Fig. 4. In the same vein, the mobile home interface of the SPM implementation as presented in Fig. 5 defines the model(s) to be applied and

serves as a link to the questioning aspects of the application. Students and stakeholders can predict the performance of a student by selecting any of the options presented on the home activity of the application. Each of these options implement an underlying model which is used for the prediction of student performance relative to their responses to questions presented.

The interface presented in Fig. 6 displays various questions which are relevant to the selected prediction perspectives. Responses to these questions are then interlinked with the underlying models. In Fig. 7, the predicted performance of the student is displayed in an alert message-box after the responses from prospective students and educational stakeholders have been substituted into the chosen model(s). This happens upon clicking the finish button which appears after the entire questions required for the prediction of student performance under the selected perspective has been duly responded to.

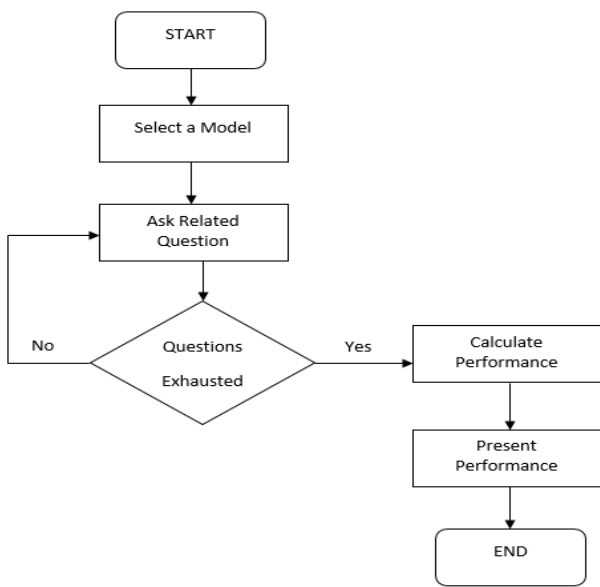


Fig. 3. Flow control of the implementation of student performance models.

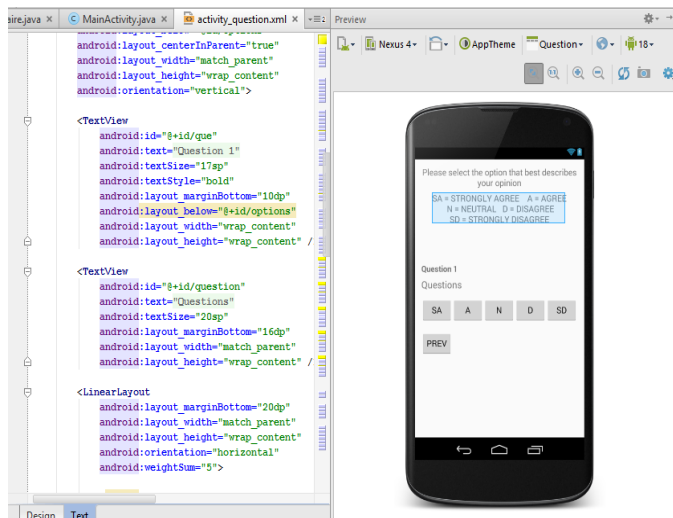


Fig. 4. Code and design interface of the student performance models.

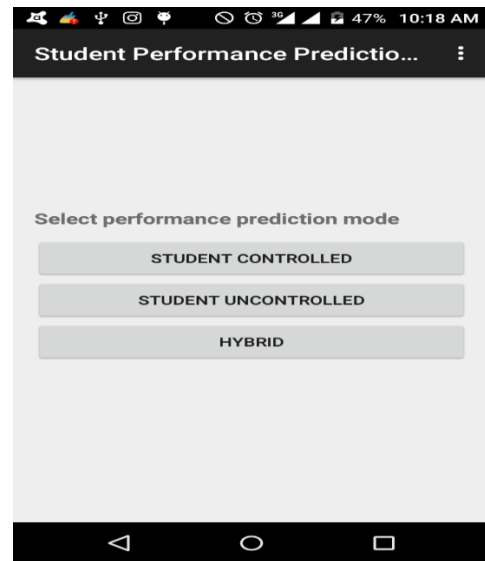


Fig. 5. Home interface of the mobile student performance evaluator.

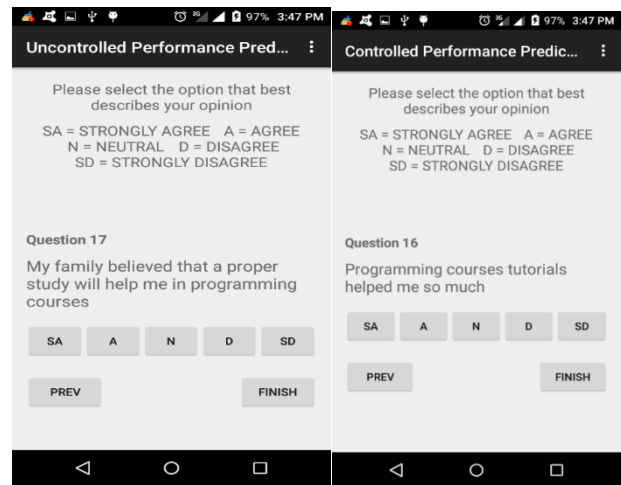


Fig. 6. Interface of the implemented SP models.

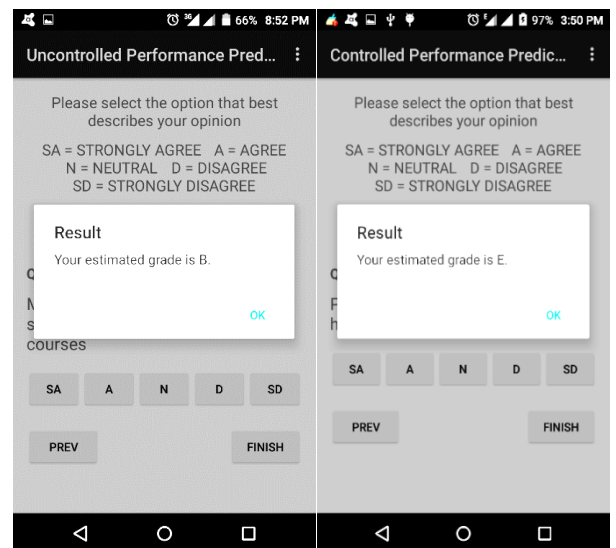


Fig. 7. Instances of predicted students' performance.

IV. RESULTS AND DISCUSSION

In this section, the performance and comparative evaluation results of the machine-learning predictive approaches and the developed student performance models are presented and discussed.

A. Results of Performance Evaluation of the Machine Learning Methods

The results regarding the mean absolute error, root mean square error, relative absolute error, root relative squared error, time taken to build and test the models for both linear regression and M5P decision tree classifiers are presented in Table III. The variable-based Linear Regression Classifier produced the best model in terms of mean absolute error, root mean squared error, relative absolute error and the root relative squared error having yielded the least values in all these metrics. This is followed by the variable-based M5P decision tree, factor-based M5P Decision Tree and the factor-based linear regression classifiers in decreasing order of performance. In terms of the time to build the model, the results obtained indicate that the factor-based M5P Decision Tree is the most computationally-efficient classifier followed by variable-based Linear Regression classifier, variable-based M5P decision tree and factor-based linear regression classifier.

TABLE III. MODEL VALIDATION INSTANCES FOR LINEAR REGRESSION AND M5P DECISION TREE CLASSIFIERS

Technique	Mean Absolute Error	Root Mean Square Error	Relative Absolute Error (%)	Root Relative Squared Error (%)	Time taken to build model (s)	Time taken to test model (s)
Linear Regression Classifier (variable-based)	0.1638	0.2386	20.307	24.8369	0.09	0.03
Linear Regression Classifier (factor-based)	0.5853	0.7273	72.5498	75.7246	0.25	0.06
M5P Decision Tree Classifier (Variable-based)	0.3054	0.4067	41.0867	47.2537	0.13	0.02
M5P Decision Tree Classifier (Factor-based)	0.3984	0.555	53.6099	64.4848	0.05	0.01

Using the model produced by the best performing classifier (variable-based LRC), three (3) out of the 70 variables investigated are found to be insignificant to student performance as presented in Table IV. However, there are 32 variables with positive significance and 35 variables with negative significance to student performance in programming courses as presented in Tables V and VI, respectively.

TABLE IV. VARIABLE-BASED LRC' SPM VARIABLES WITH INSIGNIFICANT EXPRESSIONS

S/N	Insignificant Expressions
x_{10}	I believed I could understand the programming course
x_{27}	Programming language lecturers delivered course contents well and to my understanding
x_{32}	Programming course lecturers allowed students to ask questions and take time to explain

TABLE V. VARIABLE-BASED LRC' SPM VARIABLES WITH POSITIVE EXPRESSIONS

S/N	Expressions with Positive Significance
x_1	I had enough time to study programming
x_2	Studying before attending a class aided my assimilation during programming classes.
x_3	Studying programming was never a wasted effort
x_6	I was always nervous during programming examinations
x_8	Blending in after missing a class was very easy
x_9	I was very serious with programming classes
x_{13}	Programming is relevant to my pursuit
x_{16}	Programming courses' tutorials helped me so much
x_{19}	Programming courses' lecturers were never partial in their dealings with students
x_{20}	Programming courses' lecturers were friendly during lectures
x_{21}	Programming courses' lecturers enforced discipline during their lectures
x_{22}	Programming courses' lecturers were too serious during lectures
x_{28}	Programming courses' lecturers were very clear and explicit
x_{30}	Programming courses' lecturers attended to me whenever I had difficulties with their course(s)
x_{34}	Programming courses' lecturers spent extra time to explain things during class
x_{35}	Programming courses' lecturers usually came early to class
x_{37}	Prolong usage of computer caused me headache
x_{38}	I took a few compulsory medications frequently
x_{40}	Erratic power supply reduced the effectiveness of my practice
x_{43}	I had a good background in mathematics
x_{45}	Strong background in Physics and Mathematics helped me in programming
x_{48}	Lack of computer programming facilities disrupted clear understanding of programming lessons
x_{50}	Large class population disrupted my concentration during programming lectures
x_{51}	Population of students offering programming courses debarred my commitment to learning
x_{52}	Effectiveness of the programming lecturers' teaching was reduced by huge programming class population.
x_{53}	Programming lectures were scheduled after an equally tiring lecture
x_{58}	I had a visual understanding of what the programming lecturer was implying
x_{59}	Expensive cost of living did not affect my performance in programming classes
x_{60}	My family could afford to buy enough programming textbooks
x_{63}	I had to travel to settle quarrels within my family
x_{66}	My mother is familiar with computers
x_{67}	My parents are well educated

TABLE VI. VARIABLE-BASED LRC' SPM VARIABLES WITH NEGATIVE EXPRESSIONS

S/N	Expressions with Negative Significance
x ₄	Programming sounded very scary
x ₅	I was always nervous during programming classes
x ₇	I attended programming classes regularly
x ₁₁	I had interest in programming beyond class level
x ₁₂	Programming was not confusing and did not cause headache
x ₁₄	Group discussions helped me to understand programming
x ₁₅	Attending programming tutorials was very helpful
x ₁₇	Motivation of programming lecturers encouraged my commitment towards learning programming
x ₁₈	Programming language lecturers helped me develop interest in programming
x ₂₃	Teaching methods and styles of programming lecturers inhibited lecture clarity
x ₂₄	Programming language lecturers wasted time on matters with less relevance in class
x ₂₅	Programming language lecturers were always clear, precise and communicates understandably
x ₂₆	Programming language lecturers made use of enough relevant instructional materials
x ₂₉	Programming language lecturers didn't miss classes
x ₃₁	Programming lecturers were always available
x ₃₃	Programming course lecturers came to class fully prepared
x ₃₆	I fell sick quite often
x ₃₉	It was difficult to charge my computer even within the campus
x ₄₁	Consistent power supply helped me in programming courses
x ₄₂	I had a good background in physics
x ₄₄	I had a good background in English
x ₄₆	Absence of accessible ICT facilities inhibited my programming performance
x ₄₇	The environment where we had programming lectures was not conducive
x ₄₉	The school library was not equipped with materials relevant to programming
x ₅₄	Programming courses were scheduled to non-conducive times
x ₅₅	We had programming classes at unfavorable times
x ₅₆	Programming lecture theatres were equipped with audio-visuals and learning aids
x ₅₇	Programming courses were analyzed clearly to sight
x ₆₁	My family sponsored my academic pursuit
x ₆₂	Quarrel between family members is normal
x ₆₄	Quarrel between my family members escalates a times
x ₆₅	My father is familiar with computers
x ₆₈	My parent would want me to offer programming courses
x ₆₉	I received educational advices from family members often
x ₇₀	My family believed that a proper study will help me in programming courses

B. Comparative Evaluation of the Developed Student Performance Models

The expressions of variable-based LRC model with positive significance agree with some already established variables such as students' lack of understanding, absence from class, negative attitudes towards programming, students' performance in Mathematics [29], study habit [30], review study materials, self-evaluate, rehears explaining materials, and studying in a conducive environment [31], students' class attendance (Pudaruth, Nagowah, Sungkur, Moloo and Chinia [32], Teaching Styles and Strategies [33], availability of University facilities [6] and mathematics background [34]. However, this study established the negative significance of variables such as group discussions, good background in physics and English among others on student performance in programming as against the reports of Mohd and Abdullah

[29] and Darwin *et al.* [30] for example. In general, the variable-based LRC model is an explicit extension of most existing counterparts by salient factors such as Lecturers' Teaching Style (LTS), Health (OH), Electricity (OE), Parental Education (FPE), Student Fear and Perception (SF), Tutorials and Extra Classes (ST) among others which have not been duly considered by other previous works.

V. CONCLUSION AND FUTURE WORKS

This study was conducted to explore the factors affecting the academic performance of undergraduates in programming courses and develop models with which the performance of students can be predicted. The research was conducted on a sample of students who have at one time or the other offered Web programming, C or JAVA within the Federal University, Oye-Ekiti, Ekiti State, Nigeria between 2012 and 2016. This was based on students' performance records which cut across the second and third (200-300) levels of study within the institution. Machine learning approaches were gainfully employed for the analysis of the retrieved data from a defined number of respondents. Results obtained indicate that the attitude of students and lecturers, fearful perception of students, erratic power supply, university facilities, student health, students' attendance are significant to the performance of students in programming courses. It is recommended that future research adopts improved statistical machine learning approaches to comparatively model the learning behaviour in private and public Universities of Nigeria and identify the salient factors significant to performance of students in both systems for robust evaluation of quality of training and to aid effective decision making by the government, students and University education stakeholders. Furthermore, a consideration of all programming courses being offered in the institution and a relatively larger population might graciously improve the findings reported in this study. The existing statistical machine learning approaches can also be extended while some other ones can be introduced for more accurate results.

REFERENCES

- [1] A. Kofi, T. John and P. Prince (2013), Causes of Failure of Students in Computer Programming Courses: The Teacher – Learner Perspective: International Journal of Computer Applications, Vol. 77, No. 12, pp. 27 – 32.
- [2] R. Gomes, A. Anabela and T. Mendes (2007), Learning to program - difficulties and solutions, International Conference on Engineering Education, pp. 118 – 124.
- [3] S. Akinola and K. Kazeem (2014), Factors influencing students' performance in computer programming: A Fuzzy set operations approach, International Journal of Advances in Engineering and Technology, Vol. 7, Issue 4, pp. 1141 – 1149.
- [4] B. Mustafa Baser (2013), Attitude, Gender and Achievement in Computer Programming: Middle-East Journal of Scientific Research, Vol. 14, Issue 2, pp. 248 – 255.
- [5] E. Chermahini (2013). Learning Styles and Academic Performance of Students in English as a Second - Language Class in Iran. Bulgarian Journal of Science and Education Policy, 7, 2, 322 – 333.
- [6] O. Ogbogu (2014), Institutional Factors Affecting the Academic Performance of Public Administration Students in a Nigerian University: Public Administration Research, Vol. 3, No 2, pp. 171 – 177.
- [7] E. Mushtaq and R. Khan (2012). Factors affecting students' academic performance, Global Journal of Management and Business Research, Vol. 12, Issue 9, pp. 189 – 200.

- [8] R. Masura, S. Shahrani, Rodziah, L., Faezah, M.Y., Faridatul, A.Z. & Rohizah A.R. (2012). Major problems in basic programming that influence student performance. *Procedia - Social and Behavioral Sciences*, 59, 287 – 296.
- [9] I. Irfan and N.K. Shabana (2012). Factors affecting students' academic performance. *Global Journal of Management and Business Research*, 12, 9, 189 – 200.
- [10] I. Amiroh and S. Farinda (2016). Factors Affecting the Academic Performance of Executive Diploma Students: The case of University of Malaya Center for Continuing Education.
- [11] T. M. Fagbola, Adeyanju I. A., O. M. Olaniyan and Ayodele Oloyede (2018): Development of Multi-Factored Predictive Models for Students' Performance in Programming Courses in a Nigerian Tertiary Institution, *Transylvanian Review: Vol XXVI, No. 25*, pp. 6809-6820.
- [12] A. Oyerinde and T. Chia (2017), Predicting Students' Academic Performances – A Learning Analytics Approach using Multiple Linear Regression: *International Journal of Computer Applications*, Vol. 157, No 4, pp. 37 – 44.
- [13] O. Gurmeet and W. Williamjit (2016), Prediction of Student Performance Using Weka Tool: *An International Journal of Engineering Sciences*, Vol. 17, pp. 8 – 16.
- [14] A. Ashish, Saeed, Maizatul and Hamidreza (2015), Clustering Algorithms Applied in Educational Data Mining: *International Journal of Information and Electronics Engineering*, Vol.5, No. 2, pp. 112 – 116.
- [15] S. O. Olabiyisi., T. M. Fagbola., E. O. Omidiora Elijah O. and C. A. Oyeleye (2012): "Hybrid Metaheuristic Feature Extraction Technique for Solving Timetabling Problem". *International Journal of Scientific and Engineering Research, USA*, 3(8): pp 1-6.
- [16] T. M. Fagbola, S. O. Olabiyisi, A.A. Adigun (2012): "Hybrid GA-SVM for Efficient Feature Selection in E-Mail Classification". *Journal of Computer Engineering and Intelligent Systems, International Institute of Science, Technology and Education (IISTE), New York, USA*, 3(3): pp 17-28.
- [17] C. A. Oyeleye, Olabiyisi S. Olatunde, Omidiora E. Olusayo and T. M. Fagbola (2014). Hybrid Metaheuristic Simulated Annealing and Genetic Algorithm for Solving Examination Timetabling Problem. *International Journal of Computer Science and Engineering (IJCSE)*, India, 3 (5): pp 7-22.
- [18] B. Kolo, S. Adepoju and A. Alhassan (2015), A Decision Tree Approach for Predicting Students Academic Performance: *International Journal of Education and Management Engineering*, Vol. 5, No. 5, pp. 12 – 19.
- [19] J. R. Quinlan (1992). Learning with Continuous Classes, *World Scientific*, pp. 343-348.
- [20] E. K. Onyari and F. M. Ilunga (2013). Application of MLP Neural Network and M5P Model Tree in Predicting Streamflow: A Case Study of Luvuvhu Catchment, South Africa, *International Journal of Innovation, Management and Technology*, Vol. 4, No. 1, pp. 11-15, DOI: 10.7763/IJIMT.2013.V4.347.
- [21] W. Yong Wang and W. Ian H. Witten (1997). Inducing model trees for continuous classes. In *Proc. of the 9th European Conf. on Machine Learning Poster Papers*, pp. 128-137.
- [22] H. Hieu Chi Huynh (2015). Improving M5 Model Tree by Evolutionary Algorithm, Master's Thesis in Computer Science, Ostfold University College, Norway.
- [23] B. Cestnik and Bratko I. (1991). On Estimating Probabilities in Tree Pruning. pp. 138-150.
- [24] T. Chai and R. R. Draxler (2014). Root mean square error (RMSE) or mean absolute error (MAE)? – Arguments against avoiding RMSE in the literature, *Geosci. Model Dev.*, 7, pp. 1247–1250, doi:10.5194/gmd-7.
- [25] S. C. Gupta (2013). *Fundamentals of Statistics, Seventh Revised & Enlarged Edition*, Himalaya Publishing House.
- [26] L. Jia Li (2002). *Linear Methods for Classification*, retrieved from <http://www.stat.psu.edu/~jiali>
- [27] N. Imran Naseem, T. Roberto Togneri and Mohammed Bennamoun (2010). Linear Regression for Face Recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 32, No. 11, pp. 2106-2112.
- [28] W. Cort J. Willmott and M. Kenji Matsuura (2005). Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance. *Climate Research*, Volume 30, pp. 79-82.
- [29] H. Mohd and A. Abdullah (2013), Predicting Student Performance in - Object Oriented Programming Using Decision Tree : A Case at Kolej Poly-Tech Mara, Kuantan:
- [30] Y. Darwin, Isidro and Beelinda (2015), Study Habits and the Performance of BSCs Students in Computer Programming-1: The IIER International Conference.
- [31] R. Miguel and U. Ksenia (2015), Influence of Study Habits on Academic Performance of International College Students in Shanghai: *Canadian Center of Science and Education*, Vol. 5, No. 4, pp. 42 – 55.
- [32] W. Pudaruth, T. Nagowah, E. Sungkur, Y. Moloo and E. Chiniah (2013), The Effect of Class Attendance on the Performance of Computer Science Students, *2nd International Conference on Machine Learning and Computer Science*.
- [33] I. Agba and E. Michael (2015), Electric Power Supply and Work Performance of Academic Staff in Nigeria Universities: A Synergy Analysis: *Indian Journal of Commerce & Management Studies*, Vol. 6, Issue 1, pp. 33 – 44.
- [34] S. Owusu, A. Arthur and S. Sarpong (2013), Causes of Failure of Students in Computer Programming Courses: The Teacher – Learning Perspective: *International Journal of Computer Application*, Vol. 77, No. 12, pp. 27 – 32.

New Techniques to Enhance Data Deduplication using Content based-TTTD Chunking Algorithm

Hala AbdulSalam Jasim, Assmaa A. Fahad
Department of Computer Science, College of Science
University of Baghdad
Baghdad, Iraq

Abstract—Due to the fast indiscriminate increase of digital data, data reduction has acquired increasing concentration and became a popular approach in large-scale storage systems. One of the most effective approaches for data reduction is Data Deduplication technique in which the redundant data at the file or sub-file level is detected and identifies by using a hash algorithm. Data Deduplication showed that it was much more efficient than the conventional compression technique in large-scale storage systems in terms of space reduction. Two Threshold Two Divisor (TTTD) chunking algorithm is one of the popular chunking algorithm used in deduplication. This algorithm needs time and many system resources to compute its chunk boundary. This paper presents new techniques to enhance TTTD chunking algorithm using a new fingerprint function, a multi-level hashing and matching technique, new indexing technique to store the Metadata. These new techniques consist of four hashing algorithm to solve the collision problem and adding a new chunk condition to the TTTD chunking conditions in order to increase the number of the small chunks which leads to increasing the Deduplication Ratio. This enhancement improves the Deduplication Ratio produced by TTTD algorithm and reduces the system resources needed by this algorithm. The proposed algorithm is tested in terms of Deduplication Ratio, execution time, and Metadata size.

Keywords—Data deduplication; big data compression; data reduction; Two Threshold Two Divisor (TTTD); chunking algorithm

I. INTRODUCTION

There is an explosion on the amount of digital data in the world right now, as manifest by the considerable growth in the measured amount of stored data in 2010 and 2011 from 1.2 zettabytes to 1.8 zettabytes¹, respectively [1], and the prophesied amount of data to be created in 2020 is 44 zettabytes [2], [3]. So manage the storage which is cost-effectively, has become an important task of the most challenging in the big data era. The workload studies performed by an American multinational corporation Dell, EMC, (Richard Egan, Roger Marino & John Curly the E, M & C in EMC) and Microsoft, suggest that approximately 50% and 85% of the data are redundant in their primary and secondary storage systems, respectively.

According to International Data Corporation (IDC) recent study, almost 80% of the surveyed corporations indicated that they are using in their storage systems to reduce redundant data kind of data deduplication technologies, which increased storage in an efficient way and reduced the costs of storage spaces [4].

Data deduplication does not only reduced storage space, but also decreased the transmission rate by eliminating redundant data in low bandwidth network environments. A sub file-level chunking deduplication system breaks the input data stream into multiple data “chunks” that are individually distinguished by a hash signature (e.g., SHA-1), and detects the duplicate ones by some kind of comparison method. Deduplication systems remove duplicate chunks, and store or transfer only one copy of them to achieve the goal of saving storage space or network bandwidth. In the other hand Deduplication system, suffers from the long execution time and the need of many CPU resources on its job.

Teng-Sheng Moh [5] in 2010 adds a new switch condition to enhance the execution time of TTTD algorithm with the same deduplication ratio. He reduced the value of the main divisor (D) and the second divisor (Ddash) to the half when the break point was not found before 1600 byte, this condition reduced about 6% of the running time and 50% of the large-sized chunks.

Manogar and Abirami [6] in 2014 first examined and compared different deduplication techniques, and then they concluded that variable size data deduplication is more efficient than the rest of the deduplication techniques.

AbdulSalam and Fahad [7], in 2017 performed a survey on different chunking algorithms of data deduplication. They discussed, studied the most popular chunking algorithm TTTD, and evaluated this algorithm using three different hashing functions; Rabin Finger print, Adler, and SHA1 implemented each one as a fingerprinting and hashing algorithm and then compared the execution time and deduplication elimination ratio.

In this paper a new chunking condition is added to enhance the deduplication ratio and a new four hashing function is proposed to improve the matching process by reducing the probability of hash collision occurrence. In addition, a searching technique suggested in order to reducing the time needed for deduplication process.

¹ IDC, “The 2011 digital universe study,” Tech. Rep., Jun. 2010, [Online]. Available: <http://www.emc.com/collateral/analystreports/idc-extracting-value-fromchaos-ar.pdf>

TABLE I. CHARACTERISTICS OF THE USED DATA SET

Data Set	On Line Link	Number of Files and Folders	Total Size
Versions of Emacs of GNU	http://www.gnu.org	16,296 Files, 327 Folders	580 MB
Versions of 3DLDF of GNU	https://www.kernel.org	5,795 Files, 63 Folders	1.27 GB

The input data to the system consists of a number of files with diverse sizes and types. The system process these files as one file at a time. The proposed system is developed and tested on two data sets, which belongs to the GNU file system in order to show the efficiency of the proposed system. Table I shows the characteristic of each data set.

II. DATA DEDUPLICATION SYSTEM USING TTTD CHUNKING ALGORITHM

In general, any deduplication system will pass into three stages: (Chunking, Hashing and Indexing, and Matching stage). The theory of each part will explain briefly:

A. Chunking

The first step of data deduplication is chunking, it partitioning the input data stream (file) into small and non-overlapping parts named chunks. The chunking operation is performed using certain type of rolling hash that depends on the contents of the text itself so that for two strings with the same contents it will produce the same hash value. This stage is a very important stage; the deduplication ratio depends on the chunks produced from this stage [7].

TTTD is a variable size-chunking algorithm [8], it use Rabin Fingerprint to find the hash value of substring with predefined window size (48 byte). If the hash of this substring satisfy the condition of TTTD it will considered as a breakpoint otherwise slide the window size one byte [9].

Formula (1) used to compute Rabin Fingerprint for the first substring. Then, Formula (2) used for the rest substring, worked by remove first character, and added the new one [4].

$$\text{Rabin}(B_1, B_2, \dots, B_\alpha) = \left\{ \sum_{i=0}^{\alpha} (B_i * P^{\alpha-i}) \right\} \text{Mod } D \quad (1)$$

$$\left\{ \left[\text{Rabin}(B_1, B_{i+1}, \dots, B_{i+\alpha-1}) - B_i * P^{\alpha-1} \right] * P + B_{i+\alpha} \right\} \text{Mod } D \quad (2)$$

Here: D is the average chunk size [7], Bx is the ASCII code for the substring characters, P is a prime number α is the size of the sliding window.

B. Hashing and Indexing

The main target of hashing and indexing stage is to compute the hash value for whole chunk and adding it to the lookup table or index table. When the finger print satisfy one of TTTD conditions then the whole chunk of text will be sent to the Hash function (like SHA-1 or MD5).The hash value result from the hash function will be used to compare between the chunks. The name of the chunk is the location that saved with inside the chunk container. The content of the lookup table will be a set of records consist of two fields, the first field contain the chunk name, and the second field will be contain its hash value:

D:\DataBase_test\emacs-22.1\AUTHORS\Chunk-0.txt
123456

C. Matching

In original systems, the chunk of new file will compared to the chunk of the files that have the same name and type. If there is a matching then the system will retrieve its lookup table, and compare all the chunks of the new file with the chunks of the old one. For the duplicated chunks, delete the new chunk and perform a logical reference to the old one in a final lookup table of the new file, otherwise save the chunk, and add its name and its hash to the final lookup table. A collision problem may occur during the matching operations, to solve the collision problem a byte-to-byte comparison must be performed [10]. The number of collision reduces the performance of the system due to the time needed to solve it. Reducing the number of collision is one of the aims of a good deduplication system.

III. PROPOSED SYSTEM AND METHOD

In this paper, a Content Based Two Threshold Two Divisor with Multi-Level Hashing Technique (CB-TTTD-Multi-Level Hashing Technique) based on TTTD algorithm suggested to enhance deduplication technique by speed up the deduplication operation and increase its compression ratio.

A. Chunking

CB-TTTD-Multi-Level Hashing Technique introduces a new hash functions to compute the fingerprints for each tested data stream. For each character in the first string of window size (36 byte) the fingerprint value is calculated using (3). Then for each of the following substrings, fingerprint value is calculated using (4).

$$\text{FingerPrint}(B_0, B_2, \dots, B_\alpha) = \left\{ \sum_{i=0}^{\alpha-1} \text{Val}[B_i] * 2^{i+1} \right\} \quad (3)$$

$$\text{New FingerPrint}(B_{i+1}, \dots, B_{i+\alpha+1}) = \left[\left\{ \text{FingerPrint}(B_i, \dots, B_{i+\alpha}) - \text{Val}[B_i] \right\} \div 2 \right] + \text{Val}[B_{i+\alpha+1}] * 2^{\alpha-1} \quad (4)$$

Here: α is Substring Size, B1 ... Bα: are the substring characters, Val [Bi]: is the value of index [Bi] in Fingerprint-array. The value of character taken from an array of 256 position that represent the printable characters filled with random value of (1, 2) to produce an array as in Fig. 1:

Unlike Rabin fingerprint CB-TTTD-Multi-Level Hashing Technique, uses a value retrieved from the fingerprint array instead of using the ASCII code of the character, this step speed up the computation and reduce the overhead of CPU needed for each fingerprint because it uses very small values. The system test different values for the fingerprint array such as:

[0,1], [1,0], [0, 0, 1,1], [1,1,0,0], [1,2], [1,2,3,1,2,3],
[1,2,3,4,.... 255].

Array [0]	Array [1]					Array [255]
1	2	1	2	1	2

Fig. 1. Fingerprint array.

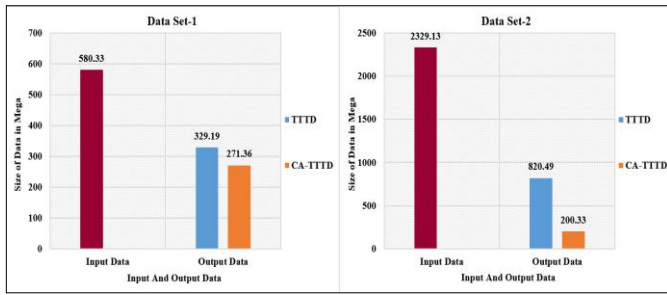


Fig. 2. Input to output data ratio of the two-compared algorithms.

The most efficient values that produces a high deduplication ratio was the sequence of [1, 2] values. This technique helps to produce different hash values for different substrings that helps in detection more redundant data and increase the deduplication ratio.

In addition, CB-TTTD-Multi-Level Hashing Technique gives a weight to the characters in the data stream. The system considers the Dot character (‘.’), as a new condition, in addition to the main and second divisor condition of TTTD. When a ‘.’ character is found followed by a (space) or (end of line) this paragraph is considered as a separated chunk. The advantage of this condition appears in case of two paragraphs considered as one chunk, then any change in one paragraph may affect the next one, but in this case, the effect will be limited with the changed paragraph only. Adding this condition increased the deduplication ratio and with the same chunking time, because it does not need any extra processing steps to compute the chunk boundary. Fig. 2 illustrates the output size using original TTTD chunking algorithm and CB-TTTD-Multi-Level Hashing Technique on the same dataset as input and the differences between them. Table II shows the result produced by implementing CB-TTTD-Multi-Level Hashing Technique on Dataset1 compared with TTTD algorithm.

B. Hashing and Indexing

The old deduplication system is suffering from the wasted time needed to solve the collision problem. CB-TTTD-Multi-Level Hashing Technique suggests a new method that uses four hashing functions rather than one to solve the collision problem. The technique will compute and save four hash values for each chunk, as shown in Table III Using the hash functions shown below:

$$Hash1(chunk) = \left\{ \sum_{i=0}^{s-1} (Array1[Bi] * 2^k) \right\} \quad (5)$$

$$Hash2(chunk) = \left\{ \sum_{i=0}^{s-1} (String[Bi] * A1) \& 0xFFFFFFFF \right\} \quad (6)$$

$$Hash3(chunk) = \left\{ \sum_{i=0}^{s-1} (Array3[Bi] * 2^k) \right\} \quad (7)$$

$$Hash4(chunk) = \left\{ \sum_{i=0}^{s-1} (String[Bi] * A2) \& 0xFFFFFFFF \right\} \quad (8)$$

TABLE II. THE EFFECT OF DOT CONDITION ON EACH SYSTEM

Algorithm	Number of Chunks	Deduplication Ratio	Size of Metadata in MB	Time in second
TTTD without Dot	621861	1.78300	82.1	2304
TTTD with Dot	1542374	1.81176	124	3349
Proposed System without Dot Condition	644933	2.03845	16.7	458
Proposed System with Dot Condition	960091	2.1386	24.1	582

TABLE III. ARRAY1 AND ARRAY2 VALUES

Array	[0]	[1]	[2]	[3]	[4]	...	[255]
Array1	0	1	2	3	4	...	256
Array2	0	1	0	1	0	...	1

Here: S is the chunk size, Array1 and Array2 is an array of 255 value as shown in Table III, K is an integer value equals to (i mod 8), A1 and A2 values is 5, 11 respectively, their values increased by one every iterator, 0xFFFFFFFF used to get limited range of value.

Using these hash functions reduces matching time by solving the collision problem in an efficient way, Also the number of bits needed to store these four hashes are about to 32 bits maximum, which is less than the number of bites needed to save the hash value in SHA-1 and MD5 which yields hexadecimal digits, SHA-1 returning 160 bit. 4 bit per character and thus equals to 40 character, and the output of MD5 hash which is 128 bits equals to 32 characters [11].

The name, the size, and the four hashes values for each chunk must save in Index-Table. The name of the chunks in CB-TTTD-Multi-Level Hashing Technique is an integer number from zero to N, where N is unlimited number increased with each chunk in the system. This information must be ordered in the Index table as shown in Table IV.

CB-TTTD-Multi-Level Hashing Technique also creates a log file for each file in the dataset. This log file will be used in reconstruction operations of the files.

TABLE IV. THE STRUCTURE OF THE INDEX-TABLE

Chunk Name	Chunk Size	Hash1	Hash2	Hash 3	Hash4
0	1268	357315 1	7173770 3	19983	7241537 3
1	466	129976 9	9843858	6961	1009266 6
2	480	133838 6	1050308 0	7454	1076210 0

TABLE V. TIME AND INDEX TABLE IMPACT OF EACH HASH

Number of Hash Used	1-Hash	2-Hashes	3-Hashes	4-Hashes
Chunking Time in Second	353	342	321	336
Deduplication Time in Sec	1050	656	627	562
Size of Log File in MB	6.048	6.048	6.048	6.048
Size of Index Table in MB	10.915	13.802	15.177	18.075
Final Meta Data Size in MB	16.963	19.85	21.225	24.123
Number of Hash Collision	21372	0	0	0

C. Matching

In deduplication matching steps, when a new file comes and passes the two previews stages, the system must detect and eliminate the duplicated chunks. It first check the hash values of the chunks, if the hash values are similar, then the algorithm will compare the two chunks byte to byte, if they are identical the system will delete the new one and add a logical reference to the location of the old one. Otherwise a collision was occurs; the chunks are difference; the system will save the new one as a new chunk. This operation takes a lot of time and overhead the system. Therefore, a new method has to add to deduplication matching process, to enhance the throughput, i.e., saving execution time and reduce CPU resources usage.

In this paper, the Multi-Level Hashing Technique was suggest to enhancing the matching process by using the four hash functions that already computed in hashing stage. If a collision occur in first hash values of the compared chunks then compare the second, third and fourth hash. The test result shows a significantly noteworthy improvement with the time needed by matching process, because comparing four numbers is faster than comparing the whole compared chunks byte-to-byte.

This solution tested with dataset1 and dataset2, the collision founds in first hash function will reduced to zero by the second hash function. For dataset1 with the first hash function the collision number was 21372 for total chunk number 960091, by using the second hash function, the number is reduce to zero. Third and fourth hash functions are uses as an extra step to be assurance the collision is determined, if a chunk overpasses the first two hashes. For each hash in the system the elapsed time and collision number is computed as shown in Table V, the proposed system found that four hashes is a balanced number between time and the size of index table.

When a new file comes, it must be chunked and hashed, and then each chunk will be compared with all chunks within the database. The previews deduplication systems search for the similarity of the chunks within a file name or type in the dataset, the proposed system search for the similarity of the chunk within the whole files in the dataset. The side effect of this method is the time needed to complete the matching operation. To enhance this method, the size of the chunk is utilized as the searching parameter, and a binary search technique is used instead of linear searching method. Because when using liner search, the system will be with $O(N)$ complexity while using the binary search reduce the complexity to $O(\log N)$ [12].

To implement a binary search in an efficient way, the new matching algorithm that proposed in this paper, divide the chunks in the workspace into 16 groups, depending on the chunk size as shown in Table VI. The first group contains the chunks with size (0 – 462) byte, and the second group is for the chunks with size (463 – 471) byte, and so on. The number of the chunks in each group is approximately equal.

Fig. 3 shows chunks distribution, the minimum chunk size of the algorithm is 460 byte and the maximum chunk size is 2800 byte as the TTTD algorithm suggested [9]. However, there are special cases where the chunk size is less than 460 byte. These cases are:

- The size of the file is smaller than the minimum chunk size (460 byte) or less than window size (36 byte).
- The size of the rest of the file from the last breakpoint is less than 460 byte or 36 byte.
- The breakpoint could not found after the last breakpoint to the end of the file.

In this paper the data in the above three cases will be preserved as one chunk.

Working with large amount of the chunk as groups is easier and faster than working with it as a one large search space Table VII shows the effect of the searching techniques with respect to time.

TABLE VI. CHUNK DISTRIBUTION ACCORDING TO PARTITIONING METHOD

Index	From	To
0	0	462
1	463	471
2	472	481
3	482	491
4	491	503
5	504	515
6	516	530
7	531	547
8	548	570
9	571	598
10	599	636
11	637	691
12	692	775
13	776	920
14	921	1291
15	1292	2800

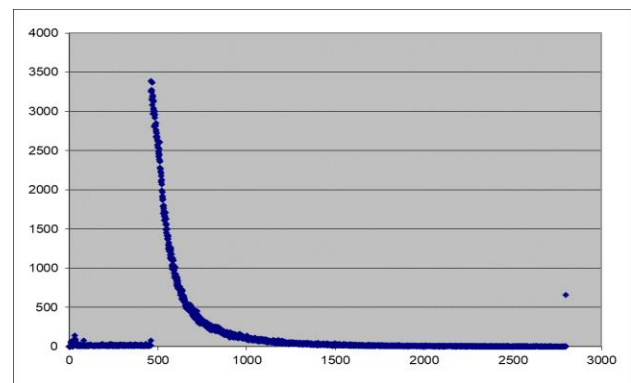


Fig. 3. Distribution of chunk with respect to size.

TABLE VII. COMPARISON BETWEEN SEARCHING WITH ONE BIG SEARCH SPACE AND PARTITIONS METHOD

Benchmarks	One Large Search Space	16 - Parts Search Space
Number of Chunk	960091	960091
Deduplication Ratio	2.1386	2.1386
Chunking Time in Second	320	339
Whole Program Time (Matching and Chunking) in Second	3230	629

To represent the chunk of dataset with a distribution-based representation that summarizes scalar information into much-reduced groups, one of statistical distribution methods should be used [13]. In this paper, CB-TTTD with Multi hashing Technique used the histogram. The disadvantages of used statistical distributions method is that the distribution representing the chunks in one dataset differs from another one.

However, in the proposed case the parts boundaries was approximately equals for all tested datasets. The histogram steps that used to rearrange the chunks in to range from (0 to16) parts depending on chunks size instead of (0-2800) range are:

- Count number of chunk for each size of range (0 – 2800).
- Compute the probability of each size with respect to other; i.e.: $P(i) = (\text{count}(i) / \sum \text{count}(i))$.
- Compute the Probability Ratio for each size using the formal: $PA(i) = (\sum P(i))$.
- Multiply the PA column with Density Slicing number which in our case is (15) and round the result to nearest integer number, that give as range from (0-15) only, see Table VI.

IV. RESULTS AND DISCUSSION

Proposed technique is implemented on a machine with configuration Intel i7 CPU with Installed memory 4.00 GB on 64bit windows OS. To implement proposed technique dataset is collected as mentioned in Table I. TTTD chunking algorithm with Rabin fingerprint and SHA-1 hashing algorithm implemented also in the same environment to compare the result of the proposed system with it.

To analyze CB-TTTD-Multi-Level Hashing Technique the following performance metrics are used. Table VII shows the result of the two algorithms.

- Data Size after Deduplication: It describes how many data remains after the data deduplication eliminates all redundant data.
- Deduplication Gain: It indicates how much unique content is present in the dataset. In this paper, it calculated as in (9).

$$\text{Deduplication Gain} = \frac{\text{The Size of Deduplicated Data Detected}}{\text{Total Output Data Size After Deduplication}} \quad (9)$$

TABLE VIII. COMPARISON TTTD AND CB-TTTD-MULTI-LEVEL HASHING

Evaluation Metrics	DataSet-1		DataSet-2	
	TTTD	CB-TTTD	TTTD	CB-TTTD
Input Data (MB)	580.33	580.33	2329.13	2329.13
Data Size after Deduplication (MB)	329.19	271.36	820.49	200.33
Duplicated Data Detected (MB)	251.14	308.97	1508.64	2128.8
Deduplication Rate	1.7628	2.1386	2.8387	11.6264
Deduplication Gain	0.4327	0.5324	0.6477	0.9139
Chunking Time in Sec	727	462	3747	1075
Deduplication Time in Sec	2632	532	9083	1539
Total Number of chunk	621861	960091	2468026	2611322
Average Chunk size	978.55	633.82	989.56	935.26

- Deduplication Ratio: The data deduplication ratio measures the effectiveness of the deduplication process, it is expressed as in (10).

$$\text{Deduplication Ratio} = \frac{\text{Total Input Data Size Before Deduplication}}{\text{Total Input Data Size After Deduplication}} \quad (10)$$

- Average Chunk size: Calculated as in (11)

$$\text{Average Chunk Size} = \frac{\text{Total Input Data Size}}{\text{Total Number of Chunks}} \quad (11)$$

- Chunking and Hashing time: It is the total time taken to perform hashing and chunking operation.

Experimental results are shown in Table VIII. These results clearly demonstrate that CB-TTTD-Multi-Level Hashing Technique satisfactorily reduces the deduplication processing time and increases its ratio.

Result charts also clearly demonstrate that our proposed approach perform better than original approach for small and big data sets, the deduplication gain is also increased.

The proposed finger print equation used in chunking stage is more efficient than Rabin fingerprint equation; it increase the number of the chunks by the way it works, especially small chunk sizes, with less CPU overhead cost which increase the deduplication ratio. In addition, using Content Based condition (Dot character), also increased the number of the small chunks leading to increasing deduplication ratio without influence chunking time Fig. 4 illustrate the average chunks size of the two algorithms.

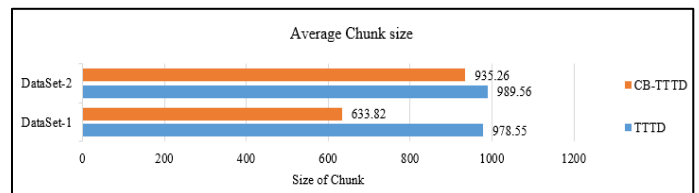


Fig. 4. Average chunk size.

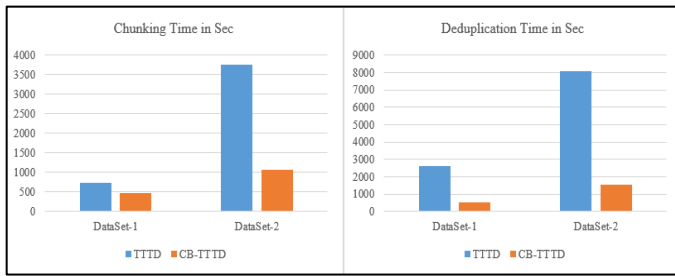


Fig. 5. Deduplication and chunking time.

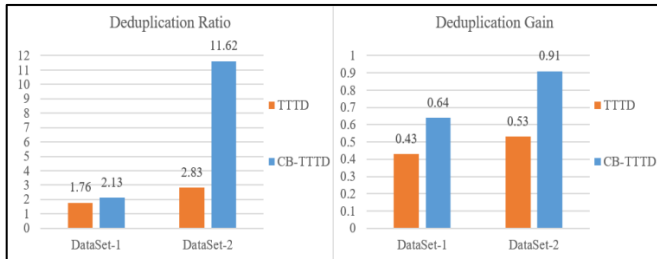


Fig. 6. Deduplication ratio and gain for both system.

The suggested comparing method improved deduplication ratio, but the execution time was increased. Nevertheless, this disadvantage of comparison was solved by partitioning the search space, and the use of mathematical modules to overcome the collision problem, which is an efficient solution that enhanced the execution time of matching process. As shown in Fig. 5, the time of chunking and the overall time (deduplication time) is less than the time of TTTD algorithm.

Fig. 6 depicts data deduplication gain and ratio. In proposed CB-TTTD-Multi-Level Hashing technique, data deduplication ratio and its gain is high as compared to traditional deduplication methods.

V. CONCLUSIONS AND FUTURE WORK

In big data storage, data is too large and efficiently store data is difficult task. To efficiently stores and de-duplicate the data, this paper suggest a new technique to reduce the deduplication ratio. This technique examined the deduplication detection and elimination system performance and explained the rationale parts, data deduplication consist of three stages, the enhancement operation involved all stages that leads to good deduplication ratio and fast execution time, The Metadata produced by CB-TTTD with Multi hashing technique is less than the one that produced by traditional

algorithms. The effectiveness of the proposed method was evaluated using two relative datasets; the preliminary results are encouraging to go forward toward developing new method for detection and elimination deduplication algorithms to meet the challenges and demands of fast and efficient deduplication systems. Moreover, we can use some kind of fast compression with the Meta data (Index Table and Log file) to saving more disk space.

REFERENCES

- [1] D. Stevenson and N. J. Wagoner, "Bargaining in the shadow of big data," Fla. Law Rev., vol. 66, no. 5, p. 66, 2014.
- [2] Gantz, John, and David Reinsel. "The digital universe in 2020: Big data, bigger digital shadows, and biggest growth in the far east." IDC iView: IDC Analyze the future 2007.2012, pp 1-16 , 2012.
- [3] Turner, V., Gantz, J. F., Reinsel, D., & Minton, S. "The digital universe of opportunities: Rich data and the increasing value of the internet of things". IDC Analyze the Future,p. 5, 2014.
- [4] Xia, Wen and Jiang, Hong and Feng, Dan and Dougli, Fred and Shilane, Philip and Hua, Yu and Fu, Min and Zhang, Yucheng and Zhou, Yukun, "A comprehensive study of the past, present, and future of data deduplication," Proc. IEEE, vol. 104, no. 9, pp. 1681–1710, 2016.
- [5] T.-S. Moh and B. Chang, "A running time improvement for the two thresholds two divisors algorithm," Proceedings of the 48th Annual Southeast Regional Conference on - ACM SE '10. p. 1, 2010.
- [6] E. Manogar and S. Abirami, "A study on data deduplication techniques for optimized storage," 6th Int. Conf. Adv. Comput. ICoAC 2014, pp. 161–166, 2015.
- [7] H. Abdulsalam and A. Fahad, "Evaluation of Two Thresholds Two Divisor chunking algorithm using Rabin fingerprint, Adler, and SHA-1 hashing algorithms," The Iraqi Journal of Science, paper 4C.58, 2017.
- [8] Nisha, T. R., S. Abirami, and E. Manohar. "Experimental study on chunking algorithms of data deduplication system on large scale data." In Proceedings of the International Conference on Soft Computing Systems, pp. 91-98. Springer India, 2016.
- [9] K. Eshghi and H. K. Tang, "A framework for analyzing and improving content-based chunking algorithms," Hewlett-Packard Labs Tech. Rep. TR, 2005.
- [10] Chen, Zhengguo and Chen, Zhiguang and Xiao, Nong and Liu, Fang, "Nf-dedupe: a novel no-fingerprint deduplication scheme for flash-based ssds," in 2015 IEEE Symposium on Computers and Communication (ISCC), pp. 588–594.
- [11] Zhang, Yang and Wu, Yongwei and Yang, Guangwen, "Droplet: A distributed solution of data deduplication," in Proc - IEEE/ACM Int Work Grid Comput. 2012;114–21.
- [12] D. S. KUSHWAHA and A. K. MISRA, "Data structures a programming approach with C," 2nd ed. PHI Learning Pvt. Ltd., 2014.
- [13] Wang, Ko-Chih and Lu, Kewei and Wei, Tzu-Hsuan and Shareef, Naeem and Shen, Han-Wei, "Statistical visualization and analysis of large data using a value-based spatial distribution," in 2017 IEEE Pacific Visualization Symposium (PacificVis), 2017, pp. 161–170.

Security Improvement in Elliptic Curve Cryptography

Kawther Esaa Abdullah, Nada Hussein M. Ali
Department of Computer Science, College of Science
University of Baghdad
Baghdad, Iraq

Abstract—This paper proposed different approaches to enhance the performance of the Elliptic Curve Cryptography (ECC) algorithm. ECC is vulnerable to attacks by exploiting the public parameters of ECC to solve Discrete Logarithm Problem (DLP). Therefore, these public parameters should be selected safely to obviate all recognized attacks. This paper presents a new generator function to produce the domain parameters for creating the elliptic curve; a secure mechanism is used in the proposed function to avoid all possible known attacks that attempts to solve the Elliptic Curve Discrete Logarithm Problem (ECDLP). Moreover, an efficient algorithm has been proposed for choosing two base points from the curve in order to generate two subgroups in a secure manner. The purpose of the aforementioned algorithm is to offer more confidence for the user since it is not built upon a hidden impairment that it could be subsequently utilized to retrieve user's private key. The Elliptic Curve Diffie Hellman (ECDH) algorithm is implemented to exchange a session key between the communicating parties in a secure manner. Beside, a preprocessing operation is performed on the message to enhance the diffusion property and consequently leads to increase the strength against cryptanalysis attack. Finally, the dual encryption/decryption algorithm is implemented using different session keys in each stage of the encryption to boost immunity against any attack on the digital audio transmission. The gained results show the positive effect of the dual elliptic curve system in terms of speed and confidentiality without needing any extra time for encryption.

Keywords—Elliptic curve cryptography; elliptic curve discrete logarithm problem; dual encryption/decryption; Elliptic Curve Diffie Hellman

I. INTRODUCTION

Elliptic curves were suggested by Neal Koblitz and Victor Miller independently in 1985 to design a public-key cryptographic system [1]. The Elliptic Curve Cryptography (ECC) is a public-key cryptosystem which playing an important role in cryptography world. The shorter key size in ECC provides an equivalent protection level for public-key algorithms which utilized the largest key size (e.g., Rivest Shamir Adleman (RSA)). In addition, the ECC offers more security compared to the RSA algorithm since it is based on DLP, while the latest algorithm based on the prime number factorization problem [2], [3]. ECC is based on an Abelian group, the main operation used in ECC is the addition operation; the multiplication is defined as a repeated addition. For example, $\times k = (a + a + \dots + a)$, addition a with times k and it is performed over an elliptic curve. Cryptanalysis includes determining k given a and $(a \times k)$, this is called DLP. The definition of elliptic curve is based on the equation, two variables and two coefficients; the values of variables and

coefficients are limited to elements of a finite field [1]. In this paper, the elliptic curve over the prime field GF_p is considered.

A. Mathematics of ECC Over Finite Field

In general, an elliptic curve E over prime field (F_p) denoted by $E(F_p)$ is given by simplified the weierstrass equation as below [1]:

$$E: y^2 \bmod p \equiv (x^3 + ax + b) \bmod p \quad (1)$$

Where, b and $x, y \in F_p$, The value of variables are sets of elements from 0 to $p - 1$. In addition, the coefficients must satisfy (2), where Δ denoted to the discriminant of E .

$$\Delta = (4a^3 + 27b^2) \bmod p \neq 0 \quad (2)$$

B. Point Addition

If two points on an elliptic curve were added to each other, the output result represents a third point that denotes the intersection of that curve. Graphically, drawing a straight line between any two points on a curve represents a tangent line and reflects a third point around the x - axis as denoted in (5). The formula $P + Q = -R$ represents the addition operation between points $P(x,y)$ with $Q(x,y)$ to produce $R(x,y)$ by applying (3) and (4) [1], [2].

C. Point Doubling

The output value of adding a point $P(x,y)$ on the curve to itself in condition that $y_p \neq 0$ will yield the point R . One could draw a tangent line where the intersection of that line on the curve represents the cross reflection point on x - axis (the R point), where $+P = 2P = -R$. Equations (3), (4) and (6) are used to compute second point $R(x,y)$ and the tangent line (slope), respectively [1],[2].

$$x_R = (\lambda^2 - x_P - x_Q) \bmod p \quad (3)$$

$$y_R = (\lambda(x_P - x_R) - y_Q) \bmod p \quad (4)$$

Where,

$$\lambda = (y_2 - y_1 / x_2 - x_1) \bmod p \quad (5)$$

$$\lambda = (3x_1^2 + a / 2y_1) \bmod p \quad (6)$$

The strength of the ECC depends on the DLP. This means that logarithm Q to base $(l = \text{Log}_G Q)$, where l is a private key, G is a base point and Q is a public key (G, Q are publicly parameters). There are some attacks on the Elliptic Curve Discrete Logarithm Problem (ECDLP) from given G, Q try to extract l . To avoid all recognized attacks on the ECDLP should be selected, the domain parameters for ECC cautiously and in a secure way.

The layout of this paper is composed of the following sections: the related work of elliptic curve cryptography is introduced in Section II, Section III describes the concepts of the advanced encryption standard and cryptographic hash function, the linear congruential generator is presented in Section IV, Section V clarifies the password based key derivation function, Section VI discusses the proposed system followed by the experimental results and discussion in Section VII, and finally, Section VIII presents the conclusions.

II. RELATED WORK

In literature, many researchers have attempted to utilize the strength of the elliptic curve to implement in different tasks of the public key cryptography. Summarized below are some of the features of the linked work.

Rahul Singh, et al. [4] in 2014 investigated an implementation for ECC encryption and decryption audio file was presented.

Artan Luma, et al. [5] in 2015 presented the encryption and decryption for audio file transported through the network - based on ECC. In this study, the scholars have been concluded that ECC is suitable for large amounts of data and also the ECC is preferred comparing with RSA since it provided the same security with small key size.

Manish Kumar, et al. [6] in 2016 proposed a new method for image security by using DNA for encoding RGB image thereafter applied encryption based on Elliptic Curve Diffie-Hellman Encryption (ECDHE). This algorithm supplied a double layer of security.

Fang, Xianjin and Wu, Yanting. [7] in 2017 studied the details of the elliptic curve cryptography, this discussion includes the basic information about ECC and how to partition a message into blocks and encoding/decoding the message into points on the curve using the koblitz method. Also, is presented the encryption/decryption with the elliptic curve. The researcher concluded that the ECC is utilized for encryption, key exchange, and the digital signature with swift and lesser memory.

Kawther, Esaa and Nada, Hussein [8] in 2018 have been investigated a new mapping method based on x -coordinate values of an elliptic curve to generate a secret lookup table, this table is used to convert samples of an audio file (or even any data type) into points on the elliptic curve and vice versa. Besides, the changing form of samples before applying the proposed method to make cryptanalysis more difficult to guess the points on the curve by an intruder (through exploiting statistical analysis) to achieve diffusion, the obtained results indicate that the proposed method is faster, more secure and less time-consuming when embedding a message into a point on the curve.

III. ADVANCED ENCRYPTION STANDARD AND CRYPTOGRAPHIC HASH FUNCTIONS

In 2001 the National Institute of Standards and Technology (NIST) has issued the Advanced Encryption Standard (AES). The AES belongs to a symmetric encryption algorithms family and the type of process is a block cipher, it replaces the Data

Encryption Standard (DES) as a criterion for an enormous domain of applications. The block size of a plain-text input to the AES is 128 bits and key size is 128,192,256 bits. The number of rounds in AES algorithm is altered rely on the key size, 10 rounds for key size 16 bytes, 12 rounds for 24 bytes and 14 rounds for 32 bytes. Each round includes four transformation functions are (SubBytes, ShiftRows, MixColumns, and AddRoundKey) but the last round comprises the three transformations except MixColumns [1], [9], [10].

Cryptographic hash functions play a significant part in the computer security and can be used in various applications such as in Message Authentication, Digital Signatures, hash-based key derivation functions and also in Pseudorandom Number Generator. Hash functions take an arbitrary size of input and produce a fixed size called (hash code or message digest). The basic features of Hash functions are preimage resistant (infeasible to obtain a message from knowing its hash code), second preimage resistant (from given a message x impossible to find a message y has the same hash code of x) and collision resistant (impossible to get any pair (x, y) has the same hash code) [11], [12]. In 2002, NIST designed a new version called the family of the Secure Hash Algorithm (SHA-2) content on three algorithms are SHA-256, SHA-384 and SHA-512 the length of bits produced for each one of them are 256, 384 and 512 bits respectively [1], [12].

IV. LINEAR CONGRUENTIAL GENERATORS

A linear congruential algorithm was suggested by Lehmer and it is mostly used for pseudorandom number generator in order to generate a sequence of numbers such as $x_j = x_0, x_1, x_2, \dots, x_n$ by relying on the following repetition equation [1]:

$$x_{j+1} = ax_j + b \text{ mod } n \quad (7)$$

Where, x_j is the seed value $0 \leq x_j < n$, a represents the multiplier $0 < a < n$, b is the increment value $0 \leq b < n$, n represents the modular $n > 0$. Choosing the values of $(a, b$ and $n)$ accurately helps to produce ideal sequence numbers. The characteristics of this generator are easy to execute and pass the next statistical tests: the frequency test, run test, serial test, poker test and etc., it can be regarded as the best selection for generating robust random numbers [1], [13]. In this study, the proposed system specifies the value of a and x are prime numbers to give a large period in the outcome of the Linear Congruential Generator (LCG).

V. PASSWORD BASED KEY DERIVATION FUNCTION

With the issue of PKCS #5 v2.0 and RFC 2898 [14], is one of the most famous key-derivation functions, Password-Based Key Derivation Function #2 indicated as PBKDF2, with a view to deriving cryptographic keys from (human entrance) passwords. PBKDF2 aims to frustrate expected attacks on password like the dictionary and brute force attacks by incrementing the time necessary for checking every password through two significant parameters in PBKDF2 are the salt and iteration count, it will be difficult to execute these attacks [15]. PBKDF2 is a pseudo-random number PRF that based on some parameters: Salt S , Iteration Count C , Password P and Len_{key} which is the length of derived key also called master key M_{key} actually $(2^{32}-1) \times$ Length of hash function (Len_{hash}). Usually, the PRF includes an HMAC (Hash Message Authentication

Code) structure depends on a cryptographic hash function, which can select from the designer. The general expression of PBKDF2 is defined as [16]:

$$M_{key} = PBKDF2_{(PRF,C)}(P, S, Len_{key})$$

The iteration count C is a constant value represents the number of iterated requisitions of the PRF in order to produce one block of the M_{key} , according to RFC 2898, a minimum C of 1000 is recommended, where the benefit of C increment the calculated cost of carrying out a dictionary attack on a password. While the value of salt offers a wide range of keys for all password. Further, information can be seen in [14].

VI. PROPOSED SYSTEM

The proposed system is divided into four main phases, the main idea of each part of this system is to strength the security against the cryptanalyst and to reduce the risks of compromised the private key in ECC algorithm. The following sections will demonstrate the procedures for each one.

A. The New H-AES-LCG Generator

The first phase of the proposed system is assigned for generating the domain parameters (Prime number (p), a , b) that used for creating the ECC. The purpose of the H-AES-LCG generator is to form an elliptic curve in a random and secure manner to avoid all possible known attacks against it. The steps of constructing the H-AES-LCG generator are described as follows:

Stage 1: Apply the family of SHA-2 includes (SHA-512, SHA-384, SHA-256) and MD5 (Message Digest 5) sequentially on a seed value (*String Password*) shown as below:

1) *The inputs to the family of SHA-2 are:*

- **String Password:** The Left Circular Shift (*LCS*) process has been done with a password before it is entering the hash algorithms. The amount of rotation is determined randomly and denoted by (R_{len}), each letter in the password is shifted with a different R_{len} through converting characters of the password into bytes to store it in a byte array denoted by ($Pass$), the length of a password denotes as P_{len} . The Pseudo-code for applying the *LCS* is shown as below:

- For** $i = 0$ **to** P_{len}
- $R_{len} = (R_{len} + i) \bmod 8$
- $Pass_{rotation}(i) = ((Pass(i) \ll R_{len}) | (Pass(i) \gg (8 - R_{len}))) \& 127$
- End**

The main purpose of the above operation is to increase the strength of the password by permutation its

- **Salt Value:** The time of recording event on the computer (Time Stamp) has been taken as a salt value, it includes the date and time, for example, (2018-02-08 02:11:55 AM). This operation is considered as a one-time pad key to increase the security.

2) Execute the family of SHA-2 (SHA-512, SHA-384, SHA-256) with the salt value on the String Password after applying the rotation process as explained in point 1.

3) Merge between the results of the family SHA-2 to obtain a string of size **1152-bit** denoted as a **Hpassword**. The encoding system implemented in the present study is Unicode which uses two bytes (16-bit) for each character, in this case the **Hpassword** composed of **72** characters.

4) Expanding the **Hpassword** to make it as an input to the MD5 algorithm through converting the **Hpassword** to characters, pick some characters from it and insert some ASCII letters chooses randomly in a specific manner to form a new string. This technique can be described as follows:

At the beginning, determine some secret parameters which are (*start, amount of characters and jump*) to extract a new string from the existing string that has been produced in point 3, the pseudo-code is illustrated as below:

- Initializing the Parameters:

- Ch:** array of character to store characters in the **Hpassword**.
- St:** represents the start value, which is a positive integer value chosen randomly and specified from a range between (0 - size of *Ch*).
- No.ofCh:** amount of characters (length of a new string) is a positive integer chosen according to the length of string that needs.
- jump₁:** a positive integer value to represent the amount of jump between characters.

- Processing:

- $ExHpassword = Ch(St) || ASCII\ Character$
- For** $i = 2$ **to** *No.ofCh*
- $St = (St + jump_1) \bmod length\ of\ Ch$
- $ExHpassword = ExHpassword || Ch(St) || ASCII\ Character$
- End**

The resultant *ExHpassword* string considered as an input to MD5 algorithm. The aforementioned mechanism is applied in order to make an attack for the MD5 algorithm is difficult to guess the string password.

Stage 2: The requirements of AES algorithm for generating random sequence bits are:

1) **Plain-text:** The output of the MD5 algorithm makes as a plain-text to the AES, the number of bits produced from the MD5 corresponds to the size of the plain-text for AES algorithm which is 128 bits.

2) **Key:** Prepare a key for the AES algorithm using PBKDF2 function, it applied to derive a master key (M_{key}) for AES. The parameters to perform PBKDF2 comprised from:

- **Password:** Use the LCG algorithm to generate the pseudo-random numbers and make it as a password denoted by ($Random_{password}$). The values of

$Random_{password}$ are specified in the range between $[0, \dots, 1024]$. The pseudo-code is shown as the following:

- a. Initialization four parameters to use in (7), two prime numbers p_1 is a multiplier, q_1 is seed value, (e.g., the value of p_1, q_1 are 10247 and 8161 respectively) and two positive integer b_1, n_1 where, b_1 is representing the increment equal to 1 and n_1 is a modulus which is equal to 1024 in the present work.
- b. **For** $i = 0$ **to** n_1
- c. $q_1 = (p_1 \times q_1) + b_1 \bmod n_1$
- d. $Random_{password}(i) = q_1$
- e. **End**

- **Salt value:** A time stamp (Date and Time) is regarded as a salt value (S).
- **Iteration Count:** Number of iterations to derive master key denoted as C , in this work C equal to 10000 iterations to increment the calculated cost of carrying out a dictionary attack on a password.
- **Length:** Indicate to the length of derived key (M_{key}) in bits and denoted by (Len_{key}), in this work Len_{key} equal to 256 bits.

The computation of both the plaintext from MD5 algorithm and the key from PBKDF2 are used as input for AES algorithm. The later algorithm is implemented to generate a single encrypted block stored in an array called $Single\ encrypted_{Block}$ of size equal to 128 bits.

Stage 3: In order to increase the randomness of crop from the AES algorithm the Exclusive-OR (XOR) bit wise operation is implemented between, the $Single\ encrypted_{Block}$ and the random numbers generator that generated from the LCG algorithm ($RandomSequece$) by applying (7). The values of the random number are in the range $[0, \dots, length\ of\ Single\ encrypted_{Block}]$, the pseudo-code of the above process is illustrated as below:

- a. Initialize four parameters include two prime numbers are p_2, q_2 (e.g., the value of p_2, q_2 are 7933 and 4093 respectively) and two positive integer numbers, $b_2 = 1, n_2 = length\ of\ Single\ encrypted_{Block}$ to satisfy Equation (7).
- b. **For** $i = 0$ **to** n_2
- c. $q_2 = ((p_2 \times q_2) + b_2) \bmod n_2$
- d. $Random_{Sequece}(i) = q_2$
- e. $RandomSequeceBits(i) = Single\ encrypted_{Block}(i) \oplus Random_{Sequece}(i)$
- f. **End**

The proposed H-AES-LCG algorithm usually generates 128-bit and can be used by any algorithm needs a random number generator. The parameters needed to extract any random number from the string produced by proposing algorithm are:

- **RandomSequeceBits** : Represents the random sequence bits that produced from the H-AES-LCG generator.
- **start**: A positive integer number specified in the range $[0 - 127]$ in order to represent an index of a bit that stored in $RandomSequeceBits$.
- **size**: Refers to how many number of bits that needs.
- **jump₂**: A positive integer number represents jump between the random sequence bits.

This paper has been applied the proposed generator (H-AES-LCG) to extract the domain parameters (p, a, b) for the ECC algorithm, these parameters are used to plot a secure elliptic curve from the output of H-AES-LCG generator. Algorithm 1 shows how to extract the domain parameters (p, a, b) for the ECC algorithm. Also, the aforementioned procedures of H-AES-LCG generator are demonstrated in Fig. 1.

Algorithm 1: Extract the Domain Parameters to Plot a Secure Curve

Input:

- **Password, R_{len} , Salt value**: Serves as input to the family of SHA-2
- **ASCII letters, St, No. of Ch, $jump_1$** : Parameters for producing $ExHpassword$
- **p_1, q_1, n_1, b_1** : Parameters for LCG algorithm to produce $Random_{password}$
- **$Random_{password}, S, C, Len_{key}$** : Parameters for PBKDF2 to produce M_{key}
- **p_2, q_2, n_2, b_2** : Parameters for LCG algorithm to produce $Random_{Sequece}$
- **start, size, $jump_2$** : Parameters for extract a value from $RandomSequeceBits$

Output:

- **p, a, b**

1: Parameters Setting

Generate $RandomSequeceBits$ using the H-AES-LCG generator
 $ML = 1$ // is a positive integer value initialization by one to make a number of multiples two.

$Count = 1$ // To avoid exceeding the number of bits that specified.

2: Extract the domain parameters from invocation the Function1

$p = \text{Function1}(RandomSequeceBits, start, size, jump_2)$
 $a = \text{Function1}(RandomSequeceBits, start, size, jump_2)$
 $b = \text{Function1}(RandomSequeceBits, start, size, jump_2)$

- 3: Check the primality of p using Miller-Rabin algorithm and values of the coefficient (a, b) are satisfied Equation (2), If p is not prime increment p by two and go to Miller-Rabin algorithm. Check values of the coefficient (a, b) are not satisfied Equation (2) go to Step 2.

Function1 (*RandomSequeceBits, start, size, jump₂*)

```

number = RandomSequeceBits(start)
For i = 2 to size

start =
(start + jump2) mod length of RandomSequeceBits
number = number + RandomSequeceBits (start) × ML
ML = ML × 2
Count = Count + 1
if (Count > size)
ML = 1

```

End

Return number

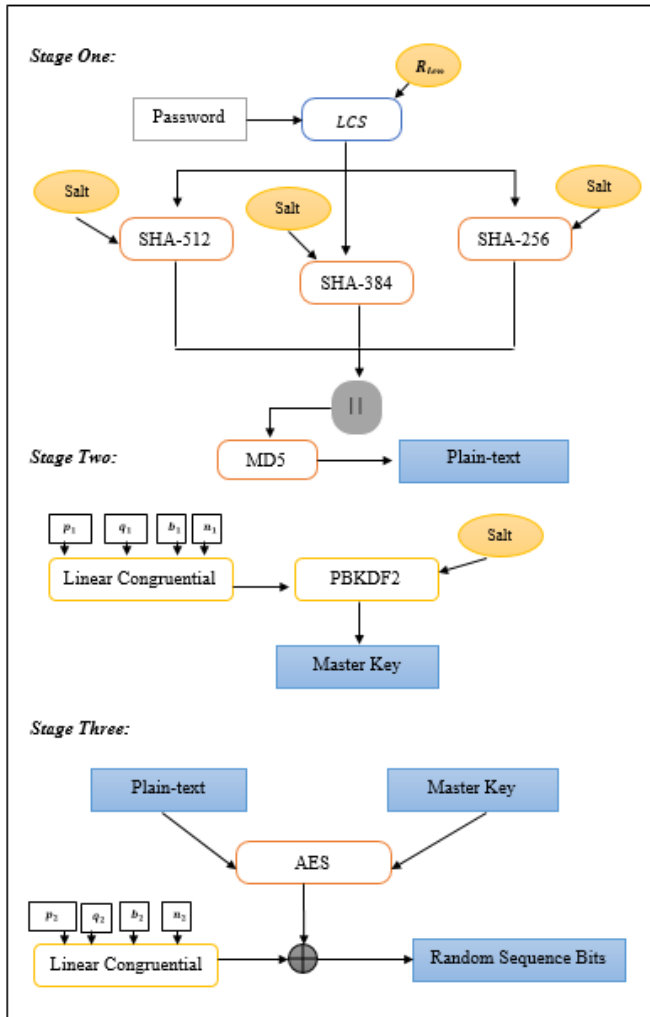


Fig. 1. The block diagram for the H-AES-LCG generator.

B. Base Point Selection from an Elliptic Curve Over Prime Field

After determining the domain parameters (p, a, b) for an elliptic curve denoted by E in a secure mechanism that based on the H-AES-LCG generator, it will be selected two base points from the points on the E . There are two important matters to improve the adequacy of ECC: the determination of the base point from E and the point multiplication operation. The choice of two base points from E that depends on an efficient technique in order to generate two subgroups to implement dual encryption. Each point on elliptic curve should be satisfied as (1) by computing the quadratic residues $\text{mod } p$ denoted as Q_p in order to obtain the value of y - coordinate. The quadratic residues obtained by substitution the value of x in the range $0 \leq x < p$ in (1) and comparing the crop from (1) with Q_p if equal then get the value of y - coordinate.

Every parity in the proposed technique must agree on the value of x - coordinate, which determine randomly in the range $[0, \dots, p - 1]$ to select two base points. Then, applying the right- hand side of (1) to check the quadratic residue of the crop which is denoted as (α) computed by Euler criteria [17], it is defined as below.

$$\alpha^{(p-1)/2} = \begin{cases} 1 & \text{is quadratic residue} \\ -1 & \text{is non - quadratic residue} \end{cases} \quad (8)$$

Where:

α is an integer belong to prime field p .

The checking quadratic residue of x - coordinate to obtain the value of y - coordinate will leads for retrieving the first base point (α, y) denoted by G_1 and second base point (x, y) labeled as G_2 , both achieved through use (1). Algorithm 2 illustrates the proposed technique to determinate two points on an elliptic curve (E).

Algorithm 2: Selection Two Base Points from an Elliptic Curve

Input:

- p : Prime number
- a, b : The coefficients values using in (2)
- x : An integer value specified in the range $[0, \dots, p - 1]$

Output:

- $(\alpha, y), (x, y)$: Two Base Points

1: Parameters Setting

$s = 0$
 $y = 1$

- 2: Compute the quadratic residue of x denoted as Q_x by implementing the right-hand side of Equation (1). After that, checking the value of x has quadratic residue using Equation (8).

$$\alpha = (x^3 + ax + b) \text{ mod } p$$

$$Q_x = (\alpha^{(p-1)/2}) \text{ mod } p$$

While $(\alpha == 0 \mid Q_p != 1)$

$$x = (x + 1) \text{ mod } p$$

go to step 2

End

3: Recover two base points

For $i = 1$ **to** p

$s = s + i \bmod p$

if $(s == \alpha)$

Return

$(\alpha, y), (x, y)$

else

$y = y + 1$

$i = i + 2$

End

Furthermore, the proposed technique provides an enhancement in checking the quadratic residue of $\bmod p$ for a value of x and recovers the value of y – coordinate at the same time. The total time complexity of recovering two base points in the proposed method is $O(\frac{p-1}{2})$.

Moreover, it would be generated two subgroups by applying the *Doubling-Addition* operations; are the basic operations in ECC; from the two base points also called generator points, all one separately. At the beginning, it executes the *Doubling* operation on the base point with itself and then, the *Addition* operation is implemented on the first base point and the outcome from doubling operation and so on until reach to the infinity point (∞), where the doubling operation implements only once in the beginning. In addition, it would be applied the same mentioned above in the second base point in order to generate two subgroups, where the first base point produces the first subgroup and the second base point generates the second subgroup. The points in the subgroup are considered as the public keys and the number of points in the subgroup (the order of subgroup) specify the range of the private keys, where each a private key scalar corresponding to a public key point in the subgroup.

C. Elliptic Curve Diffie Hellman Key Exchange

Elliptic Curve Diffie Hellman is used to exchange a session key between the parties. It has been described as a method to generate two session keys by applying ECDH algorithm in order to use them with dual ECC encryption. This procedure is illustrated below:

Step1: Both parties agree on the domain parameters $(p, a, b, G_1, order_1, G_2, order_2)$ in a secure manner depended on the H-AES-LCG generator and the method for selection the two base points. Where, $order_1$ is the number of points that is generated from the G_1 and $order_2$ represents number of points that is generated from G_2 .

Step2: Each party (Alice and Bob) selects two private keys $(N_{a1}, N_{a2}, N_{b1}, N_{b2})$ respectively, which are smaller than order. The selection of the first private keys (N_{a1}, N_{b1}) for each party is from the range of $order_1$ (the first subgroup), while the

second private keys (N_{a2}, N_{b2}) from the range of $order_2$ (the second subgroup).

Step3: Calculate two public keys for each party through multiply private key with the base point, where $(P_{a1}, P_{a2}, P_{b1}, P_{b2})$ are the public keys for Alice and Bob respectively. The (P_{a1}, P_{b1}) are points in the first subgroup and (P_{a2}, P_{b2}) are points in the second subgroup, and are computed as follows.

$$P_{a1} = N_{a1} \times G_1 \bmod p \quad (9)$$

$$P_{a2} = N_{a2} \times G_2 \bmod p \quad (10)$$

In the same way, Bob can compute his public keys as Alice.

Step4: The computation of two session keys (k_{s1}, k_{s2}) , at the beginning, exchange the public keys between the parties (Alice/Bob) and then, each one computes the session key by multiplying his/her private key with the public key for the corresponding party. This operation is demonstrated as below:

Alice:

$$k_{s1} = N_{a1} \times P_{b1} \bmod p \quad (11)$$

$$k_{s2} = N_{a2} \times P_{b2} \bmod p \quad (12)$$

Bob:

$$k_{s1} = N_{b1} \times P_{a1} \bmod p \quad (13)$$

$$k_{s2} = N_{b2} \times P_{a2} \bmod p \quad (14)$$

Analysis of the above is similar on both sides.

After the two parties agreed of the two session keys (k_{s1}, k_{s2}) , which utilize in dual encryption phase that will be discussed later. Here multiplication is not implied simple multiplication, which is an algebra, rather it is repeated addition of points by the point multiplication operation (scalar multiplication) in ECC.

$k_{s1} = (N_{a1} \times P_{b1}) \bmod p$	$k_{s1} = (N_{b1} \times P_{a1}) \bmod p$
$(N_{a1} \times N_{b1} \times G_1) \bmod p$	$\equiv (N_{b1} \times N_{a1} \times G_1) \bmod p$
$k_{s2} = (N_{a2} \times P_{b2}) \bmod p$	$k_{s2} = (N_{b2} \times P_{a2}) \bmod p$
$(N_{a2} \times N_{b2} \times G_2) \bmod p$	$\equiv (N_{b2} \times N_{a2} \times G_2) \bmod p$

D. The Pre-processing Operations

In cryptography, there are two important operations confusion and diffusion to make the cipher more secure against attacks. According to *Shannon*, confusion means that every bit on the cipher-text should depend on many parts of the key, concealing the relation between the two and make it as complex as possible. Diffusion means that if it changes only one bit (digit) of the plain-text, statistically, half of the cipher-text should change, and vice versa, therefore, making the cryptanalysis so difficult. This complexity generally implemented through of substitutions and permutations. In the present work, the *XOR* and *Circular Shift* bitwise operations are applied to achieve diffusion. The diffusion spreads any change in only one bit of the data to the entire cipher-text, so the sensitivity increased. Moreover, these operations are efficient to perform, less time consuming and provide more security against statistical analysis. These advantages are obtained by removing the characteristics that exploit by an

intruder such as repeated plain text values. The following steps demonstrate the above process:

- The Right Circular Shift (RCS) process is implemented on the samples of an audio file (WavFile). The amount of shifting is determined randomly and denoted by $Rlen$. Each value in the message is shifted with different $Rlen$. The Pseudo-code of RCS is shown as below:
 - For** $i = 0$ **to** $length\ of\ WavFile$
 - $Rlen = (Rlen + i) \bmod 8$
 - If** ($Rlen == 0$)
 - $Rlen = 3$
 - $RCS(i) = (WavFile(i) \ll Rlen) \& 255 | WavFile(i) \gg 8$
 - End**
- The XOR bit-wise operation is implemented between the first sample in the (RCS) and an initial random value (IV), the output is fed back to XOR operation with the next sample in the (RCS) and so on to produce a chaining cipher. Therefore, if one bit of the plain-audio or the initial value altered, all the cipher-text will be changed. This operation is defined as follows:

$$New_{plain-audio}(0) = RCS(0) \oplus IV \quad (15)$$

$$New_{plain-audio}(i) = RCS(i) \oplus New_{plain-audio}(i - 1) \quad (16)$$

Where:

$i : 1 \geq i \leq length\ of\ WavFile.$

IV : an initial random value.

\oplus : Exclusive-OR operation.

The post-processing, the XOR, and the Left Circular Shift (LCS) operations are implemented on the message, where the amount of shifting and the initial value for XOR are determined randomly.

E. Dual Encryption/Decryption in ECC

The purpose of the dual encryption/decryption process for the audio file is to increase the immunity against any expected attacks. This process is implemented through the encryption of the audio file in two layers. Each layer uses different key pairs (Private and Public). Also, the decryption process is executed in the same manner but in reverse order. The procedure for this process is explained in what follows:

Step 1: Dual Encryption

Alice wants to send an encrypted message to Bob, in the present work the message is an audio file. At the beginning, she converts the audio file data into points on the curve denoted by $P_m(x, y)$; then encrypt them to produce the cipher text points denoted by $C_m(x, y)$. The dual encryption process composed of two phases: the first encryption applies the addition operation between $P_m(x, y)$ and $k_{s1}(x, y)$, the second encryption implements the addition operation between $C_m(x, y)$ and $k_{s2}(x, y)$. Equations (17) and (18), respectively define the aforementioned procedure as below:

$$C_{m1} = P_m + k_{s1} \bmod p \quad (17)$$

$$C_{m2} = C_{m1} + k_{s2} \bmod p \quad (18)$$

Step 2: Dual Decryption

Bob received the dual encrypted points and needs to decrypt it, Bob implements the decryption operation on the received points by using the subtraction operation, obtain a reflect coordinate of the subtracted point along an x-axis and execute point addition (i.e., $P(x_1, y_1) - Q(x_2, y_2) = P(x_1, y_1) + Q(x_2, -y_2)$) as clarified in (19) and (20) respectively as follows:

$$C_{m1} = C_{m2} - k_{s2} \bmod p \quad (19)$$

$$P_m = C_{m1} - k_{s1} \bmod p \quad (20)$$

After the dual decryption have been implemented on the received points, converting the decryption points into data of the audio file is the next step.

Here the addition and subtraction operations does not mean simple operations, which they are in algebra, rather they are addition and subtraction, the points in ECC. So, the subtraction operation is the same addition operation, just take the inverse of $y - coordinate$ i.e, the inverse of $P(x, y) = -P(x, -y) = (x, -y + p)$.

VII. RESULTS AND DISCUSSION

The experiments run under Windows 10 professional operating system, Intel(R) Core(TM) i5-3210M CPU @ 2.50GHz, 4 GB random access memory and 64-bit system type. The visual studio 2010 (C# programming language) is used to evolve the proposed system. Table I shows the randomness of different AES key sizes. The inputs setting for generating AES key from the PBKDF2 algorithm are:

- **Password** is composed of 1024-byte that generated from the LCG algorithm, where the inputs to LCG algorithm are ($p_1 = 10247, q_1 = 8161, b_1 = 1, n_1 = 1024$).
- **Salt Value** is " 2018-02-08 02:11:56 AM ".
- Iteration Count is 10000.
- **Length of key** is (128,192,256) bits.

TABLE I. THE RANDOMNESS TESTS FOR DIFFERENT AES KEYS FROM THE PBKDF2 ALGORITHM

Size of AES keys	Frequency Test	Frequency within a block Test	Run Test
128-bit	0.193931	0.281692	0.079932
192-bit	0.288844	0.598714	0.902214
256-bit	0.900524	0.920130	0.937808

TABLE II. EXECUTION TIME AND THREE RANDOMNESS TESTS FOR THE H-AES-LCG GENERATOR

Size of output from H-AES-LCG generator	Execution Time	Frequency Test	Block Test	Run Test
128-bit	0.576 sec	0.859684	0.857480	0.725681

Table II illustrates the CPU time and three randomness tests for the H-AES-LCG generator, the setting inputs for the family of SHA and AES algorithm are:

- Password is " K@vvTHer9064Isaa ".
- The amount of rotation is 139.
- Salt Value is " 2018-02-08 02:11:56 AM ".

The inputs setting for the LCG algorithm to generate the sequence numbers are ($p_2 = 7933, q_2 = 4093, b_2 = 1, n_2 = 128$).

Fig. 2 demonstrates a comparison between the time of the traditional method for generating all points on a curve and the proposed method for selecting two base points from the curve in a secure manner, that have been tested on the different curves (i.e., different sizes of the domain parameters (p, a, b)).

Tables III and IV illustrate the execution time for dual encryption/decryption process in ECC, different sizes of audio models and various elliptic curves are used to test the results.

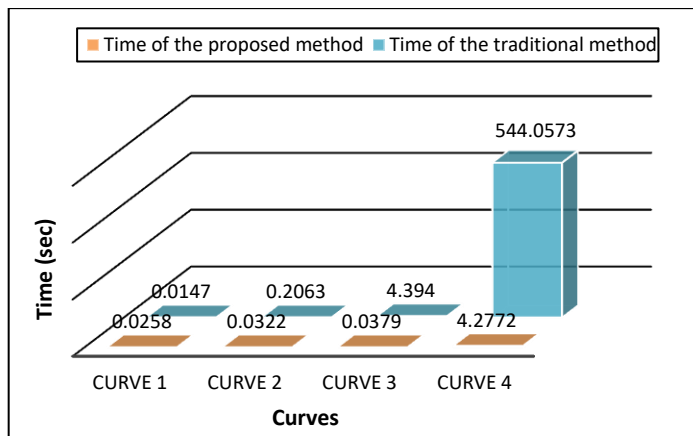


Fig. 2. A comparison between the time of the traditional and proposed method.

TABLE III. EXECUTION TIME FOR THE DUAL ENCRYPTION IN ECC

Prime Number	a,b	Record Time in sec	Data Size KB	Sample Rate Hz	Execution Time in sec
18787	3921,479	0.13	2.79	22050	0.0052
		0.251	5.4	22050	0.0103
		0.323	3.46	11000	0.0062
		0.420	4.5	11025	0.00805
		0.551	5.9	11025	0.008
		0.655	7.04	11025	0.0125
210487	155,180	0.13	2.79	22050	0.0052
		0.251	5.4	22050	0.0111
		0.323	3.46	11000	0.00705
		0.420	4.5	11025	0.009
		0.551	5.9	11025	0.0117
		0.655	7.04	11025	0.0144

TABLE IV. EXECUTION TIME FOR THE DUAL DECRYPTION IN ECC

Prime Number	a, b	Record Time in Sec	Data Size KB	Sample Rate Hz	Execution Time in sec
18787	3921,479	0.13	2.79	22050	0.0055
		0.251	5.4	22050	0.0108
		0.323	3.46	11000	0.0063
		0.420	4.5	11025	0.0081
		0.551	5.9	11025	0.0102
		0.655	7.04	11025	0.0125
210487	155,180	0.13	2.79	22050	0.006
		0.251	5.4	22050	0.01105
		0.323	3.46	11000	0.00705
		0.420	4.5	11025	0.0094
		0.551	5.9	11025	0.0119
		0.655	7.04	11025	0.0151

A. Security Analysis

Security analysis is a fundamental process to guarantee the power of the cryptography mechanism; the following metrics are used to evaluate the performance of the proposed system:

1) Entropy

The entropy is considered as one of the most significant measurements to measure the degree of secrecy (randomness). Where, a maximum entropy value in the first order entropy reach to 8 and 16 in the second order entropy, to compute the first order entropy $1^{st} E(\tau)$ and second order entropy $2^{nd} E(\tau)$ of a source τ are formulated as in the following equations [18][19]:

$$1^{st} E(\tau) = - \sum_{i=0}^{n-1} P(\tau_i) \log(P(\tau_i)) \quad (21)$$

$$2^{nd} E(\tau) = - \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} P(\tau_i, \tau_j) \log(P(\tau_i, \tau_j)) \quad (22)$$

Where, the total number of τ denotes by n , $P(\tau_i)$ referred to the probability of appearance of τ_i . Table V demonstrates the entropy values for different audio files where the curve parameters are: $p=210487, a=155, b=180$.

TABLE V. MEASURE THE ENTROPY OF AUDIO FILES BEFORE AND AFTER DUAL ENCRYPTION

Audio Data Size in KB	Sample Rate in Hz	Entropy Before Encryption		Entropy After Encryption	
		1 st Entropy	2 nd Entropy	1 st Entropy	2 nd Entropy
2.79	22050	5.30	10.58	7.93	15.81
3.46	22050	6.98	13.79	7.97	15.92
4.5	11025	6.01	11.87	7.94	15.80
5.4	11000	6.82	13.44	7.94	15.83
5.9	11025	5.84	11.50	7.91	15.69
7.04	11025	6.10	12.18	7.96	15.91
16.8	11025	6.26	12.38	7.98	15.97
28.6	16000	5.93	11.84	7.99	15.98
34.2	11000	6.32	12.61	7.99	15.98
47.06	22254	5.93	11.84	7.98	15.94

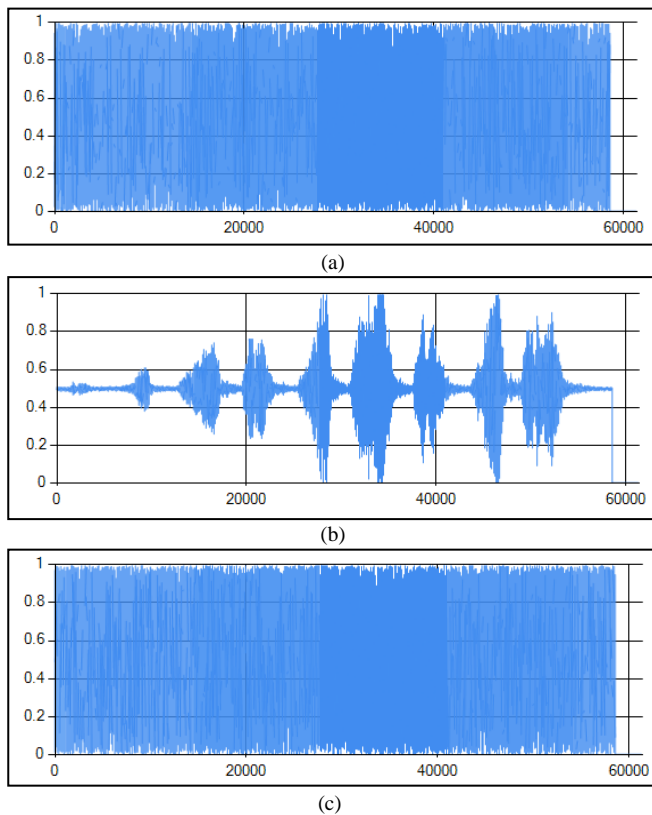


Fig. 3. (a) Cipher-audio; (b) Decrypted with right keys (k_{s1} , k_{s2});(c) Decrypted with key as ($k_{s1}-1$, $k_{s2}-1$).

2) Key Sensitivity Analysis

A fundamental feature of a good cryptosystem is a key sensitivity that assures the cryptosystem is secure contra the brute force attack. The key sensitivity is any simple change in the encrypted key gives a different result in recovering the plain-text from a cipher-text. Assume the user A sends a message (audio file) to the user B, the original session keys are $KS_1 = (2481,1143)$, $KS_2 = (4071,391)$. A slight alteration in the private keys (i.e., $n_{b1} - 1$, $n_{b2} - 1$) produce the different sessions keys ($KS_1 = (2937,907)$, $KS_2 = (3657,864)$) that lead to recovering a different message. Fig. 3 clarifies the decryption operation with the correct and wrong session keys, respectively.

3) Cryptanalysis Attacks

The intruder attempts an analysis of encrypted message by exploiting some characteristics in cipher-text such as repeated cipher or knowing some pairs of plain-text/cipher-text. In addition, he/she tries to identify the nature of the algorithm to recover a plain-text or key [1].

Our Suggested Solution: The diffusion property enhanced in the proposed system by using the Exclusive-OR and Circular Shift that performed to increase strength against cryptanalysis attack. Besides, it uses the dual encryption to increase immunity against any expected attack.

4) Attacks on ECDLP

a) *Brute-Force Attack (Exhaustive Search):* An intruder seeks to compute all the points that generated from the base

point (G) until a public key (P_m) is obtained to recognize the private key (P_r). The Seeking time in this type of attacks depends on the order (the number of points in a subgroup (G)) denoted by O_r , the large the order ($O_r \geq 2^{80}$) makes the computational infeasible [1],[2].

b) *Pohling-Hellman and Pollard's rho Attacks:* These attacks in order to accelerate the calculation of the ECDLP, the countermeasure of these attacks the (O_r) should be prime and $O_r > 2^{160}$ [2].

Our Suggested Solution: It is choosing the public elements in a secure manner to avoid all recognized attacks, instead of a large number of points; it requires a high computational cost.

5) Man in the Middle Attack

When two parties exchange the public keys for each them an intruder intercepts the transmission to occupy the public keys in order to generate a fake shared key between it and each party through exploits the public parameters (prime number, base point, their public keys) [1].

Our Suggested Solution: In our methodology, it frustrates this attack by selecting the global parameters randomly and in a secure manner this helps us even if an intruder knows a public key impossible to obtain the private key because there are unknown parameters to solve ECDLP.

VIII. CONCLUSIONS

This paper implements a new design for the ECC cryptosystem in random, efficient and secure manner based on the H-AES-LCG generator function, besides it chooses the domain parameters of the ECC within a given safe mechanism in order to defeat all organized attacks on the ECDLP. The ECDH method is used to make the communication between two parties more secure during the key exchanged process. In addition, the encryption process implemented in double or dual stages, the aim of this is to provide secure transmission for the audio messages and increase immunity against any attack. The proposed methodology is faster, more secure and provides many positive aspects such as enhancements in the key exchange compared with Diffie-Hellman key exchange and ECC performance. In addition, the (H-AES-LCG) is useful for generating encryption key for some algorithms, a slate value for the hash function, a prime number for the RSA algorithm or generates the domain parameters for ECC in a random and safe style.

REFERENCES

- [1] W. Stallings, Cryptography and Network Security Principles and Practice, 6th ed., Pearson Education: United States of America, 2014.
- [2] D. Hankerson, S. Vanstone, and A. J. Menezes, Guide to elliptic curve cryptography, Springer Science & Business Media: New York, 2004.
- [3] A. Soleymani, M. J. Nordin, and Z. M. Ali, "A novel public key image encryption based on elliptic curves over prime group field," Journal of Image and Graphics, vol. 1, no. 1, 2013.
- [4] R. Singh, R. Chauhan, G. Ritu, K. Vinit, and P. Singh, "Implementation of elliptic curve cryptography for audio based application," International Journal of Engineering, vol. 3, no. 1, 2014.
- [5] A. Luma, B. Selimi, and L. Ameti, "Using elliptic curve encryption and decryption for securing audio messages," in Transactions on Engineering Technologies, GC.Yang, SI. Ao, L. Gelman (eds.), Springer, Dordrecht, 2015.

- [6] M. Kumar, A. Iqbal, and P. Kumar, "A new RGB image encryption algorithm based on DNA encoding and elliptic curve Diffie-Hellman cryptography," *Signal Processing*, vol. 125, pp. 187–202, 2016.
- [7] X. Fang and Y. Wu, "Investigation into the elliptic curve cryptography," In *Information Management (ICIM)*, 2017 3rd International Conference on. IEEE, pp. 412–415, 2017.
- [8] K. E. Abdullah, N. H. M. Ali "A secure enhancement for encoding/decoding data using elliptic curve cryptography," *Iraqi Journal of Science*, vol. 59, no. 1A, pp. 189–198, 2018.
- [9] A. Hafsa, N. Alimi, A. Sghaier, M. Zeghid, and M. Machhout, "A hardware-software co-designed AES-ECC cryptosystem," In *Advanced Systems and Electric Technologies (IC_ASET)*, 2017 International Conference on. IEEE, pp. 50–54, 2017.
- [10] Y. Lin, K. Kang, and Y. Shi, "Research on encryption model based on AES and ECC in RFID," In *Computer Sciences and Applications (CSA)*, 2013 International Conference on. IEEE, pp. 9–13, 2013.
- [11] Q. Dang, "Changes in federal information processing standard (FIPS) 180-4, secure hash standard," *Cryptologia*, vol. 37, no. 1, pp. 69–73, 2013.
- [12] R. P. McEvoy, M. F. Crowe, C. C. Murphy, and P. W. Marnane, "Optimisation of the SHA-2 family of hash functions on FPGAs," In *Emerging VLSI Technologies and Architectures*, 2006. IEEE Computer Society Annual Symposium on. IEEE p. 6-pp, 2006.
- [13] T. Van, C. Henk, and S. Jajodia, *Encyclopedia of Cryptography and Security*, 2nd ed., Springer Science + Business Media: New York, 2011.
- [14] B. Kaliski, "PKCS# 5: Password-based cryptography specification version 2.0," 2000. [Online] Available: <https://tools.ietf.org/pdf/rfc2898.pdf>
- [15] M. S. Turan, E. B. Barker, W. E. Burr, and L. Chen, "Recommendation for password-based key derivation part 1: storage applications," NIST Special. Publication, vol. 800, pp.132, 2010.
- [16] M. Dürmuth, T. Güneysu, M. Kasper, C. Paar, T. Yalcin, and R. Zimmermann, "Evaluation of standardized password-based key derivation against parallel processing platforms," in *European Symposium on Research in Computer Security*, pp. 716–733, 2012.
- [17] A. J. Menezes, T. Okamoto, and S. A. Vanstone, "Reducing elliptic curve logarithms to logarithms in a finite field," *IEEE Transactions on Information Theory*, vol. 39, no. 5, pp. 1639–1646, 1993.
- [18] J. Payingat and D. P. Pattathil, "Pseudorandom bit sequence generator for stream cipher based on elliptic curves," *Mathematical Problems in Engineering*, vol. 2015, 2015.
- [19] A. A. A. El-Latif and X. Niu, "A hybrid chaotic system and cyclic elliptic curve for image encryption," *AEU-International Journal of Electronics and Communications*, vol. 67, no. 2, pp. 136–143, 2013.

Classification of Affective States via EEG and Deep Learning

Jason Teo, Lin Hou Chew, Jia Tian Chia, James Mountstephens

Faculty of Computing and Informatics

Universiti Malaysia Sabah

Jalan UMS, 88400 Kota Kinabalu, Sabah, Malaysia

Abstract—Human emotions play a key role in numerous decision-making processes. The ability to correctly identify likes and dislikes as well as excitement and boredom would facilitate novel applications in neuromarketing, affective entertainment, virtual rehabilitation and forensic neuroscience that leverage on sub-conscious human affective states. In this neuroinformatics investigation, we seek to recognize human preferences and excitement passively through the use of electroencephalography (EEG) when a subject is presented with some 3D visual stimuli. Our approach employs the use of machine learning in the form of deep neural networks to classify brain signals acquired using a brain-computer interface (BCI). In the first part of our study, we attempt to improve upon our previous work, which has shown that EEG preference classification is possible although accuracy rates remain relatively low at 61%-67% using conventional deep learning neural architectures, where the challenge mainly lies in the accurate classification of unseen data from a cohort-wide sample that introduces inter-subject variability on top of the existing intra-subject variability. Such an approach is significantly more challenging and is known as subject-independent EEG classification as opposed to the more commonly adopted but more time-consuming and less general approach of subject-dependent EEG classification. In this new study, we employ deep networks that allow dropouts to occur in the architecture of the neural network. The results obtained through this simple feature modification achieved a classification accuracy of up to 79%. Therefore, this study has shown that the use of a deep learning classifier was able to achieve an increase in emotion classification accuracy of between 13% and 18% through the simple adoption of the use of dropouts compared to a conventional deep learner for EEG preference classification. In the second part of our study, users are exposed to a roller-coaster experience as the emotional stimuli which are expected to evoke the emotion of excitement, while simultaneously wearing virtual reality goggles, which delivers the virtual reality experience of excitement, and an EEG headset, acquires the raw brain signals detected when exposed to this excitement stimuli. Here, a deep learning approach is used to improve the excitement detection rate to well above the 90% accuracy level. In a prior similar study, the use of conventional machine learning approaches involving k-Nearest Neighbour (kNN) classifiers and Support Vector Machines (SVM) only achieved prediction accuracy rates of between 65% and 89%. Using a deep learning approach here, rates of 78%-96% were achieved. This demonstrates the superiority of adopting a deep learning approach over other machine learning approaches for detecting human excitement when immersed in an immersive virtual reality environment.

Keywords—*Neuroinformatics; emotion classification; preference classification; excitement classification;*

electroencephalography (EEG); deep learning; virtual reality; dropouts.

I. INTRODUCTION

We have conducted a number of prior investigations into the use of electroencephalography (EEG) as a method for passively monitoring the brainwaves of users as they are exposed to 3D visual stimuli as well as immersive stimuli and then using different machine learning algorithms to predict their preferences among the various visual stimuli [1], [2]. In the first part our study, we focus on human preference classification. The ability to passively identify the preferences of users as they are being presented with different stimuli will have novel and significant applications in various choice-based domains such as neuromarketing, affective entertainment, virtual rehabilitation and forensic neuroscience.

In our early work with a small set of five test subjects, good classification rates of up to 80% were attained using simple k-nearest neighbor (kNN) classifiers [1]. However, when the number of test subjects was increased to 16, the noise arising from inter-subject variability became a substantial factor which made the classification process significantly more challenging [2]. While most studies generally deal only with intra-subject variability where for each user, retraining is required before classification testing. We attempt a cohort-wide classification to enable direct applications to new users without the need for per-person pre-training before classification usage. In the our expanded study, classification rates for the large majority of conventional classifiers such as kNN, support vector machines, Naïve Bayes, Random Forest, C4.5 and other rule-based classifiers were only between 56-60%. The best classification result obtained from this comparative study was using deep neural networks at 64% [2].

The second part of our study focuses on excitement detection in immersive environments since much less is known about human emotion recognition in fully immersive environments such as virtual reality (VR). VR environments provide an arguably more effective emotion stimulating environment since users are fully immersed in the stimulus environment without any distracting views and/or other stimuli such as those present when using conventional displays such as computer and TV screens. Furthermore, users are free to move their heads to fully view their VR environments, which is more akin to their real-world viewing experience, hence suggesting the possibility of greater emotional response correlation with real life experiences. Additionally, as VR continues to garner

widespread adoption among everyday consumers, the ability to incorporate an effective emotion recognition system for VR applications will open up a wealth of novel interactions between the user and the VR experience particularly in the video gaming, live events, video entertainment, retail, real estate, healthcare, education, military and engineering domains [3].

As such, the main objective of the study is to investigate the various architectural tuning of the deep neural networks for improving the classification rates of our EEG-based preference classification as well as excitement classification task. Section II presents the background on emotion classification. Section III presents our approach to EEG-based preference and excitement classification using visual stimuli. Section IV presents the results of our investigations and Section V concludes the paper with some future avenues for expanding upon the current work.

II. BACKGROUND

A. Emotion Modeling and Classification

Emotion classification entails the use of various physiological signals and markers in an attempt to identify different emotions such as the user being in a state of anger, disgust, happiness, sadness, fear, anxiety, excitement and surprise among other [4], [5]. Some commonly measured bio-signals include the heart rate, skin conductance, pupil dilation, respiration rate and also brainwaves, which is also known as EEG [6], [7].

EEG-based emotion classification typically involves the measurement of the millivolt-range electrical signals through the placement of a number of electrodes on the scalp of the user, the waveforms of which are then spectrally transformed into features used by machine learning algorithms trained on labelled data to predict the emotion currently being sensed. Numerous studies have shown that classifications for various emotions can be reliably obtained using EEG.

B. Emotion Classification of Preferences

Preference classification can be considered a sub-task of emotion classification. This more specific task entails the identification of a user's like or dislike when presented with a stimulus. Preference classification is generally considered to be more challenging to classify compared to other emotions that are more strongly evoked such as anger or sadness.

The very large majority of EEG-based preference classification has been conducted using music as the stimulus [8], [9]. There have been very limited studies done using 2D images [10], [11] whereas our earlier studies were the first to implement rotating virtual 3D images as the stimulus [1], [2]. Furthermore, preference classification, which is already more challenging compared to other forms of emotion classification due to its comparatively weaker evocation, is rarely studied as a cohort-wide classification task. EEG-based emotion classification with large-sized cohorts will typically yield significantly lower accuracy rates due to inter-subject [12] and as well as intra-subject variability [13]. Doing so requires the classifier to be able to overcome inter-subject variability in addition to intra-subject variability of the users' EEG signal.

Consequently, the weak signal evocation and inter-subject variability make EEG preference classification a very challenging classification task.

C. Extraction of Features from EEG Signals

Emotion modeling using machine learning approaches can be categorized into three broad domain classes: (i) time, (ii) frequency, and (iii) time and frequency combination. Time-based emotion modeling employs the detection of event-related potentials (ERPs). Of these, they can be further divided into groups that are detected based on whether they are having short, medium or long post-latency exposures after stimuli presentation. Emotion classification for valence and arousal produced accuracy rates of 55.7% for arousal and 58.8% for valence [14] when using these ERP-based methods.

The classification of emotions based on the frequency domain is achieved through the learning of features obtained power spectrum analysis, producing the canonical delta, theta, alpha, beta and gamma frequency bands. Emotion classification for the preference of music produced an accuracy of 74.8% with linear support vector machines (SVMs) using the preprocessed features obtained through the Common Spatial Patterns (CSPs) method [15]. Emotion classification for the preference of music via preprocessed features obtained from a using a conventional Fast Fourier Transform (FFT) produced a classification accuracy rate of 85.7% using SVMs [16]. Radial SVMs were used in the only published emotion classification of preferences not using music stimuli, in this case for 2D image preferences using power spectrum analysis where the classification outcome produced an accuracy of 88.5% [17].

From the perspective of using a combination of time and frequency (TF) leverages on the power spectrum analysis at predefined time periods that encompass the whole duration of the post stimuli period for measuring brain activity. Several conventional machine learning algorithms were used to conduct emotion classification tasks employing three distinctly different TF analysis methods were studied to identify the preference for music. Here, it was observed that the k-Nearest Neighbors (kNN) machine learning approach produced the overall best outcome with an accuracy of 86.5% [18]. The same group of researchers then conducted a follow-on investigation utilizing a much finer-grained approach which attempted to categorize the emotion stimuli into two groups: (i) familiar versus, (ii) unfamiliar music. In this later study, using a kNN machine learning approach, they managed to produce a much higher emotion classification accuracy of 91.0% [19]. Emotion modeling for the preference of music using TF approaches used a Short-Time Fourier Transform (STFT) and using a kNN machine learning approach produced emotion classification accuracy rates of 98.0% [20].

D. Preference Classification using Deep Learning Approaches

The preferences of 32 participants for the viewing of music video clips was attempted using deep learning via the Deep Belief Networks (DBNs) approach [21]. DBNs accomplish deep learning through the stacking of various Restricted Boltzmann Machines (RBMs) on top of each other. In this method of deep learning, the output obtained from a lower-

level RBM is subsequently utilized to serve as the input to a higher-level RBM. This process is continued progressively through deeper and deeper layers thus forming a multi-layer stacking of these so-called RBMs. An average emotion classification accuracy of 77.8% was obtained where this method performed significantly better than various types of different SVMs as well as standard non-stacked RBMs.

A limited study involving only 6 subjects was reported for the emotion classification of participants when presented with the stimuli of viewing a number of short video clips for the elicitation of emotions with positive or negative valences [22]. In a novel approach for the emotion classification task which utilizes only the top five EEG recording electrodes, the investigation produced emotion classification accuracies of 87.6% using DBN's with this novel critical feature channel selection method. These results were observed to perform better than Extreme Learning Machines (ELMs) as well as SVMs and at the same time was observed to perform significantly better than the kNN machine learning approach. However in both of these two reported studies, it is important to point out that the training and classification prediction tasks were accomplished on a per-subject basis and not over the entire cohort of participants, which means that this only caters for intra-subject variability and not inter-subject variability. In other words, these two studies utilized an approach that requires the retraining of machine learning classifiers during the training phase whenever there is a new participant before the emotion classification prediction task can be performed. Essentially what this intra-subject or subject-dependent method employs is an approach that bypasses the difficulty of handling inter-subject variability and only caters for intra-subject variability, which means that it will not work for subject-independent classification tasks.

From the literature survey, there was only one paper found in which the deep learning approach was used in emotion modeling to classify preferences in a subject-independent methodology. Here it was reported that using a combination of unsupervised learning employing stacked autoencoders (AEs) in conjunction with the supervised learning of softmax machine learning classifiers was able to perform prediction of the emotional states for 32 participants for valence and arousal. Nonetheless, this paper reported the requirement of utilizing an extremely large number of hidden neurons in the deep learning classifier. It is interesting to note that the authors themselves alluded to the fact that an extended amount of computational time was utilized during the training phase with such an approach. Subsequently the authors hybridized this approach with feature preprocessing routines employing Principal Component Analysis (PCA) as well as Covariate Shift Adaptation (CSA) during the pre-learning process. However, even with the extended processing time and numerous augmentations with supplementary preprocessing, the emotion modeling was only able to produce very low prediction accuracy rates of 53.4% and 52.0% for valence and arousal classification, respectively from this subject-independent approach using leave-one-out cross-validation (LOOCV) [23]. What this study clearly demonstrates is the fact that inter-subject variability very significantly and critically adds tremendous difficulty to the classification of emotions based on

preferences when compared against the much more common and significantly easier prediction task of subject-dependent studies that only caters for intra-subject variability in the learning of the EEG-based emotion modeling.

E. Classification of Affective States in Virtual Reality and Mixed Reality Environments

There have been very few studies that have conducted human emotion recognition that have used virtual reality environments as the stimulus. To the best of our knowledge, there has yet to be any study that uses solely EEG to detect human emotions using purely VR stimulus.

Wu et al. [24] used a Virtual Reality Stroop Task (VRST) from the Virtual Reality Cognitive Performance Assessment Test (VRCPAT) to detect arousal levels in their attempt to identify various affective/cognitive states. A number of VR stimuli were presentations with various levels of arousal were selected from the VRST. It was shown that a relatively high classification accuracy rate of 96.5% using support vector machines (SVM) could be achieved through VR stimuli. However, the study used an elaborate and involved sensor setup with a wide range of psychophysiological responses which included skin conductance level, respiration, ECG, as well as EEG were used to conduct the emotion recognition task. As such, it remains unknown if a much simpler setup involving EEG alone would be feasible in achieving successful emotion recognition.

Massari et al. [25] and Kovacevic et al. [26], respectively used mixed reality stimuli to conduct brain state recognition based solely on EEG signals as input to the classification system. Massari et al. utilized their proprietary eXperience Induction Machine (XIM) as the mixed reality stimuli system to classify different brain states for spatial navigation, reading and calculation, achieving the best results of 86% using linear discriminant analysis [25]. Kovacevic et al. implemented an EEG-based mental state recognition system as part of an immersive and interactive multi-media science-art installation using the recognition of relaxation and concentration mental states of its participants to determine the audio-visual output of a dome-based artistic installation comprising video animations that were projected on to the 360° surface of the semitransparent dome as well as the generation of soundscapes based on pre-recorded sound libraries and live improvisations [26]. Although both these studies utilized EEG solely as the feature input, these studies were not specifically classifying emotional states and both were utilizing mixed reality stimuli rather than pure virtual reality stimuli. As such, it remains unknown whether a purely VR-based stimulus system could be successfully used for emotion recognition.

III. METHODOLOGY

A. Experimental Setup for Preference Classification

Emotion classification entails the use of various physiological signals and markers in an attempt to identify different emotions such as the user being in a state of anger. In this, first of this investigation for preference classification, 16 subjects (8 female and 8 male, mean age = 22.44) were involved where all the participants had corrected-to-normal or normal vision. Furthermore, they were asked and confirmed

to be free of any known history of psychiatric illnesses prior to the participation in the study. The participants were briefed on what to expect in terms of the BCI equipment that was to be used during the data acquisition phase before the actual experimentation was to proceed. The EEG acquisition device was a brain-computer interface (BCI) headset called the ABM B-Alert X10, which has nine active electrodes, namely the POz, Fz, Cz, C3, C4, F3, F4, P3 and P4 channels according to the standard 10-20 naming convention where a subject participant wearing the said BCI headset is depicted in Fig. 1. MATLAB, Java and R were the three programming languages used. The visual stimuli were developed and displayed using the Java programming language. Integration between the visual stimuli and the BCI headset was accomplished by implementing the MATLAB programming language with the B-Alert X10's SDK. Finally, the statistical programming language R was used for the signal preprocessing phases, feature extraction, and finally for the training and prediction classification tasks.

The data acquisition processes experienced by the participants are as shown in Fig. 2 where during the commencement of the data acquisition process, a blank screen of three seconds is shown to the participant to obtain the base resting brain signal in order to avoid any brain activities related to the previous stimuli during the actual emotion modeling trial phase. After this blank screen, there will be between five to fifteen of actual viewing time for the 3D stimuli where the minimum viewing time and maximum viewing time is set between five and fifteen seconds respectively. The participant is allowed to commence to the following rating state based on their own choosing after the minimum viewing duration time of five seconds while once the maximum viewing duration time is up, the system will proceed by default to the next rating state. The purpose of implementing this particular method of the data acquisition process flow is to allow the participant to decide on their own accord during the stimuli viewing time so as to mitigate the possibility of boredom from setting in and making the participant fatigued while viewing the stimuli during the data acquisition process since requiring the participant to continuously view only at fixed intervals in a repetitive manner for the purposes of rating the stimuli could possibly cause the participant to experience boredom which will subsequently lead to further fatigue towards the end of the data acquisition process. As such, since the participant is no longer required to just wait until the maximum set and fixed time in order to conduct the rating, this essentially provides the participant with the freedom and ability to shift to the following visual stimuli, which will potentially save some overall viewing time and at the same time prevent the participant from fatiguing. A rating system containing a discrete scale of 1-5, where 1 represents like very much; 2 represents like; 3 represents undecided; 4 represents do not like; and finally 5 represents do not like at all, is shown to the participant at the conclusion of the visual shape stimuli viewing period.



Fig. 1. The medical-grade 9-channel EEG acquisition device is shown being worn by a participant in the study.

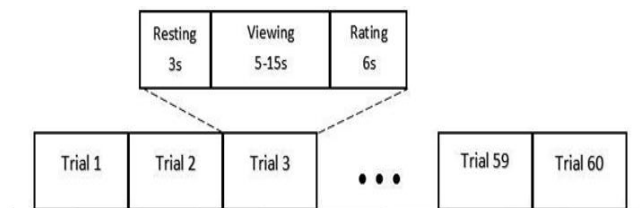


Fig. 2. The flow of the data acquisition process as experienced by the participants during the experimentation.

$$r(\theta) = \frac{1}{n1 \sqrt{\left(\left(\frac{1}{a} \cos\left(\frac{m}{4} \theta\right)\right)\right)^{n2} + \left(\left(\frac{1}{b} \sin\left(\frac{m}{4} \theta\right)\right)\right)^{n3}} \quad (1)$$

The Gielis Superformula is used to generate three-dimensional shapes which were used as the visual stimuli in this study and had the visual appearances of a bracelet-like virtual generated [27], the mathematical formula of which is as shown in (1). Our main reason for choosing this shape as the three dimensional visual stimuli for evoking emotions is to determine the aesthetic quality of jewelry-type objects since visual aesthetic quality is primarily the key motivating factor when one decides whether or not make a purchase of such an item. By modifying the various superformula parameters, the generation of different and myriad natural three dimensional virtual shapes can be generated.

Sixty different bracelet-like shapes generated and used in this study is as shown in Fig. 3, which were generated by utilizing different parameters with randomly generated values in the superformula. Through preliminary testing, different ranges of suitable parameter values were chosen to synthesize virtual three dimensional shapes that possess visual characteristics of a bracelet-like shape. These three dimensional bracelet-like shapes were then shown to the participants virtually on a computer. The visual system allowed the presentation of the three dimensional virtual shapes with rotations on different axes of the presented stimuli so that it could be viewed at different angles in order for the participant to be able to fully visualize the generated three dimensional bracelet-like shapes.

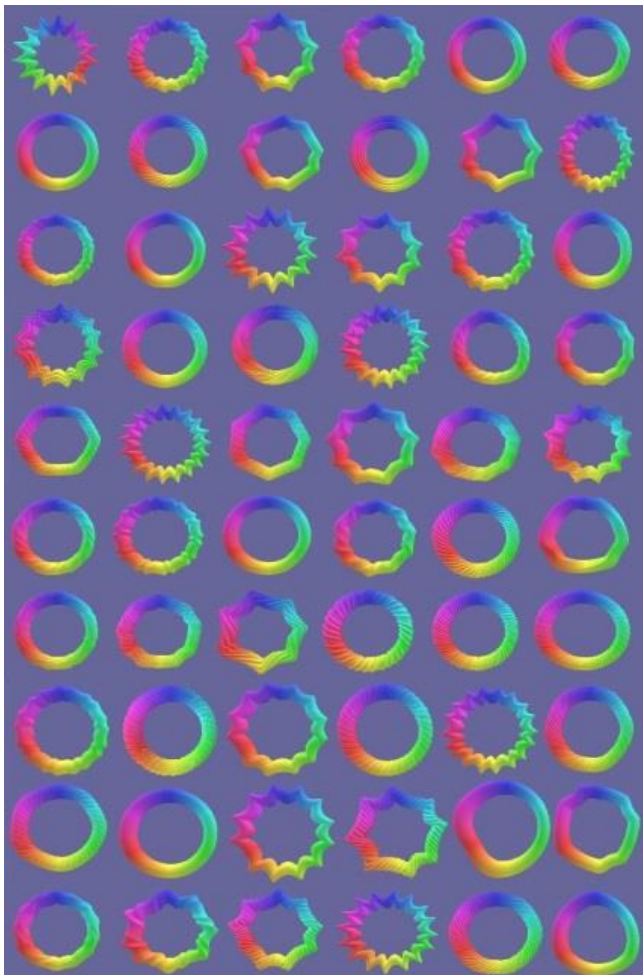


Fig. 3. The Gielis Superformula used to produce 60 bracelet-like shapes.

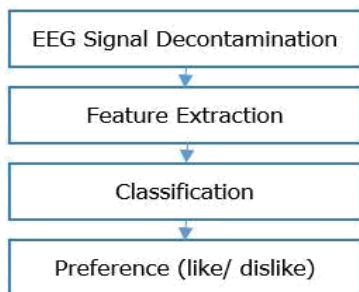


Fig. 4. Summary of the signal processing process flow.

The major processes for preference classification are as shown in Fig. 4. Firstly, environmental and physiological artifacts are always present in EEG signal recordings and require decontamination. The SDK in the MATLAB programming language provided by ABM for the B-Alert X10 BCI headset automatically provides this decontamination function. A 50Hz notch filter removes environmental artifacts while five physiological artifacts comprising electromyography (EMG), eye blinks, excursions, saturations, and spikes are similarly removed automatically in real-time. The eye excursions, saturations, and spikes are replaced by zero values where they are later filled in using spline interpolation.

Subsequently, a Short-Time Fourier Transform (STFT) is then used to transform the decontaminated EEG signals into the TF domain where it decomposes each of the nine BCI channels into five spectral bands, which are the delta 1-3Hz, theta 4-6Hz, alpha 7-12 Hz, beta 13-30 Hz, and gamma 31-64Hz bands. These five bands across the nine channels thereby provides a total of forty-five input features. The brainwave recordings from the 16 participants where each viewed the sixty 3D visual stimuli of the bracelet-like shapes generated 960 observations altogether. However, only 208 observations were used during the training and prediction classification process. These were the strongest ratings on the ratings scale of 1, which represented like very much, and 5, which represented do not like at all, respectively. A final dataset matrix comprising forty-seven feature columns consisting of the observation ID reference, participant rating, and each of the forty-five TF features, over two hundred and eight rows of selected observations served as the training and testing data for the respective machine learning classifiers. Moreover, the subjects' baseline readings acquired while in the resting state were subtracted from the stimuli viewing state values before the values were utilized in the prediction classification process.

The deep neural networks utilized were set to two hundred hidden neurons within each of the two hidden layers using the uniform adaptive method [28] for weight matrix initialization. Preliminary experimentation showed that this setup with the number of hidden layers as well as the number of hidden neurons per layer provided the optimal settings for this preference prediction task. Cross-entropy [29] was used as the error function during the 10-fold cross-validation, which was conducted for 10 epochs in each of the cross-validation steps.

B. Experimental Setup for Excitement Classification

Fig. 5 describes the overall approach adopted in conducting this second study of excitement classification in virtual reality consisting of a number of distinct phases, each of which will be explained in the following subsections below.

1) Experiment Stimuli

In this project, the immersive stimuli were created using Google's Cardboard VR technology and a 360° video available on YouTube.com. The selected video was the experience video of the 360-degree ride on a roller coaster that is promising in eliciting excitement emotion with providing sensations of a roller coaster ride which drops from high peaks and high speed 360-degree turns. Screenshots of these stimuli is as shown in Fig. 6 and 7.

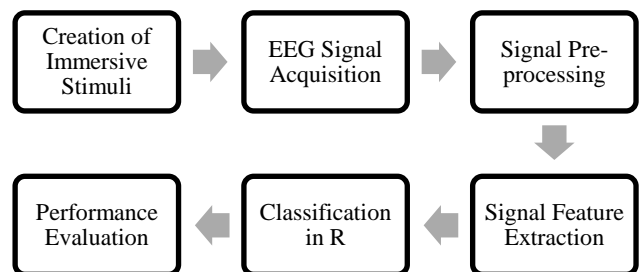


Fig. 5. Overall process of human mental states classification.



Fig. 6. Scenes in the stimuli video.



Fig. 7. Scenes in the stimuli video.

2) Experimental Test Subjects

A total of 24 human subjects (12 females, 12 males) participated in this study and had normal or corrected-to-normal vision with no history of psychiatric illness. The age of the subjects was in the range of 20 to 28 years old. During the experimental session, subjects were advised to sit comfortably on the chair without any restriction to head movements which being immersed in the virtual reality stimuli. An image of a test session in progress with the test subject wearing the VR headset and EEG headband is as shown in Fig. 8.



Fig. 8. Experiment setup using human test subject.

3) Data Acquisition Device

To increase the applicability of EEG-based predictive analytics in human mental state classification, the Muse brain sensing headband from Interaxon was used as the data acquisition device since it is trivial to set up and comes at a much lower cost compared to medical-grade conventional EEG devices. Conventional EEG devices such as the B-Alert X10 by ABM that uses adhesive sponge discs with the requirement of applying electrode gels require a significantly longer set up time for individual electrodes and such a setup often limits the behavioural freedom on the participants since there are numerous connecting wires to connect to the electrodes which severely restricts head movements as necessary in immersive VR environments. The Muse headband is extremely accessible as it is wireless, lightweight, flexible, adjustable and easy application. The wearer that puts on the Muse headband will not experience any limitations on their mobility as the Muse headband is connected through wireless Bluetooth technology for data transmission. The Muse headband has four dry electrode channels at international standard 10-20 coordinates of TP9, AF7, AF8 and TP10 as illustrated in Fig. 9. The earpieces of Muse are adjustable but the headband area with channel electrodes AF7, AF8 and reference channel Fpz are not flexible. Although the Muse headband is still new in the market as a commercialised EEG device, there have been other studies that have reported on its potential to be used as a research tool despite its limited numbers of electrodes and low signal resolution [30], [31].

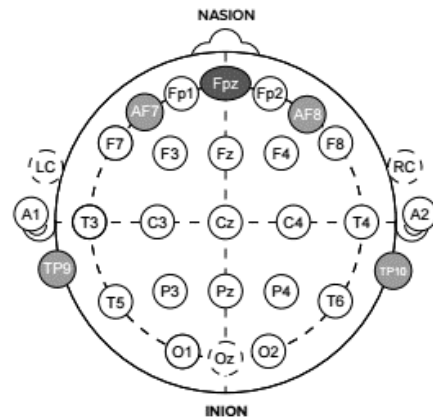


Fig. 9. Muse electrode locations by 10-20 International Standards.

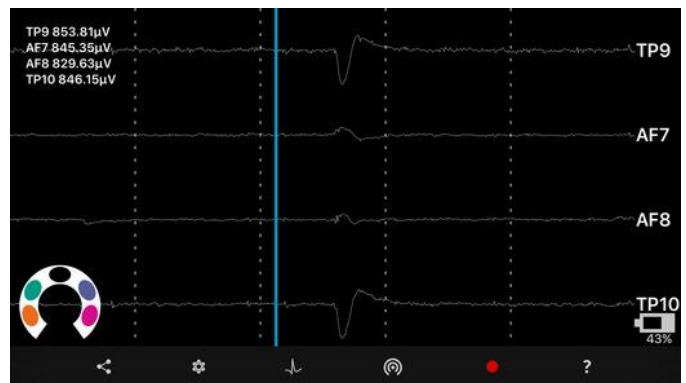


Fig. 10. Signal acquisition and recording screen of the Muse Monitor app.

4) EEG Recording

The EEG recordings were acquired using the Muse Monitor app available from Google’s Android Play Store. The EEG recordings were exported in CSV file format from Muse Monitor. The real-time EEG signal is recorded with an interval of 0.5 seconds, providing 256 data points per second for raw values to capture the minor changes of the brain rhythms. Fig. 10 shows the recording screen of the Muse Monitor app with the view of raw EEG values captured from each of the sensors (TP9, AF7, AF8, TP10) in microvolts (μV). In this study, two recordings were recorded from each experimental test subjects. One recording is for the “Rest” state where there is no video stimulus and subjects were asked to keep calm and breathe normally. The second recording is for the “Excited” state where subjects were wearing VR headset being immersed and stimulated by the 360° roller coaster video experience.

C. Signal Pre-processing and Feature Extraction

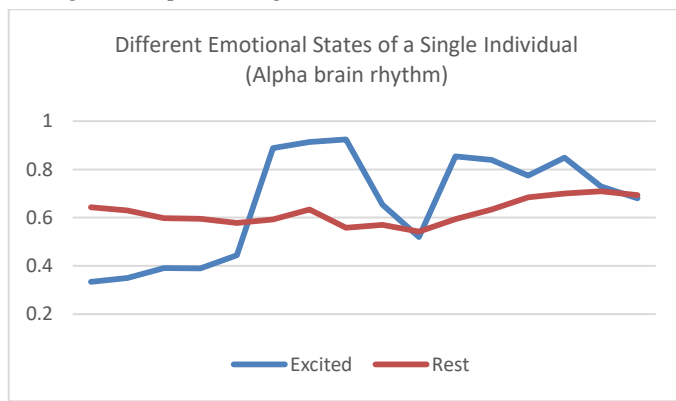


Fig. 11. Different emotional states of a single individual in Alpha brain rhythm representation.

Recorded EEG signals are always subject to artefacts and noise during acquisition. The common artefacts found are electromyography (EMG), eye-blinks, excursions, saturations, and muscle spikes. It is important to perform signal pre-processing to enhance the signal-noise power ratio [32]. The band of interest in this study are the frequency bands: delta (δ), theta (θ), alpha (α), beta (β) and gamma (γ), each reflecting different brain states of human experimental subject. Fast Fourier Transform (FFT) was used to convert the obtained EEG signal to a representation in the frequency domain based on Butterworth’s 4th Order Filter with different cut-off frequency thresholds to extract the five frequency bands [32] as in (2):

$$X_k = \sum_{n=0}^{N-1} x_n e^{-i2\pi k \frac{n}{N}} \quad (2)$$

where $k = 0, 1, 2, \dots, N-1$, X_k is the FFT coefficients, N is the total number of input EEG samples, n is the total number of points in FFT. There are two EEG recordings per subject, one is the “Rest” state, and the other one is the “Excited” state. Fig. 11 shows the two Alpha brain rhythms of different states from a single individual. In the “Rest” state recordings, a length of 16 data points was extracted for classification. While in the “Excited” state recordings, two sets of 16 data points were extracted in accordance with the two excitement eliciting

events in the video stimulus: (1) the drop of the roller coaster from the highest peak and (2) the high speed 360° degree turns of the roller coaster. In conclusion, 3 sets of data points were extracted from each experimental subject, giving a total of 72 objects for classification. The extracted data was then tabulated according to each band as features of interest for classification.

The classification work in this project was performed in the R environment as R is the leading statistical analysis tool that includes a large collection of packages that provides a wide variety of linear and non-linear modelling, classification function and etc. The main package used to build classification models was the ‘caret’ package as it has a consistent syntax for various machine learning methods. Additionally, the ‘caret’ package also provides an easy implementation to perform the 10-fold cross validation on the classification model. Since the ‘caret’ package made it easy to expand the range of tuning parameters of the machine learning methods, this experiment had systematically investigated various parameter settings for each classifier used.

1) *K-Nearest Neighbour (KNN)*: KNN is a simple and intuitive method of classifier used in many research works typically for classifying signals and images. KNN classifies objects based on the similarity between two instances to locate the nearest neighbour. The classifier will compare a newly labelled sample with the baseline data. The decision rule applied will vote where the new labelled sample will be assigned based on the class of the majority of the k-nearest neighbours.

2) *Support Vector Machine (SVM)*: SVM function attempts to find a hyperplane in between the groups of objects to classify them. The SVM operates by minimising the loss function as in (3):

$$\min_w \frac{1}{2} w^T w + C \sum_{i=1}^l \max(0, 1 - y_i w^T \phi(x_i)) \quad (3)$$

where w is the vector of weights, C is cost parameter, $\phi(x)$ is a kernel function applied on the input data.

The Radial Basis Function (RBF) kernel will be applied with SVM to enable operations performed in the input space rather than the potentially higher dimensional feature space [33] as in (4):

$$K(x, x') = \exp(-\gamma \|x - x'\|^2) \quad (4)$$

where $\|x - x'\|^2$ is the square of the Euclidean distance between the two vectors, γ is the kernel parameter, equivalent to $\frac{1}{2\sigma^2}$ where σ is a free parameter: the inverse kernel width for RBF kernel. There are two tuneable parameters in this function used: C and σ .

3) *Random Forest (RF)*: RF is an ensemble classifier that operates by constructing a multitude of decision trees [34]. The final predicted class for a test example is obtained by combining the predictions of all individual trees. The decision tree with controlled variance was constructed through a combination of bootstrap aggregation (bagging) and random feature selection. Each node in RF is split using the best

among a subset of predictors that are randomly chosen at the node. This strategy makes the classifier perform better compared to other classifiers such as SVM, neural networks and Linear Discriminant Analysis (LDA) and is robust against over-fitting [35]. There is only one tuning parameter for RF: mtry (number of variables randomly sampled as candidates at each split).

4) *Feed-forward Neural Network (NN)*: NN is the first and simplest form of artificial neural network developed. The network's information only moves in one direction, from the input nodes through the hidden nodes (if any) and to the output nodes in forward directions. There exist no cycles or loops in the network. The simplest design of NN is a single-layer perceptron network that consists only one layer of output nodes. The inputs are fed directly to the outputs via a series of weights and the output units are of the same form but with an output function [36]:

$$y_k = \phi_0 \left(\alpha_k + \sum_h \omega_{hk} \phi_h \left(\alpha_h + \sum_i \omega_{ih} x_i \right) \right) \quad (4)$$

The activation functions ϕ_0 and ϕ_h are taken to be the logistic function:

$$l(z) = \frac{\exp(z)}{1 + \exp(z)} \quad (5)$$

There are two tuning parameters in NN: size (number of hidden units) and decay (weight decay).

5) *C5.0 Decision Tree & Rule-based Model (C5.0)*: This algorithm was developed based on the C4.5 algorithm. C5.0 can be applied for classification as a decision tree or rule-based model. It supports boosting with any number of trials and can automatically winnow the attributes to remove those attributes that may be obstructive. For high-dimensional applications, this winnow features can lead to smaller classifiers and higher predictive accuracy while minimising the time required to generate rule sets. C5.0 has three tuning parameters: model (choose between decision tree and rule-based model), winnow (decision on whether predictor winnowing should be used) and trials (number of boosting iterations).

IV. RESULTS AND DISCUSSION

A. Preference Classification Result

Four distinct deep net architectures were tested, which were the standard deep nets, deep nets with dropouts only, deep nets with L1 regularizations only and finally deep nets with both dropouts and L1 regularizations. In L1 regularizations, λ is set at 10^{-5} . For dropouts, we set the hidden layer dropout probability at 0.5. For each of these architectures, we also

paired them with different activation functions for the hidden layers, which were the tanh, maxout and rectified linear unit (ReLU) activation functions. The rectified linear activation function [37] was used with an adaptive learning rate method [38]. The results of this specific part of the study have been previously published [39].

Table I presents the 10-fold cross-validation results obtained from using the various deep net architectures as well as with dropouts and L1 regularization terms. The best classification was obtained using the deep net with dropout architecture using rectified linear units for activation at 79.76%. The second best classification result was also obtained using the deep net with dropout architecture but using the tanh activation at 74.38%. This was followed next with the deep net architecture using both dropouts and L1 regularization with the rectified linear unit and tanh activations, respectively at 72.44% and 72.43%. The lowest classification obtained was 54.92% using the deep net with L1 regularization and maxout activation. As can be seen from Fig. 12, a very significant improvement in classification accuracy was attained using the deep net with dropouts compared to the earlier work which did not make use of any dropouts and/or regularization, which was only between 61.15-67.68%. This is an improvement of over 10% and clearly shows the benefits of using dropouts to improve the generalization ability of deep nets.

TABLE I. EEG PREFERENCE CLASSIFICATION RESULTS

Deep Net Architecture	Hidden Layer Activation Function	Classification Accuracy (%)
Standard Deep Net	Tanh	67.68
Standard Deep Net	Maxout	61.15
Standard Deep Net	ReLU	63.99
Deep Net with Dropout	Tanh	74.38
Deep Net with Dropout	Maxout	67.71
Deep Net with Dropout	ReLU	79.76
Deep Net with L1 Regularization	Tanh	71.86
Deep Net with L1 Regularization	Maxout	54.92
Deep Net with L1 Regularization	ReLU	63.02
Deep Net with Dropout and L1 Regularization	Tanh	72.43
Deep Net with Dropout and L1 Regularization	Maxout	67.16
Deep Net with Dropout and L1 Regularization	ReLU	72.44

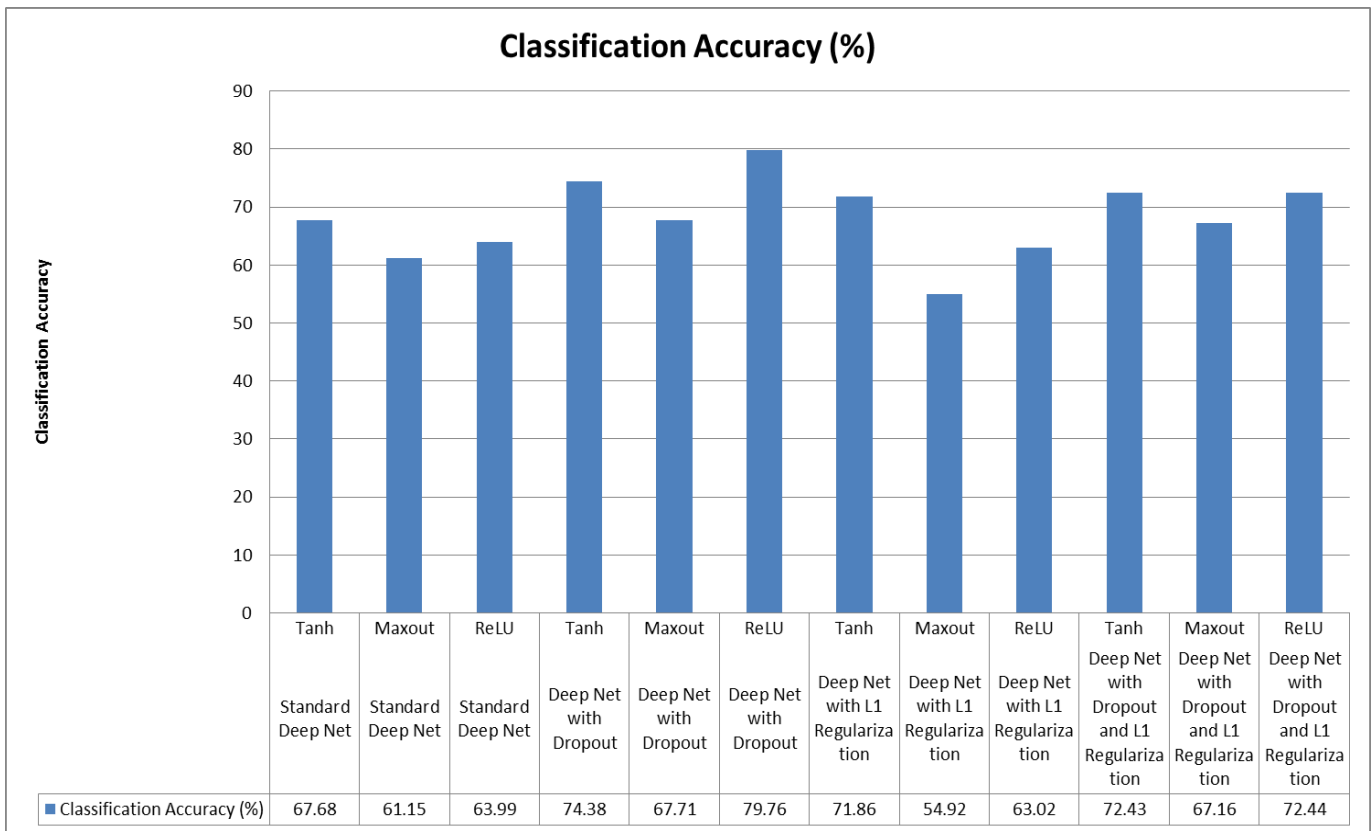


Fig. 12. Summary comparison of various deep learning architectures used.

B. Excitement Classification Results

Based on results shown in Table II, the SVM classifier achieved the overall best accuracy result of 89.36% using the Alpha band, while the second highest was 84.82%, achieved by the KNN classifier using the Theta band. Taking into account all datasets, the KNN classifier had the best performance as it held the most number of highest accuracy from 4 datasets (Delta, Theta, Beta, Gamma). However, the SVM classifier had a better average performance (SVM: 80.58%, KNN: 79.41%, RF: 78.92%, NN: 77.77%, C5.0: 75.54%).

Alpha band had shown the highest classification accuracy from two different classifiers (SVM and RF). This suggests that the Alpha band that represents the relaxed awareness of human contains some features that are useful to be used to classify the human emotion of “Excitement”. Moreover, the Theta band also showed a similar behaviour as the Alpha band. The Theta band tops the accuracy results on KNN and C5.0 classifiers and it represents the emotional stress, drowsiness and sleeps in adults.

In contrast as shown in Table III, the Gamma band had the worst overall results across all of the classifiers except NN classifier. This suggests that Gamma band that represents consciousness is not suitable to be used to classify human emotion of “Excitement” when the subject id immersed into the virtual stimuli.

For deep learning, preliminary testing yielded deep neural networks that performed best for this excitement classification

task using three hidden layers with 200 nodes each with weights initialized using the uniform adaptive method [29]. The deep neural networks were run using 10-fold cross-validation for 10 epochs each time using cross-entropy [29] as the error function and having a softmax output layer. Six different deep neural network architectures with different activation functions were tested, namely, tanh, maxout, and rectified linear (ReLU) [33], with and without dropout respectively, with dropout set at 0.5 and an adaptive learning rate method [34] applied when ReLU was used. The results obtained are tabulated below in Table IV.

TABLE II. SUMMARY OF TOP RESULTS OF 5 CLASSIFIERS

Classifier	Best Accuracy (%)	Band of Interest
SVM	89.36	Alpha
KNN	84.82	Theta
RF	81.96	Alpha
RF	81.96	All
NN	81.07	Beta
C5.0	80.89	Theta

TABLE III. SUMMARY OF BEST RESULTS OF ALL BANDS

Classifier	Accuracy of Band of Interest (%)					Combined
	Delta	Theta	Alpha	Beta	Gamma	
KNN	81.25	84.82	79.11	82.32	76.79	72.14
SVM	78.04	82.14	89.36	82.32	73.93	77.68
RF	79.62	79.29	81.96	78.21	72.50	81.96
NN	71.25	80.71	79.11	81.07	74.11	80.36
C5.0	72.50	80.89	80.36	75.36	65.36	78.75

TABLE IV. EXCITEMENT CLASSIFICATION RESULTS (%)

Activation	Delta	Theta	Alpha	Beta	Gamma	Combined
Tanh	89.77	93.71	82.57	88.21	86.21	92.07
Maxout	88.27	79.95	88.60	93.41	93.67	80.70
ReLU	85.44	88.18	88.67	87.82	86.44	91.01
Tanh with dropout	77.94	80.42	90.11	92.57	91.39	92.41
Maxout with dropout	83.68	83.12	82.82	88.96	88.71	91.42
ReLU with dropout	80.95	91.87	88.17	87.88	84.58	95.55

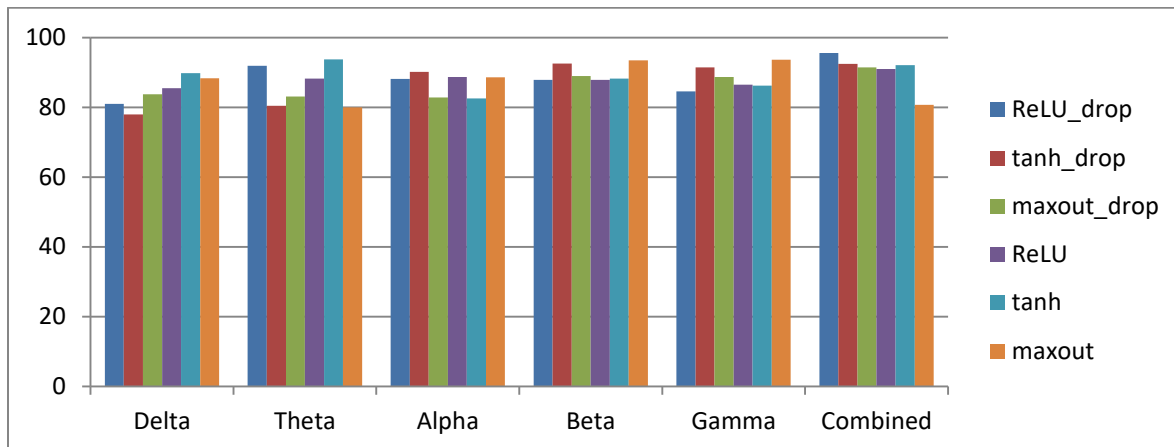


Fig. 13. Classification results grouped according to spectral band(s).

The best classification result of 95.55% was obtained using the ReLU with dropout deep neural network architecture using the combination of all of the available spectral bands. The next best result of 93.71% was obtained using the tanh with the theta band as the only input feature. The worst result of 77.94% was given by the tanh with dropout using the delta band. From the results, it appeared that there were no clear trends in terms of the architecture used but in terms of the spectral bands used, the combined approach appeared to provide an advantage whereby five out of six results yielded more than 90% accuracy as shown in Fig. 13. This suggests that, at least in terms of the excitement emotion, detecting this emotion benefits from looking at all spectral bands and not just at one or two specific bands such as alpha and beta which are commonly adopted for classifying EEG signals during active cognition.

V. CONCLUSION AND FUTURE WORK

Firstly, this study has comprehensively tested dropout and L1 regularization approaches to deep net architectures in an effort to improve the classification performance of deep learning neural networks in EEG-based preference classification. We have shown that using a deep net with dropouts using rectified linear units for activation was able to achieve a gain of more than 13%-18% at 79.76% accuracy compared to standard deep nets without such approaches at only between 61.15%-67.68% using various activations.

Secondly, this study has also investigated the use of deep learning for the detection of excitement while being immersed in virtual reality stimuli. To the best of our knowledge, this represents the first reported work that uses EEG solely as the input feature for the classification with the stimuli being virtual reality. It has been shown that a relatively high classification accuracy can be achieved with the best result yielding close to

96% accuracy. The results also suggest that using a combination of all EEG spectral bands as the input features provided more reliable classification results in general compared to using any other single EEG spectral alone.

For future work, due to the significant noise typically encountered in inter-subject EEG variations, we intend to investigate the use of autoencoders to pre-train the features extracted in order to further improve classification accuracy. Also, with the significant improvement in classification accuracy obtained through this study, we also plan to embark on application-based investigations into the use of EEG-based preference classification to guide automated generation of affective entertainment content in games, music and storytelling. It would be worthwhile to expand this line of work to include other emotions such as fear, boredom, frustration among others in view of expanding the potential applications of this EEG-based emotion classification in virtual reality approach particularly in the field of affective entertainment.

VI. ACKNOWLEDGEMENTS

This project is supported by the FRGS research grant scheme ref: FRG0349 & FRG0435 from the Ministry of Higher Education, Malaysia.

REFERENCES

- [1] L.H. Chew, J. Teo, and J. Mountstephens, "Aesthetic preference recognition of 3D shapes using EEG", *Cognitive Neurodynamics*, 10(2), pp.165-173, 2016.
- [2] J. Teo, L.H. Chew, and J. Mountstephens, "Deep learning for EEG-based preference classification", *International Conference on Applied Science and Technology (ICAST 2017)*, IEEE, April 2017.
- [3] H. Bellini, W. Chan, M. Sugiyama, M. Shin, S. Alam, and D. Takayama, excerpt from "Virtual & Augmented Reality: Understanding the Race for the Next Computing Platform", Profiles in Innovation, Goldman

- Sachs Equity Research, pp. 1-30, Feb. 2016, www.goldmansachs.com/our-thinking/pages/technology-driving-innovation-folder/virtual-and-augmented-reality/report.pdf, accessed 24 October 2017 (2016).
- [4] X.W. Wang, D. Nie, and B.L. Lu, "Emotional state classification from EEG data using machine learning approach", *Neurocomputing*, 129, pp.94-106, 2014.
- [5] A. Dhall, R. Goecke, J. Joshi, K. Sikka, and T. Gedeon, "Emotion recognition in the wild challenge 2014: Baseline, data and protocol", *Proceedings of the 16th International Conference on Multimodal Interaction*, pp.461-466, ACM, 2014.
- [6] G.K. Verma and U.S. Tiwary, "Multimodal fusion framework: A multiresolution approach for emotion classification and recognition from physiological signals", *NeuroImage*, 102, pp.162-172, 2014.
- [7] E.H. Jang, B.J. Park, S.H. Kim, M.A. Chung, M.S. Park, and J.H. Sohn, "Emotion classification based on bio-signals emotion recognition using machine learning algorithms", *2014 International Conference on Information Science, Electronics and Electrical Engineering (ISEEE)*, IEEE, vol. 3, pp.1373-1376, 2014.
- [8] S.K. Hadjidimitriou, A.I. Zacharakis, P.C. Doulgeris, K.J. Panoulas, L.J. Hadjileontiadis and S.M. Panas, "Revealing action representation processes in audio perception using fractal EEG analysis", *IEEE Transactions on Biomedical Engineering*, 58(4):1120-1129, 2011.
- [9] D.A. Adamos, S.I. Dimitriadis, and N.A. Laskaris, "Towards the bio-personalization of music recommendation systems: A single-sensor EEG biomarker of subjective music preference", *Information Sciences*, 343, pp.94-108, 2016.
- [10] M. Yadava, P. Kumar, R. Saini, P.P. Roy, and D.P. Dogra, "Analysis of EEG signals and its application to neuromarketing", *Multimedia Tools and Applications*, pp.1-25, 2017.
- [11] L.C. Chen, P. Sandmann, J.D. Thorne, C.S. Herrmann, and S. Debener, "Association of concurrent fNIRS and EEG signatures in response to auditory and visual stimuli", *Brain Topography*, 28(5), pp.710-725, 2015.
- [12] S. Goncalves, J. De Munck, P. Pouwels, R. Schoonhoven, J. Kuijter, N. Maurits, J. Hoogduin, E. Van Someren, R. Heethaar, F. L. Da Silva, "Correlating the alpha rhythm to bold using simultaneous EEG/fMRI: inter-subject variability", *Neuroimage* 30(1):203-213, 2006.
- [13] G. Pfurtscheller, C. Brunner, A. Schlogl, F. L. Da Silva, "Mu rhythm (de) synchronization and EEG single-trial classification of different motor imagery tasks", *Neuroimage* 31(1):153-159, 2006.
- [14] A. Yazdani, J. S. Lee, J.-M. Vesin, T. Ebrahimi, "A ECT recognition based on physiological changes during the watching of music video", *ACM Transactions on Interactive Intelligent Systems 2 (EPFL-ARTICLE-177741)*, pp. 1-26, 2012.
- [15] Y. Pan, C. Guan, J. Yu, K. K. Ang, T. E. Chan, "Common frequency pattern for music preference identification using frontal EEG", in: *Neural Engineering (NER)*, 2013 6th International IEEE/EMBS Conference on, IEEE, pp. 505-508, 2013.
- [16] K. C. Tseng, B.-S. Lin, C.-M. Han, P.-S. Wang, "Emotion recognition of EEG underlying favourite music by support vector machine", in: *Orange Technologies (ICOT)*, 2013 International Conference on, IEEE, pp. 155-158, 2013.
- [17] Y. Kim, K. Kang, H. Lee, C. Bae, "Preference measurement using user response electroencephalogram", in: *Computer Science and its Applications*, Springer, pp. 1315-1324, 2015.
- [18] S. K. Hadjidimitriou, L. J. Hadjileontiadis, "Toward an EEG-based recognition of music liking using time-frequency analysis", *IEEE Transactions on Biomedical Engineering* 59(12):3498-3510, 2012.
- [19] S. K. Hadjidimitriou, L. J. Hadjileontiadis, "EEG-based classification of music appraisal responses using time-frequency analysis and familiarity ratings", *IEEE Transactions on Affective Computing* 4(2):161-172, 2013.
- [20] J. Moon, Y. Kim, H. Lee, C. Bae, W. C. Yoon, "Extraction of user preference for video stimuli using EEG-based user responses", *ETRI Journal* 35(6):1105-1114, 2013.
- [21] K. Li, X. Li, Y. Zhang, A. Zhang, "Affective state recognition from EEG with deep belief networks", in: *Bioinformatics and Biomedicine (BIBM)*, 2013 IEEE International Conference on, IEEE, pp. 305-310, 2013.
- [22] W.-L. Zheng, J.-Y. Zhu, Y. Peng, B.-L. Lu, "EEG-based emotion classification using deep belief networks", in: *2014 IEEE International Conference on Multimedia and Expo (ICME)*, IEEE, pp. 1-6, 2014.
- [23] S. Jirayucharoensak, S. Pan-Ngum, P. Israsena, "EEG-based emotion recognition using deep learning network with principal component based covariate shift adaptation", *The Scientific World Journal*, 2014.
- [24] D. Wu, C. Courtney, B. Lance, S. Narayanan, M. Dawson, K. Oie, and T. Parsons, "Optimal Arousal Identification and Classification for Affective Computing Using Physiological Signals: Virtual Reality Stroop Task", *IEEE Transactions on Affective Computing*, 1(2): 109-118 (2010).
- [25] D. Massari, D. Pacheco, R. Malekshahi, A. Betella, P. Verschure, N. Birbaumer, and A. Caria, "Fast Mental States Decoding in Mixed Reality," *Frontiers in Behavioral Neuroscience*, 8:415, doi:10.3389/fnbeh.2014.00415 (2014).
- [26] N. Kovacevic, P. Ritter, W. Tays, S. Moreno, and A. McIntosh, "My Virtual Dream: Collective Neurofeedback in an Immersive Art Environment", *PLoS ONE*, 10(7): e0130129. doi.org/10.1371/journal.pone.0130129 (2015).
- [27] J. Gielis, "A generic geometric transformation that unifies a wide range of natural and abstract shapes", *American Journal of Botany*, 90:3, pp.333-338, 2003.
- [28] D. Nguyen, B. Widrow, "Improving the learning speed of 2-layer neural networks by choosing initial values of the adaptive weights", in: *Neural Networks, 1990. 1990 IJCNN International Joint Conference on, IEEE*, pp. 21-26, 1990.
- [29] P.T. De Boer, D.P. Kroese, S. Mannor, and R.Y. Rubinstein, "A tutorial on the cross-entropy method", *Annals of operations research*, 134(1):19-67, 2005.
- [30] P. Bashivan, I. Rish and S. Heisig, "Mental State Recognition via Wearable EEG," *arXiv preprint arXiv: 1602.00985* (2016).
- [31] D. Surangsrirat and A. Intarapanich, "Analysis of the meditation brainwave from consumer EEG device," *SoutheastCon 2015*, pp. 1-6 (2015).
- [32] M. Murugappan and S. Murugappan, "Human Emotion Recognition Through Short Time Electroencephalogram (EEG) Signals Using Fast Fourier Transform (FFT)," *IEEE 9th International Colloquium on Signal Processing and its Applications*, pp. 289-294 (2013).
- [33] S. Gunn, "Support Vector Machines for Classification and Regression," 1998.
- [34] S. Zainudin, S. D. Jasim and A. A. Bakar, "Comparative Analysis of Data Mining Techniques for Malaysian Rainfall Prediction," *International Journal on Advanced Science, Engineering and Information Technology*, vol. 6, no. 6, pp. 1148-1153, 2016.
- [35] C. Lehmann, "Application and comparison of classification algorithms for recognition of Alzheimer's disease in electrical brain activity (EEG)," *Journal of Neuroscience Methods* 161, pp. 342-350, 2007.
- [36] M. Soleymani, M. Pantic and P. Thierry, "Multi-Model Emotion Recognition in Response to Videos," *IEEE transactions on affective computing*, vol. 3, no. 2, pp. 211-223, 2011.
- [37] V. Nair, G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines", in: *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pp. 807-814, 2010.
- [38] M. D. Zeiler, "Adadelta: an adaptive learning rate method", *arXiv preprint arXiv:1212.5701*.
- [39] J. Teo, L.H. Chew, and J. Mountstephens, "Improving Subject-Independent EEG Preference Classification using Deep Learning Architectures with Dropouts", *Future of Information and Communication Conference (FICC 2018)*, IEEE, April 2018.

Rainfall Prediction using Data Mining Techniques: A Systematic Literature Review

Shabib Aftab, Munir Ahmad, Noureen Hameed, Muhammad Salman Bashir, Iftikhar Ali, Zahid Nawaz

Department of Computer Science
Virtual University of Pakistan
Lahore, Pakistan

Abstract—Rainfall prediction is one of the challenging tasks in weather forecasting. Accurate and timely rainfall prediction can be very helpful to take effective security measures in advance regarding: ongoing construction projects, transportation activities, agricultural tasks, flight operations and flood situation, etc. Data mining techniques can effectively predict the rainfall by extracting the hidden patterns among available features of past weather data. This research contributes by providing a critical analysis and review of latest data mining techniques, used for rainfall prediction. Published papers from year 2013 to 2017 from renowned online search libraries are considered for this research. This review will serve the researchers to analyze the latest work on rainfall prediction with the focus on data mining techniques and also will provide a baseline for future directions and comparisons.

Keywords—Rainfall prediction; data mining techniques; SLR; systematic literature review

I. INTRODUCTION

Analysis of time series data is one of the important aspects of modern research in the domain of knowledge discovery [28]. Time series data is collected over a specific period of time such as hourly, daily, weekly, monthly, quarterly or yearly [23], [40]. Data mining techniques can use this data to predict upcoming situations in various domains such as climate change, education, and finance etc. These techniques can be used to extract hidden knowledge from time series data for future use [23], [27], [29], [40]. Weather forecasting is very beneficial but challenging task [26]. Weather data consists of various atmospheric features such as wind speed, humidity, pressure and temperature etc. Data mining techniques have the capacity to extract the hidden patterns among available features of past weather data and then these techniques can predict future weather conditions by using extracted patterns [40]. Rainfall is a complex atmospheric process, which depends upon many weather related features. Accurate and timely rainfall prediction can be helpful in many ways such as planning the water resources management, issuance of early flood warnings, managing the flight operations and limiting the transport & construction activities [24], [25]. Accurate rainfall prediction is more complex today due to climate variations. Researchers consistently have been working to predict rainfall with maximum accuracy by optimizing and integrating data mining techniques [41]. Data mining algorithms are classified as supervised and unsupervised. Supervised methods get trained first with pre-classified data (training data) and then classify the input data

(test data) [7], [38], [39]. Un-supervised methods on the other hand do not require any training, instead of pre-classified data these techniques use algorithms to extract hidden structure from un-labeled data. It has been observed from latest research that for high accuracy, researchers prefer the integrated techniques for the rainfall prediction. To reflect the latest research, this study provides a systematic literature review by focusing on latest papers, which are published in last five years (2013-2017). Three renowned online search libraries are selected for literature extraction: Elsevier, IEEE and Springer. Initially 4844 papers are extracted and then through a systematic research process 8 most relevant research articles are selected for critical review.

Further organization of this paper is as follows. Section II elaborates the related work. Section III presents the research protocol, which is followed in this research. Section IV presents the review of shortlisted articles. Section V discusses the review findings. Section VI finally concludes this study.

II. RELATED WORK

Researchers have been working to improve the accuracy of rainfall prediction by optimizing and integrating data mining techniques. Some of the selected studies are discussed in this section. In [1], author performed a comparative analysis of Support Vector Machine (SVM), Artificial Neural Networks (ANN), and Adaptive Neuro Fuzzy Inference System (ANFIS) on rainfall prediction. The authors have compared the prediction models in four terms: (i) by using different lags as modeling inputs; (ii) by using training data of heavy rainfall events only; (iii) performance of forecasting for 1 hour to 6 hours and; (iv) performance analysis in peak values and all values. According to results ANN performed better when trained with dataset of heavy rainfall. For 1 to 4 hour ahead forecasting, the previous 2-hour input data was suggested for all three modeling techniques (ANN, SVM and ANFIS). ANFIS reflected better ability in avoiding information noise by using different lags of inputs. And finally during peak values, SVM proved to be more robust under extreme typhoon events. Researchers in [2] performed a comparative analysis of various data mining techniques for rainfall prediction in Malaysia such as: Random Forest, Support Vector Machine, Naive Bayes, Neural Network, and Decision Tree. For this experiment, dataset was obtained from various weather stations in Selangor, Malaysia. Before classification process, Pre-processing tasks were applied to deal with the noise and missing values in dataset. The results showed significant

performance of Random Forest as it correctly classified large amount of instances with small amount of training data. In [3], author performed a survey on various Neural Network architectures which were used for rainfall prediction in last 25 years. The authors highlighted that most of the researchers got significant results in rainfall prediction by using Propagation Network, moreover the forecasting techniques which used SVM, MLP, BPN, RBFN, and SOM are more suitable than other statistical and numerical techniques. Some limitations have also been highlighted. Researchers in [4] used Artificial Neural Network for rainfall prediction in Thailand. They used Back Propagation Neural Network for prediction which reported an acceptable accuracy. For future direction it was suggested that few additional features would be included in input data for rainfall prediction such as Sea Surface Temperature for the areas around Andhra Pradesh and Southern part of India. Researchers in [5] predicted monthly rainfall by using Back Propagation, Radial Basis Function and Neural Network. For prediction, the dataset was collected from Coonoor region in Nilgiri district (Tamil Nadu). Performance was evaluated in terms of Mean Square Error. According to results higher accuracy was reported in Radial Basis Function Neural Network with smaller Mean Square Error. Moreover the researchers also used these techniques for future rainfall prediction. Researchers in [6] presented a Hybrid Intelligent System by integrating Artificial Neural Network and Genetic Algorithm. In ANN, MLP works as the Data Mining engine to perform predictions whereas the Genetic Algorithm was utilized for inputs, the connection structure between the inputs, the output layers and to make the training of Neural Network more effective. Researchers in [8] discussed rainfall pace in previous years with respect to various crops seasons like rabi, Kharif, zaid and then predicted (rainfall) for future seasons via Linear Regression Method. For prediction, input dataset was selected according to particular crops seasons of previous years. In [9], one month and two month forecasting models were developed for rainfall prediction by using Artificial Neural Network (ANN). The input dataset was selected from multiple stations in North India, spanned on past 141 years. Feed Forward Neural Network using Back Propagation and Levenberg-Marquardt training function were used in these models. Performance of both models was evaluated by using Regression Analysis, Mean Square Error and Magnitude of Relative Error. The results showed that one month forecasting model can predict the rainfall more accurately than two month forecasting model. Researchers in [10] presented an algorithm by integrating Data Mining and Statistical Techniques. The proposed technique predicted the rainfall in five different categories such as: Flood, Excess, Normal, Deficit and Drought. The predictors were selected with highest confidence level, based on association rules and derived from local and global environment. From local environment: wind speed, sea level pressure, maximum temperature, and minimum temperature were taken. From global environment: Indian ocean dipole conditions and southern oscillation were taken.

In [11], researchers predicted the rainfall by using proposed Wavelet Neural Network Model (WNN), an integration of Wavelet Technique and Artificial Neural Network (ANN). To analyze the performance, monthly rainfall prediction was performed with both the techniques (WNN and ANN) by using dataset of Darjeeling rain gauge station in India. Statistical techniques were used for performance evaluation and according to results WNN performed better than ANN. In [12], researchers provided a detailed survey and performed a comparative analysis of various neural networks on rainfall forecasting. According to survey RNN, FFNN, and TDNN are suitable for rainfall prediction as compared to other statistical and numerical forecasting methods. Moreover TDNN, FFNN and lag FFNN performed well for yearly, monthly and weekly rainfall forecasting respectively. This research also discussed the various measures of accuracy used by different researchers to evaluate the ANN's performance.

III. RESEARCH PROTOCOL

High quality SLR is one which attains its objective by providing the compact information of required research topic for a particular time span. A detailed research methodology with step by step guidance is needed to conduct an effective SLR. In this research a systematic research process is formulated by following the guidelines extracted from [13]-[18]. Usually SLR consists of three basic steps: plan review, conduct review and document review moreover further nested steps can be included from modern and state of the art research papers for an effective presentation. For this study, a step by step systematic review process is extracted from the latest review articles of software engineering domains [19]-[22]. The systematic review process of this research consists of the following steps: A) Identification of research questions, B) Keywords selection for query string, C) Selection of search space, D) Outlining the selection criteria, E) Literature extraction, F) Quality assessment, G) Literature Analysis and H) Results and Discussion (Fig. 1).

A. Identification of Research Questions

Research objectives are identified and presented in the form of research questions. The ultimate purpose of SLR is to find the answers of those questions via critical review. Flowing are the research questions identified for this research.

RQ1: Which data mining techniques are used / proposed for rainfall prediction?

RQ2: How the performance of prediction techniques is evaluated?

RQ3: Which type of data is used for prediction?

RQ4: For which location the rainfall prediction is performed?

RQ5: Which factors affect the prediction results?

RQ6: Which are the latest research trends in the domain of rainfall prediction?

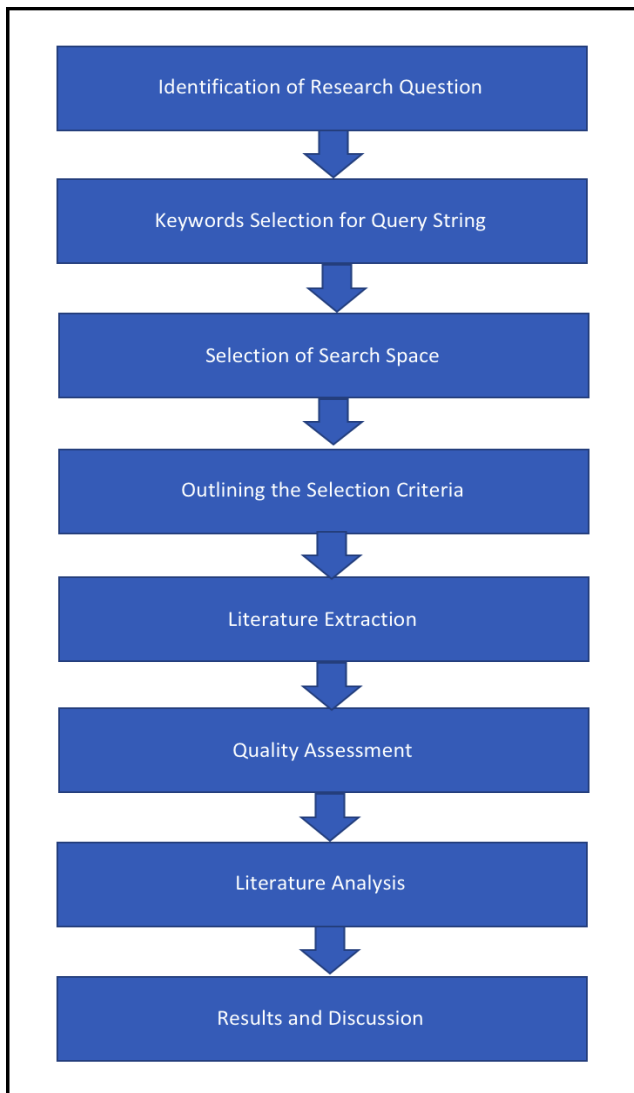


Fig. 1. SLR process.

B. Keywords Selection for Query String

Second step is to formulate the query string and for that purpose, keywords are extracted first from the research questions and then arranged in a particular sequence to form a query. Following keywords are extracted for query:

Improved, Customized, Integrated, Data Mining, Techniques, Methods, Algorithms, Rainfall, Prediction, Forecasting, Estimation, Performance, Evaluation, Assessment.

The finalized query string is given below:

("Performance" AND ("Evaluation" OR "Assessment") AND / OR ("Improved" OR "Customized" OR "Integrated") AND ("Data Mining") AND ("Techniques" OR "Methods" OR "Algorithms") AND "Rainfall" AND ("Prediction" OR "Forecasting" OR "Estimation")).

C. Selection of Search Space

This step deals with the selection of libraries from where the related literature will be extracted through query string.

Three well known and widely used online libraries are selected to extract the literature: IEEE, Elsevier and Springer. All three libraries have different options to search the relevant material, so few adjustments were made in query strings to extract the appropriate and most relevant literature. The Query was searched multiple times with various combinations of key-words. Results of search queries are available in Table I.

TABLE I. SEARCH SPACE AND QUERY RESULTS

Sr. #	Digital Library	Date Searched	Total Results
1	Elsevier	2018-24-02	1819
2	IEEE	2018-24-02	1119
3	Springer	2018-24-02	1906

D. Outlining the Selection Criteria

This step aims to outline the selection boundary so that most relevant research papers can be selected. This activity consists of two steps, IC (inclusion criteria) and EC (exclusion criteria).

1) Inclusion Criteria

Below are the rules of Inclusion criteria.

IC1: Papers which are published from 2013 till 2017.

IC2: Papers which are available in journals, conferences, proceedings of conferences or workshops.

IC3: Papers which have predicted the rainfall using data mining techniques.

IC4: Papers which have performed comparison of data mining techniques on rainfall prediction.

IC5: Papers which have presented improved or customized data mining techniques to predict rainfall.

IC6: Papers which have integrated data mining technique with any other technique.

2) Exclusion Criteria (EC)

Below are the rules of exclusion criteria.

EC1: Papers which are not in English.

EC2: Papers published before 2013 or after 2017.

EC3: Papers which did not perform rainfall prediction.

EC4: Paper which did not use any data mining technique in proposed model/method?

EC5: Paper which did not use any weather data for prediction.

EC6: Papers which did not evaluate the performance of used/proposed technique.

E. Literature Extraction

The purpose of selection criteria is to extract the most relevant literature for the review. After applying IC and EC, 18 articles were shortlisted. Complete process of literature extraction is given in Fig. 2.

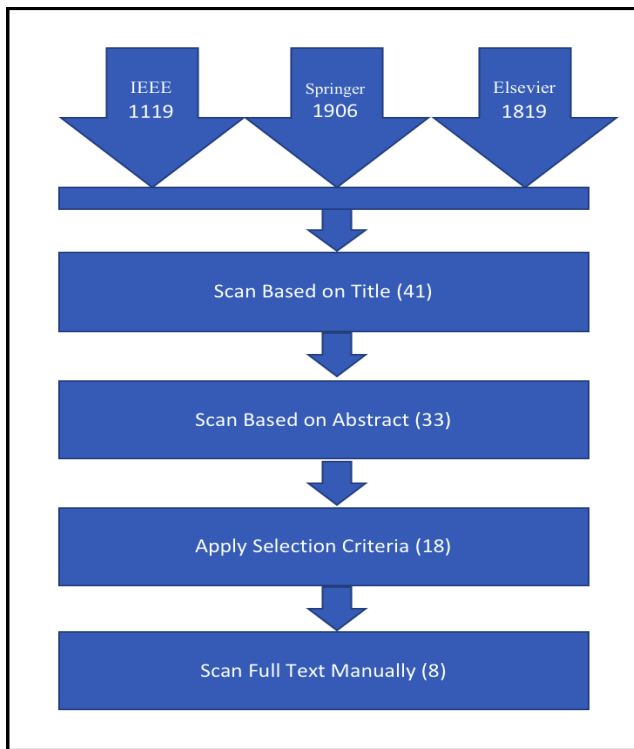


Fig. 2. Search process.

F. Quality Assessment

To meet the research objectives, it was make sure to follow the quality parameters throughout the systematic research process. To ensure the quality of results, following measures were taken.

- Authentic and renowned online libraries were selected to extract research articles.
- Latest research papers were selected to reflect latest research.
- The process of selection was un-biased.
- Complete steps of Systematic Research Process were followed in the true sense.

IV. LITERATURE ANALYSIS

Full text of 18 selected articles were analyzed and then 8 most relevant research papers are shortlisted for critical review as shown in Table II. The Review of shortlisted articles is given below.

TABLE II. MOST RELEVANT RESEARCH LITERATURE

Sr. #	Digital Library	Selected Research Literature	No. of Researches
1	Elsevier	[30]-[33]	4
2	IEEE	[34]	1
3	Springer	[35]-[37]	3

A. Indian Summer Monsoon Rainfall (ISMR) Forecasting using Time Series Data: A Fuzzy-Entropy-Neuro based Expert System

In [30], authors presented a model to forecast Indian Summer Monsoon Rainfall on the basis of monthly and seasonal timescales. To forecast, time series dataset was used, spanning from 1871 till 2014. The dataset was classified in two parts (1) 1871-1960 used as training data, and (2) 1961-2014 used as test data. Statistical analysis reported the dynamic nature of rainfall in monsoon, which could not be predicted effectively with mathematical and statistical models. So, the authors in this research recommended to use three techniques for this type of prediction: Fuzzy Set, Entropy and Artificial Neural Network. By using these three techniques, a forecasting model is developed to deal with the dynamic nature of the ISMR. In proposed model, fuzzy set theory is used to handle uncertainties which are inherited in dataset. The entropy computational concept was modified in this model and used to provide the input as a degree of membership in the entropy function. That entropy function was referred as Fuzzy Information-Gain (FIG). Then, each fuzzified rule was defuzzified using the ANN. The value of FIG of each fuzzy- set was then used as input into ANN. The proposed model was named as “Fuzzy-Entropy-Neuro Based Expert System for ISMR Forecasting” because it is the integration of fuzzy set, entropy and ANN. To evaluate the performance of proposed model following accuracy measures were used: Standard Deviations (SDs), Correlation Coefficient (CC), Root Mean Square Error (RMSE) and Performance Parameter (PP). According to results the proposed model is effective and efficient in comparison with other existing models.

B. An Extensive Evaluation of Seven Machine Learning Methods for Rainfall Prediction in Weather Derivatives

The researchers in [31] compared the predictive performance of latest and state of the art method named “Markov chain extended with rainfall prediction” with the other widely used machine learning techniques: Support Vector Regression, Genetic Programming, M5 Rules, M5 Model trees, Radial Basis Neural Networks, and k-Nearest Neighbours. Daily rainfall datasets were collected from 42 cities of two continents, with very diverse climatic features. 20 cities were selected from around the Europe and 22 from around the USA. There were two reasons of choosing two continents for data extraction, first is to perform the experiment on different climates having diverse weather and second was the geographical locations as the selected cities were very far apart from each other. The ultimate goal was to not bias the experiment to particular climate type or for particular geographic location. According to results the accumulating rainfall amounts can bring good results as compared to prediction using daily rainy data. While using the accumulated data, Support Vector Regression, Radial Basis Functions, and Genetic Programming overall performed well however Radial Basis Functions performed better then modern technique of “Markov chain”. For all selected datasets, each technique used the same parameters so it was not guaranteed that the best possible set of parameters was used for all the techniques. During the experiment, the researchers have noted

a relationship between predictive accuracy and climatic attributes such as: volatile nature of rainfall, amount of maximum rainfall and the interquartile range of rainfall. Moreover no significant difference was noted in algorithms' prediction error among the cities of both the continents (USA and Europe). Issue regarding the discontinuity in rainfall data was solved with the help of accumulated rainfall amounts.

C. A Hybrid Model for Statistical Downscaling of Daily Rainfall

Authors in [32] proposed a hybrid technique to downscale the daily rainfall by integrating two methods: 1) Random Forest, and 2) Support Vector Machine. RF was selected due to its robustness in classification and it was used to predict whether it will be rain or not whereas SVM were selected due to its feature to fit in non-linear data and it was used to predict the amount of rainfall, if it will occur. The proposed model was evaluated by downscaling daily rainfall at three stations, Dungun, Besut, and Kemaman on the east coast of peninsular, Malaysia. Daily rainfall time series data from 1961 till 2000 was collected from Department of Irrigation and Drainage Malaysia. Total of 26 climatic features were collected from National Centre for Environmental Prediction re analysis dataset, which were used as predictors for downscaling the models. To assess homogeneity in rainfall time series, various quality control activities were performed. Histograms for the dataset were created to reflect the problems moreover Student's t test was also used to identify any variance in the means between two segments of dataset which finally found homogeneous at all three locations. According to results the hybrid technique is capable to downscale the rainfall with Nash-Sutcliffe efficiency within range of 0.90-0.93, which is much higher than RF and SVM models.

D. Prediction of Monthly Rainfall in Victoria, Australia: Clusterwise Linear Regression Approach

In [33], researchers presented a technique named Clusterwise Linear Regression for monthly rainfall prediction in Victoria, Australia. The proposed CLR is an integrated method of clustering and regression techniques. CLR incrementally extracted the subsets from dataset and then those subsets could be easily estimated with linear function one by one. Dataset which was used for prediction obtained from eight different weather stations for the period of 1889 - 2014 and consisted of five meteorological variables. The selected weather stations included three from east region, two from central region and three from the west region of Victoria. The ultimate goal for the selection of geographical apart stations was to evaluate the performance of proposed model on multiple locations having different atmospheres. The meteorological variables which were used as predictors included Vapor Pressure, Solar Radiation, Evaporation, Minimum Temperature, and Maximum Temperature. This proposed technique was compared with following: SVM Reg, ANNs, CLR with CR-EM, and MLR. The model was developed first for each weather station with each technique using training data and then evaluated with test data. To analyze the performance of proposed technique, observed and predicted rainfall measures were compared and four accuracy parameters were used for evaluation: Mean Absolute Scaled Error, Mean Absolute Error, Root Mean Squared Error, and

coefficient of efficiency. According to results, the proposed technique outperformed other prediction methods in most of the locations.

E. Prediction and Anomaly Detection of Rainfall using Evolving Neural Network to Support Planting Calendar in Soreang (Bandung)

Authors in [34] proposed Evolving Neural Network for the prediction and anomaly detection of rainfall to Support Planting Calendar in Soreang. Dataset was obtained from Department of Agriculture and Department of Water Resources spanning from 1999-2013. The proposed ENN used Artificial Neural Networks and Genetic Algorithm to identify the best weights and biases. The proposed framework consisted of various steps starting from the obtaining of raw data which then gone through the pre-processing phase which consisted of following steps: Integration, Transformation, Reduction and Cleaning of data. Dataset was divided in three scenarios: scenario 1 as dry season from April to September, scenario 2 as wet season from October to March and scenario 3 as the complete data from January to December. Each scenario was further sub divided for training and test data as 9, 12, 14 years for training data and 6, 3, 1 years for testing data, respectively. Learning process of proposed framework used integrated technique and then the result was used for rainfall prediction and anomaly detection followed by the final result which was the predicted starting time for planting. The starting week of January, April and October was selected as beginning time for planting activity in year 2014. According to results, by using all data from 1999-2013 shown the accuracy of 84.6%, for dry season the reported accuracy was 66.02% and for wet season the accuracy was 79.7%.

F. Rainfall Prediction: A Deep Learning Approach

In [35], authors presented a Deep Learning based architecture to predict the daily accumulated rainfall for next day. Proposed architecture consists of two techniques: Auto encoder Network and the Multilayer Perceptron Network. Auto encoder is an unsupervised network which performed the feature selection activity and the Multilayer Perceptron Network was assigned the classification and prediction tasks. Dataset for prediction was obtained from Instituto de Estudios Ambientales (IDEA) of Universidad Nacional de Colombia which is located in Manizales Colombia. Dataset spanned from 2002 to 2013 and consisted of 47 weather attributes. IDEA extracted the data from a meteorological station located in the central area of same city and stored in an environmental DWH. As ETL steps were performed on data so pre processing was not needed. Obtained 2952 data samples were classified into subsets for the purpose of training, validation and testing, with 70%, 15% and 15%, respectively. Normalization process was then performed to keep the values of data in to the range of 0 to 1 for better working. Results of the experiment were compared with other methods such as: naive approach which predicts the accumulated rainfall of $t - 1$ for t , MLP with optimized parameters for training & validation set and with some other published techniques. Performance was evaluated in terms of measurement errors: Mean Square Error and Root Mean Square Error.

G. A novel approach for Optimizing Climate Features and Network Parameters in Rainfall Forecasting

Authors in [36] presented a Genetic Algorithm-based approach to identify the best combination of input features and Neural Network parameters to achieve most accurate result. Dataset for prediction spanning of 107 years, from 1908 to 2015, taken from Innisfail, Queensland, Australia and consisted of various weather attributes including rainfall values, mean maximum temperature, mean minimum temperature, and Southern Oscillation Index etc. Data went through a preprocessing stage where couple of tasks was performed. In preprocessing, missing values were replaced with the mean of that attribute and when not applicable the value of that record was taken from closely available weather station. Genetic algorithm usually picks the best chromosome from last iteration but in proposed approach it is customized to select the best chromosome in each of the iteration. The best network which was saved in current iteration was compared to the other generated networks in each coming iteration. The proposed model reflected the highest scores, when compared to climatology and alternative selection methods. Selection of Climatic attributes and network parameters by using proposed hybrid genetic algorithm reflected better performance with 141.67 mm RMSE for a location with 3553.0 mm annual average rainfall whereas climatology, climate input parameters selection-based genetic algorithm, and climate features selection-based genetic algorithm showed 200.32, 171.34, and 178.22 mm consecutively.

H. Early Prediction of Extreme Rainfall Events: A Deep Learning Approach

Authors in [37] presented a framework for the prediction of extreme rainfall by using past climatic features. The proposed model consisted of following phases: Feature Learning, Feature Compression, and the classification process. Stacked Auto-encoder was used for the compression of feature-set. Support Vector Machines and Neural Network were used for classification. Parameters of selected classifier were tuned for the best performance and the issue of biased dataset was dealt effectively by Cost-Sensitive SVM. Presented technique showed the ability to predict extreme rainfall before 6 to 48 hours from occurrence; however some false positives were also reported. The proposed technique also reduced the false alarms which were raised due to the rainfall in surroundings. This method had the capability to generate warnings for rain in surroundings as well. Dataset for rainfall prediction was collected from National Centers for Environmental Prediction/National Center for Atmospheric Research (NCEP/NCAR), for the following months: June, July, August and September. The obtained dataset spanned from 1969 to 2008 for Mumbai, and from 1980 to 2000 for Kolkata. Rainfall data was also obtained for the same period from India Meteorological Department. Weather variables for prediction were taken for entire Indian sub-continent region which was divided in to 255 grids. Total of 21 variables were obtained for each grid; 4725 for entire region (255*21) in case of daily data which could further increased in case of 24 h and 48 h data. The results of experiment were compared with other methods from literature and found the proposed one much better.

V. RESULTS AND DISCUSSIONS

Eight research papers are finally shortlisted by applying the literature extraction criteria, explained in Section III. Below are the answers of Research Questions which are extracted during in-depth analysis and review of shortlisted papers.

RQ1: Which data mining techniques are used/proposed for rainfall prediction?

Authors in all selected papers presented customized/integrated/modified mining techniques for effective rainfall prediction. In each research, multiple climatic attributes/variables from past weather data were used as predictors for the purpose of prediction/forecasting. The ultimate purpose of each research was to increase the accuracy of rainfall prediction. Detail review of selected papers is available in previous section.

RQ2: How the performance of prediction techniques is evaluated?

The selected papers [30]–[37] have compared the proposed technique/model with one or more published techniques. The performance was evaluated by comparing the predicted results with the observed (actual) measures. Information retrieval metrics and statistical techniques were used for performance analysis of proposed techniques in comparison with other methods from the literature.

RQ3: Which type of data is used for prediction?

Each of the selected paper used past weather data for rainfall prediction and for the training purpose of used supervised data mining techniques. Un-supervised data mining techniques were also used in combination of supervised techniques. Various climatic attributes were used as predictors including rainfall polarity, rainfall measure, minimum temperature, maximum temperature, wind speed, and humidity etc. According to researchers, using more features is not the guarantee for more accuracy in prediction instead irrelevant attributes could affect the performance. So the combination of relevant attributes is needed for accurate rainfall prediction moreover these combinations varies upon case to case.

RQ4: For which location the rainfall prediction is performed?

According to shortlisted articles, rainfall was predicted in locations situated in India, Australia, Columbia, Indonesia, Malaysia. USA and Europe.

RQ5: Which factors affect the prediction results?

After the critical review of shortlisted papers, it has been observed that following factors could affect the rainfall prediction results: Past weather data: which is selected for training the mining algorithm, climatic attributes: which are used as predictors, location: for which the rainfall prediction has to be performed, surrounding environment, pre-processing techniques and most importantly the used model/technique/method.

RQ6: Which are the latest research trends in the domain of rainfall prediction?

The ultimate goal of all the shortlisted articles was to improve the prediction accuracy, for this purpose some researchers have explored the correlation among weather features and prediction accuracy and tried to find the best combinations of those features to tune the performance. Few researchers on the other hand worked to train the mining technique well to achieve the high accuracy in prediction. Few have compared the modern techniques with the conventional ones. However most of the researchers presented/used integrated techniques and focused on using the combination of two or more techniques for prediction and claimed that this could bring more accurate results. Each research has provided the justification for the presented/proposed/used technique by means of performance evaluation through quality metrics.

A. Limitations of Research:

This research has the following limitations.

1) *The literature was extracted with a rigorous and thorough research process which indicates the quality and completeness of this study however some important relevant work might have been missed.*

2) *Most of the integrated and modified techniques were evaluated by authors, so the real results may not be as accurate as explained. This may affect the analysis and results of this study.*

VI. CONCLUSION AND FUTURE WORK

Rainfall prediction is a beneficial but challenging task. Data mining techniques have the ability to predict the rainfall by extracting and using the hidden knowledge from past weather data. In the last decade, many researchers have worked to increase the accuracy of rainfall prediction by optimizing and integrating data mining techniques. Various models and techniques are available today for effective rainfall prediction but still there was a lack of a compact literature review and systematic mapping study which could reflect the current problems, proposed solutions and the latest trends in this domain. This research provided a comprehensive systematic mapping as well as the critical review of latest research from 2013 till 2017 in the area of rainfall prediction by focusing on data mining techniques. In this research a list of significant research questions was identified and then a systematic research process was followed to extract and shortlist the most relevant research articles from renowned digital search libraries. Answers of the identified questions were explored by critically reviewing the shortlisted articles. The research focus on the domain of rainfall prediction has been increasing since last decade and so are the problem areas. So it was concluded that enhancements, optimizations and integrations of data mining methods are vital to explore and solve these problems.

REFERENCES

- [1] S. Zhang, L. Lu, J. Yu, and H. Zhou, "Short-term water level prediction using different artificial intelligent models," in 2016 5th International Conference on Agro-Geoinformatics, Agro-Geoinformatics 2016, 2016.
- [2] S. Zainudin, D. S. Jasim, and A. A. Bakar, "Comparative Analysis of Data Mining Techniques for Malaysian Rainfall Prediction," *Int. J. Adv. Sci. Eng. Inf. Technol.*, vol. 6, no. 6, pp. 1148–1153, 2016.
- [3] D. Nayak, A. Mahapatra, and P. Mishra, "A Survey on Rainfall Prediction using Artificial Neural Network," *Int. J. Comput.*, vol. 72, no. 16, pp. 32–40, 2013.
- [4] B. K. Rani and A. Govardhan, "RAINFALL PREDICTION USING DATA MINING TECHNIQUES - A SURVEY," pp. 23–30, 2013.
- [5] N. Tyagi and A. Kumar, "Comparative analysis of backpropagation and RBF neural network on monthly rainfall prediction," *Proc. Int. Conf. Inven. Comput. Technol. ICICT 2016*, vol. 1, 2017
- [6] N. Solanki and G. P. B., "A Novel Machine Learning Based Approach for Rainfall Prediction," *Inf. Commun. Technol. Intell. Syst. (ICTIS 2017) - Vol. 1*, vol. 83, no. Ictis 2017, 2018.
- [7] M. Ahmad, S. Aftab, and I. Ali, "Sentiment Analysis of Tweets using SVM," *Int. J. Comput. Appl.*, vol. 177, no. 5, pp. 25–29, 2017
- [8] C. S. Thirumalai, "Heuristic Prediction of Rainfall Using Machine Learning Techniques," no. May, 2017.
- [9] N. Mishra, H. K. Soni, S. Sharma, and A. K. Upadhyay, "Development and Analysis of Artificial Neural Network Models for Rainfall Prediction by Using Time-Series Data," *Int. J. Intell. Syst. Appl.*, vol. 10, no. 1, pp. 16–23, 2018.
- [10] H. Vathsala and S. G. Koolagudi, "Prediction model for peninsular Indian summer monsoon rainfall using data mining and statistical approaches," *Comput. Geosci.*, vol. 98, pp. 55–63, 2017.
- [11] R. Venkata Ramana, B. Krishna, S. R. Kumar, and N. G. Pandey, "Monthly Rainfall Prediction Using Wavelet Neural Network Analysis," *Water Resour. Manag.*, vol. 27, no. 10, pp. 3697–3711, 2013.
- [12] M. P. Darji, V. K. Dabhi, and H. B. Prajapati, "Rainfall forecasting using neural network: A survey," 2015 *Int. Conf. Adv. Comput. Eng. Appl.*, no. March, pp. 706–713, 2015.
- [13] P. Brereton, B. A. Kitchenham, D. Budgen, M. Turner, and M. Khalil, "Lessons from applying the systematic literature review process within the software engineering domain," *J. Syst. Softw.*, vol. 80, no. 4, pp. 571–583, 2007.
- [14] B. a. Kitchenham et al., "Preliminary guidelines for empirical research in software engineering," *IEEE Trans. Softw. Eng.*, vol. 28, no. 8, pp. 721–734, 2002.
- [15] B. Kitchenham and S. Charters, "Guidelines for performing Systematic Literature reviews in Software Engineering Version 2.3," *Engineering*, vol. 45, no. 4ve, p. 1051, 2007.
- [16] K. Petersen, R. Feldt, S. Mujtaba, and M. Mattsson, "Systematic Mapping Studies in Software Engineering," 12th *Int. Conf. Eval. Assess. Softw. Eng.*, pp. 1–10, 2008.
- [17] B. Kitchenham, O. P. Brereton, D. Budgen, M. Turner, J. Bailey, and S. Linkman, "Systematic literature reviews in software engineering – A systematic literature review," *Inf. Softw. Technol.*, vol. 51, pp. 7–15, 2008.
- [18] M. Ahmad, S. Aftab, M. S. Bashir, and N. Hameed, "Sentiment Analysis using SVM: A Systematic Literature Review," *Int. J. Adv. Comput. Sci. Appl.*, vol. 9, no. 2, pp. 182–188, 2018.
- [19] F. Selleri Silva et al., "Using CMMI together with agile software development: A systematic review," *Inf. Softw. Technol.*, vol. 58, pp. 20–43, 2015.
- [20] F. Anwer and S. Aftab, "Latest Customizations of XP: A Systematic Literature Review," *Int. J. Mod. Educ. Comput. Sci.*, vol. 9, no. 12, pp. 26–37, 2017.
- [21] S. Ashraf and S. Aftab, "Scrum with the Spices of Agile Family: A Systematic Mapping," *I.J. Mod. Educ. Comput. Sci.*, vol. 9, no. 11, pp. 58–72, 2017.
- [22] S. Ashraf and S. Aftab, "Latest Transformations in Scrum: A State of the Art Review," *Int. J. Mod. Educ. Comput. Sci.*, vol. 9, no. 7, pp. 12–22, 2017.
- [23] N. Mishra, H. K. Soni, S. Sharma, and A. K. Upadhyay, "A Comprehensive Survey of Data Mining Techniques on Time Series Data for Rainfall Prediction," *J. ICT Res. Appl.*, vol. 11, no. 2, p. 168, 2017.

- [24] K. W. Chau and C. L. Wu, "A hybrid model coupled with singular spectrum analysis for daily rainfall prediction," *J. Hydroinformatics*, vol. 12, no. 4, p. 458, 2010.
- [25] J. Wu, J. Long, and M. Liu, "Evolving RBF neural networks for rainfall prediction using hybrid particle swarm optimization and genetic algorithm," *Neurocomputing*, vol. 148, pp. 136–142, 2015.
- [26] W. C.L. and K.-W. Chau, "Prediction of Rainfall Time Series Using Modular Soft Computing Methods," *Eng. Appl. Artif. Intell.*, vol. 26, no. 852, pp. 1–37, 2012.
- [27] D. Gupta and U. Ghose, "A Comparative Study of Classification Algorithms for Forecasting Rainfall," pp. 0–5, 2015.
- [28] M. A. Nayak and S. Ghosh, "Prediction of extreme rainfall event using weather pattern recognition and support vector machine classifier," *Theor. Appl. Climatol.*, vol. 114, no. 3–4, pp. 583–603, 2013.
- [29] M. Ahmad, S. Aftab, and S. S. Muhammad, "Machine Learning Techniques for Sentiment Analysis: A Review," *Int. J. Multidiscip. Sci. Eng.*, vol. 8, no. 3, pp. 27–32, 2017.
- [30] P. Singh, "Indian summer monsoon rainfall (ISMR) forecasting using time series data: A fuzzy-entropy-neuro based expert system," *Geosci. Front.*, vol. 2002, 2017.
- [31] S. Cramer, M. Kampouridis, A. A. Freitas, and A. K. Alexandridis, "An extensive evaluation of seven machine learning methods for rainfall prediction in weather derivatives," *Expert Syst. Appl.*, vol. 85, pp. 169–181, 2017.
- [32] S. H. Pour, S. Shahid, and E. S. Chung, "A Hybrid Model for Statistical Downscaling of Daily Rainfall," *Procedia Eng.*, vol. 154, pp. 1424–1430, 2016.
- [33] A. M. Bagirov, A. Mahmood, and A. Barton, "Prediction of monthly rainfall in Victoria, Australia: Clusterwise linear regression approach," *Atmos. Res.*, vol. 188, pp. 20–29, 2017.
- [34] Gunawansyah, T. H. Liong, and Adiwijaya, "Prediction and anomaly detection of rainfall using evolving neural network to support planting calender in soreang (Bandung)," 2017 5th Int. Conf. Inf. Commun. Technol. ICoIC7 2017, vol. 0, no. c, 2017.
- [35] F. Martínez-Álvarez, A. Troncoso, H. Quintián, and E. Corchado, "Rainfall Prediction: A Deep Learning Approach," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 9648, pp. 151–162, 2016.
- [36] A. Haidar and B. Verma, "A novel approach for optimizing climate features and network parameters in rainfall forecasting," *Soft Comput.*, 2017.
- [37] S. G. B, S. Sarkar, P. Mitra, and S. Ghosh, "Early Prediction of Extreme Rainfall Events: A Deep Learning Approach," vol. 9728, pp. 154–167, 2016.
- [38] M. Ahmad and S. Aftab, "Analyzing the Performance of SVM for Polarity Detection with Different Datasets," *Int. J. Mod. Educ. Comput. Sci.*, vol. 9, no. 10, pp. 29–36, 2017.
- [39] M. Ahmad, S. Aftab, I. Ali, and N. Hameed, "Hybrid Tools and Techniques for Sentiment Analysis: A Review," *Int. J. Multidiscip. Sci. Eng.*, vol. 8, no. 4, 2017.
- [40] S. Aftab, M. Ahmad, N. Hameed, M. S. Bashir, I. Ali, and Z. Nawaz, "Rainfall Prediction in Lahore City using Data Mining Techniques," *Int. J. Adv. Comput. Sci. Appl.*, vol. 9, no. 4, pp. 254–260, 2018.
- [41] M. Ahmad, S. Aftab, M. S. Bashir, N. Hameed, I. Ali, and Z. Nawaz, "SVM Optimization for Sentiment Analysis," *Int. J. Adv. Comput. Sci. Appl.*, vol. 9, no. 4, 2018.

An Intelligent Bio-Inspired Algorithm for the Faculty Scheduling Problem

Sarah Al-Negheimish¹, Fai Alnuhait², Hawazen Albrahim³, Sarah Al-Mogherah⁴, Maha Alrajhi⁵, Manar Hosny⁶

College of Computer and Information and Sciences
King Saud University
Riyadh, Saudi Arabia

Abstract—All universities have faculty members who need to be assigned to teach courses. Those members have various specialties, preferences and different levels of experience. The manual assignment of courses is a very tedious and time-consuming task that the scheduling committee frequently faces in every department. To solve this timetabling problem, we proposed a novel approach using the Bees Algorithm (BA), which is inspired from bees' foraging behavior, hybridized with Demon algorithm and Hill Climbing for more extensive search. The scheduling process took into consideration all constraints and variables associated with scheduling courses, according to the requirements of the Computer Science department in our college. The results showed that the schedules produced from the algorithm outperformed the manual schedules in terms of achieving the objective function and satisfying the constraints. In addition, the hybridized version produced better results than the standard BA version without hybridization. The hybridized algorithm is designed for faculty scheduling, but can be further generalized to solve various timetabling problems.

Keywords—Faculty scheduling; faculty assignment problem; Bees Algorithm; Demon algorithm; timetabling; scheduling

I. INTRODUCTION

Assigning faculty members to teach courses is a tedious process that must be done by almost every university in the world each semester. Similar to other scheduling problems, faculty scheduling is an NP-hard problem [1] that is very difficult to solve optimally using conventional search methods. The reason behind this is the presence of many constraints that should be taken into consideration, such as the clash of times between courses, the maximum and minimum number of workload hours for each faculty, the preferences and specialties of the faculty, and many more hard requirements that can affect the quality of the solution. For such a difficult problem, the available variables and constraints play a significant role in choosing the method that will solve the problem.

Over the past few years, many different methods have been proposed to solve this problem, some of them are more efficient than others. Meta-heuristic optimization algorithms are among the most effective methods for this type of problem, because they keep improving the proposed solution until it reaches a certain satisfactory quality, although not necessarily the optimum [2].

Bio-inspired algorithms are meta-heuristic algorithms that are widely used in many different fields, because of their effectiveness in solving real life difficult problems that cannot

be solved to optimality given the current computing resources. However, they are not widely used in the faculty scheduling problem, probably due to its complexity and difficulty of formulating its constraints that are involved in the solution method, which makes designing a bio-inspired solution method fairly complex. This paper presents a new intelligent bio-inspired meta-heuristic algorithm, namely the Bees Algorithm (BA), for solving the faculty scheduling problem.

The remainder of this paper is organized as follows: Section 2 provides a brief overview of some related work. Section 3 describes the problem in terms of variables, constraints, and the objective function. Section 4 presents the approach used to design and implement the algorithm. Section 5 demonstrates the results obtained from applying the algorithm using various evaluation techniques. The discussion is presented in Section 6. Finally, Section 7 provides the conclusion and some future work.

II. RELATED WORK

The faculty scheduling problem has been studied for many years as an independent problem or combined with the more general university course scheduling problems [3]–[6]. In the literature, various methods were used by researchers to handle this problem. We will go through some of the different algorithms proposed to solve the problem below.

A. Gunawan and K. Ng [7] solved the teacher assignment problem using Simulated Annealing (SA) [8]–[10] and Tabu search (TS) [11]. The problem was divided into two phases. Phase one is concerned with finding a feasible schedule that satisfies all the hard constraints. Phase two aims to balance the credit hours between the faculty. Their algorithm starts by generating a random schedule, then applying SA as well as TS to improve the initial solution during both phases. The algorithm was tested using two real datasets, and yielded better results compared to the genetic algorithm [12]–[14], and manual allocation. E. Ayca and T. Ayav [15] also used SA to solve the course scheduling problem, which consisted of assigning courses to classrooms and timeslots, in addition to assigning their instructors. Their methodology focused purely on using SA as a strategy to find the best schedule. Based on the performance, that was measured by execution time and quality; they concluded that using SA gave better results compared to manual scheduling.

Parera et al. [16] used a Genetic Algorithm [13], [14] (GA) to solve a bigger problem which includes assigning faculty

members to sections and assigning sections to classrooms. Their goal was to produce a valid schedule free from time clashes, which the algorithm was able to provide. As specified earlier, A. Gunawan and K. Ng [12] used SA and Tabu search to solve the teacher assignment problem. In addition, they also developed a GA to tackle the mentioned problem. After assigning teachers to sections randomly, they used the GA in the reassignment process, where the crossover operator selects randomly two chromosomes (course with its list of teachers) and chooses crossover points to exchange the combination of the two chromosomes. In addition, the mutation operator only changes one gene (teacher) in a randomly chosen chromosome (course). The algorithm provided better schedules than the manual ones. Similarly, Y. OuYang and Y. Chen [17] used GA to solve the course scheduling problem, but with the addition of graph coloring algorithm to generate the initial population, and they used Tabu search in the mutation operation. They concluded that their algorithm performed better, in terms of timetabling speed, than the manually allocated one.

Gunawan et al. [18] proposed a solution for the course scheduling [19], and teacher scheduling problems using greedy heuristic and SA. It starts with assigning teachers to sections using a mathematical approach, then moves on to timetabling the sections into time periods using a greedy algorithm. Finally, SA is used to make improvements. The algorithm was able to solve the two problems simultaneously while providing feasible solutions.

Lastly, M. Hosny [20] proposed a heuristic approach to solve the faculty assignment problem. The designed algorithm can be divided into two phases, the first one starts by iterating over the list of teachers, assigning them sections, and reordering the list according to their assigned hours. The second phase is only needed if there are some sections that remained unassigned. After choosing the best schedule generated, a Hill Climbing Optimization algorithm is used to improve the chosen schedule. Although the algorithm was able to provide satisfying results, it was restricted to assigning labs to Teaching Assistants in the assignment process.

In the coming section, we will describe in detail the problem under consideration in this research, as it paves the first steps into solving the faculty scheduling problem.

III. PROBLEM DESCRIPTION

The faculty assignment problem is defined as assigning teachers to courses, while adhering to a number of pre-specified constraints. In this paper, we define the problem of faculty scheduling in terms of the requirements provided by the Computer Science Department in King Saud University, as variables, constraints, and objective function.

A. Variables

There are only two variables under consideration in our problem: courses, and teachers. Starting with the courses, each course has a number of sections, each of which has a unique number, its specified hours, its type which can be: lecture, tutorial, or lab, and its time slots. Whereas for teachers, they are divided into PhD holders who teach only lectures, BSc holders who teach only tutorial and lab sections, and finally, MSc holders who can teach any type of sections. The final two

categories are further divided into: students (i.e., those who are currently studying postgraduate degrees while working) and full-time.

B. Constraints

The constraints are divided into hard and soft constraints. Violating any of the hard constraints makes the schedule infeasible, whereas violating soft constraints affect the quality of the solution. In our problem, we only have one hard constraint which is the clash of courses' time for the same instructor. Having time conflicts in any teacher's schedule will make the whole solution infeasible. As for the soft constraints, we have five constraints. First, ensuring that each teacher's assigned hours are within the boundaries of a certain predefined minimum and maximum workload, which is stated for each instructor based on their rank. Second, fulfilling teachers' preferences for courses represented as a wish list. Third, balancing the workload amongst teachers within the same rank, so that there would be no significant differences. Forth, ensuring that each instructor gets at least a day off in their schedule for course preparation and research. Lastly, minimizing the number of instructors per course, as it can ease the teaching and coordination process.

C. Objective Function

It is crucial to define how to measure the quality of the generated schedules in order to evaluate the outcomes of the algorithm. First, we assume that any violation of our hard constraint (i.e., time clashes in a teacher's schedule) is not acceptable. Thus, we only measure the quality of feasible solutions. To measure the quality of the solution we count each soft constraint violation, multiplied by its weight, where the weight is determined based on the importance of the soft constraint. We also add to this, the number of unassigned sections in the current solution, since our main target is to assign all sections to instructors. The objective function is represented in (1) below. Intuitively, the closer the objective function is to zero, the better the solution would be.

$$\text{Minimize } u + \sum_{i=1}^n (w_i c_i). \quad (1)$$

Where, u is the number of unassigned sections in the solution, w_i is the penalty weight of the i^{th} (soft) constraint, c_i is the number of violations of the i^{th} (soft) constraint, n is the number of constraints.

To ensure the integrity of the objective function, the measurement of the violations for each constraint was normalized within the range [0, 1].

The following section demonstrates our proposed approach to develop a bio-inspired meta-heuristic for solving the faculty scheduling problem.

IV. METHODOLOGY

In this section, we propose a bio-inspired approach, based on the bees' foraging behavior, to solve the stated problem. The described algorithm hybridizes the Bees Algorithm [21] with the Demon Algorithm [2], and Hill Climbing [2]. The algorithm is inspired from N. Alhuwaisel and M. Hosny [22], where a hybrid Bees-Demon Algorithm was used to solve the University Course Timetabling problem. However, to the best

of our knowledge this kind of hybridization has not been attempted before for the faculty scheduling problem.

The Bees algorithm provides a variety of solutions as it starts with a population of solutions called the initial population. This approach is needed to explore the search space more; thus, gives a higher probability of finding near optimal solutions [23]. The basic idea of the Bees algorithm involves the generation of multiple initial solutions then focusing on exploiting the best and elite (best of the best) solutions, whilst performing modification to increase the chance of finding near optimal solutions. the steps used to apply the Bees algorithm, shown in Fig. 1, can be described as follows [21]:

- a) Construct the initial population of schedules.
- b) Evaluate fitness of the population using the objective function.
- c) Selecting the best and elite solutions.
- d) Perform neighborhood search on the selected sites (Hill Climbing), with more search to be done around elite solutions (Demon Algorithm).
- e) Repopulate the region after removing discarded solutions.

The algorithm consists of two major parts: initial population construction, and neighborhood search. These subparts will be combined to design the overall algorithm.

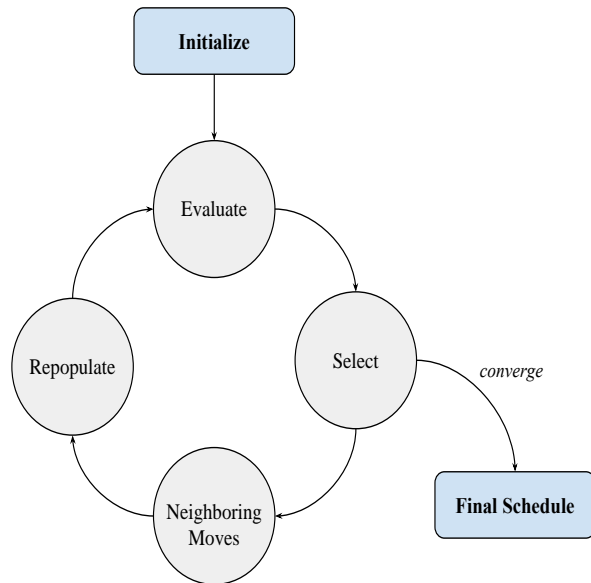


Fig. 1. Abstract representation of the hybrid algorithm.

A. Initial Population Algorithm

The first step in our proposed algorithm is to construct a population of feasible schedules that do not violate the hard constraint. To accomplish this, we designed a greedy randomized heuristic algorithm to populate the search space. After categorizing the sections into separate time slots to avoid time clashes during the assignment, the algorithm proceeds in three phases. The first phase tries to assign only the minimum

workload for each faculty from their wish list, while the second phase tries to assign the remaining workload for the faculty, but this time not necessarily from their wish list. The third phase, on the other hand, is considered only if there are still sections not assigned to faculty after the first and second phases. In Phase three, assigning faculty hours above their maximum workload (within a certain predefined percentage) is attempted. In each of these phases, we order the list of faculty based on their rank, which is calculated by their job position, administrative posts, and seniority. Within the second and third phases, the number of assigned hours is taken into consideration. This is intended to give higher priority in assignment to those faculty on top of the list. For those whose ranking criteria are equivalent, the list is ordered randomly. Using these procedures, we obtain different schedules to create the population. The steps of the algorithm are described in Algorithm 1:

Algorithm 1: Initial population generation

Input: list of faculty f , list of sections s , the percentage c to increase workload, population size $n_{SEP}^{[1]}$
Output: population of feasible solutions P

- 1: $P \leftarrow \{ \}$
- 2: **repeat** while $|P| < n_{SEP}^{[1]}$
- 3: **categorize**(s)
/ Phase 1: assigning minimum workload */*
- 4: **order**(f)
- 5: **for** $i = 1: |f|$
- 6: **assign**(f_i, s, min)
/ Phase 2: assigning maximum workload */*
- 7: **order**(f)
- 8: **for** $i = 1: |f|$
- 9: **assign**(f_i, s, max)
/ Phase 3: assigning above maximum workload */*
- 10: **if** $s \neq \{ \}$
- 11: **for** $i = 1: |f|$
- 12: $f_i(workload) \leftarrow f_i(workload) \times c$
- 13: **assign**(f_i, s, max)
- 14: $P \leftarrow P + f$

The algorithm starts with an empty set of population P , then repeats the following steps until it reaches the desired number of schedules n . At the first step, the list of sections s is divided into buckets according to their timeslot using **categorize**. This is done to ensure that any time clashing incident does not occur, thus preventing infeasible solutions. The **order** method takes the list of faculty f and sorts it according to the member's position, seniority, held admin post, and assigned hours. In cases where more than a member share the same criteria, their order will be chosen randomly. When it comes to the **assign** task, the method tries to assign f_i with a matching section according to either the maximum or minimum workload for f_i . It is important to note that in phase one's call of **assign**, s is sorted according to f_i 's preference. However, that is not taken into consideration within the subsequent calls to **assign** in phases two and three. At the end, the algorithm checks whether s still has sections that need to be assigned. If so, every faculty's workload will be increased by percentage c . Fig. 2 represents these procedures.

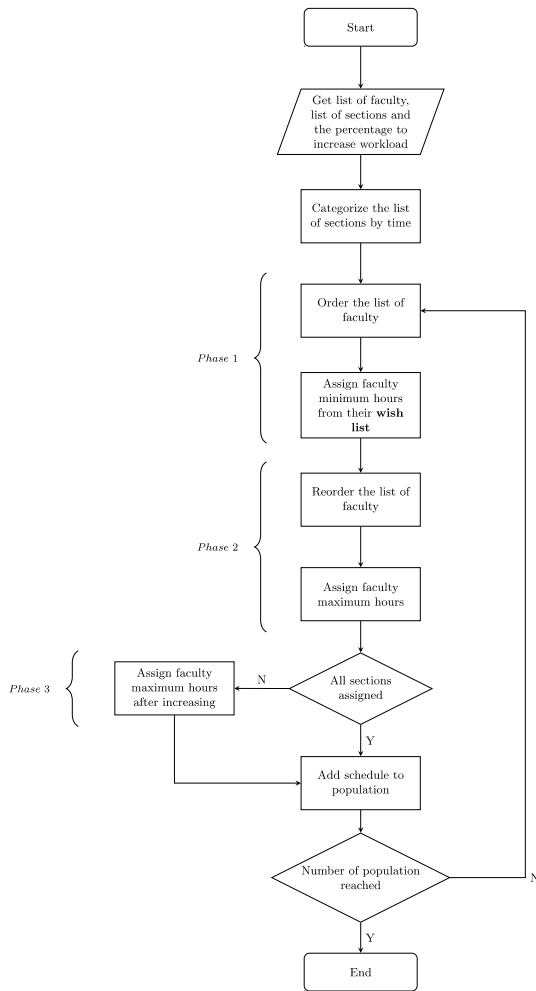


Fig. 2. Initial population generation flowchart.

B. Neighborhood Search

Neighboring moves are functions applied on already existing solutions to generate new slightly different ones, hoping to find better solutions in their vicinity. The most important characteristic of a neighboring move is locality, which is its effect on the original solution. Locality means that small changes applied to the original solution representation (e.g. changing the location of some variables) should correspond to small changes in the actual solution (i.e., real schedule). Otherwise, the new solution will have drastic changes different from the original solution, which, in extreme cases, could cause the search to converge towards a random solution [2]. Listed below are the proposed neighboring moves used to generate new solutions, which are intended to reduce the violations of the soft constraints in our problem.

1) Balancing the Workload

This procedure is applied when two different instructors, who belong to the same rank and position, are not equivalent in the number of hours assigned. This is accomplished by calculating the average hours of each group then segregating them into above-average and below-average. The algorithm

then tries to transfer sections from the above-average group to the below-average group, if no violations in constraints occur.

2) One-day off

This function aims to increase the number of faculty members who have a day off, by swapping sections between faculty members if it ensures that both individuals enjoy a day off, have no conflicts in assigned courses, and can be assigned the type of section.

3) Max number of Faculty Members per Course

Each course has a limit on the number of different instructors that should not be exceeded. Therefore, this neighboring move tries to lessen the number of faculty members for these courses. This is achieved by setting a limit to each course, then it investigates the possibility of consolidating a course. The approach used in the algorithm explores potential faculty members that share more than one course together, then checks if the course sections could be swapped to reduce the number of instructors for either course.

C. Proposed Hybrid Algorithm

As previously mentioned, the main algorithm we use is inspired from the Bees Algorithm [21]. The algorithm starts by generating a collection of schedules named the population; then it selects the highest scoring schedules according to their fitness (objective function value). Then it divides the group into best schedules and elite schedules, which are the best of the best, creating two mutually exclusive sets. In our proposed method, we have also incorporated a Hill Climbing [2] approach for accepting the new schedule after applying the neighboring moves on the best schedules. In other words, if the new solution produced after applying a neighboring move on each of the best solutions is better in terms of the objective function, it replaces the previous solution. However, we took a different approach for searching around the elite schedules, where we apply the Demon Algorithm (DA). The DA is a variant of the famous Simulated Annealing (SA) approach but with a deterministic acceptance function [2]. The DA potentially accepts worse solutions, based on a certain credit value called the demon, in hopes of finding better solutions in the next iterations. Thus, the idea is to do a more intelligent search around elite solutions, in an attempt to discover even better solutions as the search progresses. The intensifying phase of the algorithm continues for each schedule until the schedule's fitness ceases to change for five consecutive iterations.

The first step in the proposed algorithm is to create a population of schedules using the heuristic algorithm described in Section 4.1. After the initialization of all the variables, including the demon credit d , we repeat the following process until convergence.

Firstly, the program evaluates every schedule in the population then selects the best and elite schedules from the population. Secondly, we intensify the search on the best and elite schedules using the neighboring moves for best and elite respectively. When applying neighboring search on the 'best' schedules, the new schedule is accepted if the overall fitness is better. However, when applying the search on 'elite' schedules, and based on the principle of the DA, the program calculates

the difference in fitness value of the two schedules, the old schedule and the new schedule. If the difference is positive (+), the new schedule is accepted and the difference value is credited to the demon value. On the other hand, if the difference is negative (-), the new schedule is accepted only if the demon credit can withstand the difference and then we update the demon value after subtracting the difference value. This process is considered as the intensification part of the algorithm as it intensifies the search on each individual schedule to improve it.

Finally, we repopulate the search space keeping the best and elite schedules from the previous iteration, and keep the current best schedule in hand. Diversification is incorporated in the algorithm as it helps to generate a population with various fitness values. It is worth mentioning that the selection of best set is done in a stochastic manner [24] using a tournament selection approach, as the algorithm chooses five schedules randomly from the population, and then selects the best schedule among them. This process is repeated until reaching the required number of best set. This approach gives better results in terms of diversification and producing different schedules, as deterministic methods tend to stick with the same best set regardless of the iterations [2].

Once the best schedule ceases to change for five consecutive iterations, the algorithm stops repetition and outputs the result. The details are described in Algorithm 2.

The algorithm starts by receiving a population of schedules from Algorithm 1 then evaluates every schedule in the population. Method *eval* does precisely that by taking a schedule and producing the objective function value for each corresponding schedule. After that, a set of *best* and *elite* schedules is selected respectively based on the fitness value in our case. Once the sets are established, the algorithm first iterates over the *best* set and attempts to improve the schedules using the method *neighborhood* then moves on to the *elite* schedules. *Neighborhood* implements the techniques discussed in Section 4.2 to improve the quality of schedules in both sets. However, improving criteria must be chosen for the schedules in both cases; this is selected to be ‘Hill Climbing’ for *best* and ‘Demon’ for *elite*, which are passed to the method *neighborhood*. After these steps have been conducted, the algorithm repopulates the space with the inclusion of the *best* and *elite* sets, then proceeds with the next iteration.

Algorithm 2: Hybridized Algorithm

Input: demon credit d

Output: an optimal schedule s

```
1:  $p \leftarrow \text{Initial Population}$  (Algorithm 1)
2: repeat until convergence:
3:   for  $i = 1: |p|$ 
4:     eval( $p_i$ )
5:    $b \leftarrow \min(p)$  //best schedules
6:    $e \leftarrow \min(b)$  //elite schedules
7:    $b \leftarrow b - e$ 
8:   for  $i = 1: |b|$ 
9:     repeat until convergence:
10:    neighborhood( $b_i$ , hill climbing)
```

```
11: for  $i = 1: |e|$ 
12:   repeat until convergence:
13:     neighborhood( $e_i$ , demon,  $d$ )
14:    $p \leftarrow \text{Initial Population} + b + e$ 
15:   for  $i = 1: |p|$ 
16:     eval( $p_i$ )
17:    $s \leftarrow \min(p)$ 
18: return  $s$ 
```

In the next section, details of how the results were evaluated will be illustrated to conduct analysis on the performance of the designed algorithm.

V. EXPERIMENTAL RESULTS

In this section, we will discuss in details the dataset that we used to test our algorithm, and show our results and the different criteria used to evaluate these results.

A. Characteristics of the Dataset

To test our algorithm, we used a real dataset obtained from the Computer Science department at King Saud University. In addition, we created another theoretical (i.e., synthesized) dataset inspired from the real dataset. We conducted the testing on one real dataset and two different theoretical datasets. Both types were saved on (.csv) files. Table I summarizes the details of the three datasets. Each dataset is split into two different parts:

- The dataset of the teachers to be allocated: This includes: professors, associate professors, assistant professors, lecturers, and teaching assistants.
- The dataset that contains the information about the courses and sections to be taught. This includes: lectures, labs, and tutorials.

B. Parameter Tuning

The objective function was designed to measure the quality of a schedule. A crucial part of it was to decide the weights used to penalize the violations of soft constraints. The weights were selected to be in the range [0, 1] and the totality of them would equal to one. We prioritized the soft constraints with their corresponding weights depending on the department's needs:

- Minimum and maximum workload, $w_1 = 0.4$.
- Balancing the workload, $w_2 = 0.3$.
- One-day off, $w_3 = 0.2$.
- Minimize the number of instructors per course, $w_4 = 0.1$.

After designing our algorithm, we had to test different values for each of the parameters: population size, best size, elite size, and demon credit.

For each of these parameters we used the same dataset. In parameter tuning [25], [26], we applied grid search technique to test a variety of values for one parameter (approximately 13 different values) keeping the rest of values constant. Each value was tested five times then we compared the resulting average fitness with that of the rest of the values.

TABLE I. DATASET SUMMARY

Variable	Datasets		
	Real Dataset	Theoretical 1	Theoretical 2
Courses	22	7	10
Sections	167	70	120
Professors	0	1	0
Associate Prof.	1	1	1
Assistant Prof.	8	1	6
Lecturers	16	9	13
Teaching Assist.	20	5	14

In addition, we chose to start our tuning with the population size, since it was observed to have the most effect on the fitness. Then we tested the number of best solutions, and then we tuned the number of elite solutions, since it is a subset of the best. Finally, we tuned the demon credit. The final parameter values we obtained were, population size = 50, best size = 8, elite size = 4, and demon credit = 0.25, which were deemed fit for our algorithm based on the quality of solutions obtained.

C. Evaluating the Results by Analyzing the Objective Function

We ran our algorithm twenty times on each of the three previously defined datasets, documenting the fitness and the execution time (in minutes). The program was executed on a Macbook Pro with OS X Yosemite, 2.9 GHz Intel Core i7 processor, 8 GB DDR3 memory.

In Fig. 3, a sample of the fitness (objective function) of the best schedule at hand is illustrated against the iteration sequence. It is clear that in every iteration, the intensification phase of the algorithm provides a major contribution to the fitness as well as finding another potential candidate to substitute the current schedule through the diversification phase.

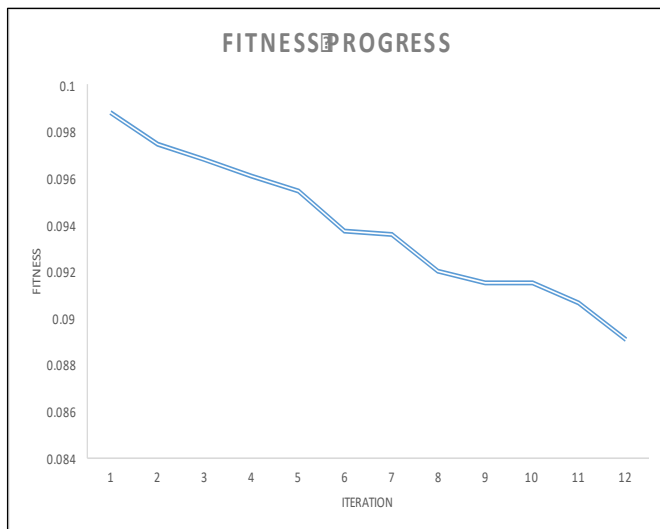


Fig. 3. Progress of best schedule over time.

TABLE II. DATASET RESULTS SUMMARY

Measure	Datasets		
	Real Dataset	Theoretical 1	Theoretical 2
Average Fitness	0.062	0.033	0.078
Standard Deviation	0.013	0.010	0.007
Minimum Fitness	0.029	0.002	0.059
Maximum Fitness	0.084	0.049	0.089

TABLE III. DATASET RUNTIME SUMMARY

Criteria	Datasets		
	Real Dataset	Theoretical 1	Theoretical 2
Number of Sections	167	70	120
Avg. Exec. Time (min)	5.17	1.775	4.35

As demonstrated by the results in Table II, all of the obtained values on all data sets are very close to zero, indicating very small number of violations of the soft constraints. Also, the values obtained fall into the same range of values, meaning that the algorithm produces a good result in each run. This is confirmed by measuring the standard deviation to ensure the stability of our algorithm. On the other hand, the maximum result shows the worst case, which is still very close to zero.

Understandably, the algorithm could not give a schedule with a zero-fitness value, as there are many factors that prevent the algorithm from reaching the ideal solution. For instance, different sections have different hours such as labs and lectures. Thus, it is nearly impossible to have the workload distributed perfectly between the faculty members. Nonetheless, considering all the constrains, we deemed important in the algorithm, the results shown are indeed very satisfactory.

Regarding the execution time, it is noticeable that the execution time differs from run to run, and that is understandable as well, since the time it takes for the algorithm to converge is distinct depending on the population in hand and the number of iterations it will take to improve various schedules in the dataset during the intensification phase. Table III above summarizes the relationship between the dataset size and execution time in minutes. Overall, the algorithm produces excellent results in a reasonable processing time.

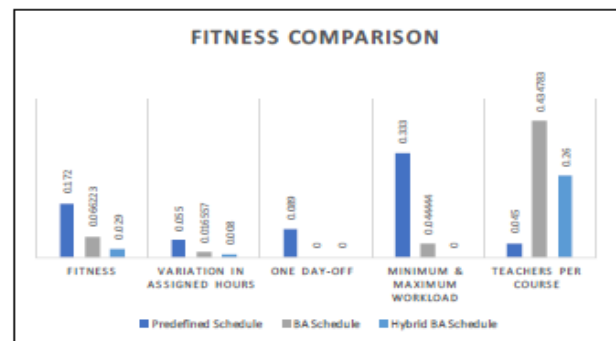


Fig. 4. Comparison between manual, BA and hybrid BA schedules.

D. Evaluating the Results by Comparing the Manual Schedule with our Algorithm's Schedule

The manually designed schedule was compared with the generated schedule by our algorithm on the real dataset to assess the latter. The assessment is based on fulfilling the soft constraints, since the hard constraint is guaranteed to be satisfied to make a valid schedule. Table V shows the results of the comparison in details.

As can be seen from Table V, our algorithm provided a schedule that eliminated any violations of 'min & max workload' constraint as well as the 'day off' constraint. In addition, it managed to balance the workload by a significant amount, since it provided an overall total variance from the average workload of just 0.008, compared to a total variance of 0.055 for the manually allocated schedule. Evaluating the overall fitness of the schedule, our schedule has a fitness of approximately 0.03 while the manual schedule has a fitness of approximately 0.17. This shows that our schedule improved the manual schedule by a remarkable 83.4% value. However, the only constraint that our algorithm was not able to provide satisfactory results for was minimizing the number of instructors per course shown evidently in Fig. 4. This is probably due to the small penalty weight that we assigned for this constraint, since we considered it less important in the schedule than the other constraints.

E. Evaluating the Results by Comparing the Proposed Algorithm with the Classic BA

To assess the effectiveness of hybridizing BA with Demon and Hill Climbing algorithms, we compared the results of our algorithm to the results produced from the classic BA using the real dataset following the same criteria used in the previous

section. As shown in Table V, the hybrid algorithm outperforms the classic BA in generating the best schedule. The two algorithms performed equally well in avoiding the violation of the day off constraint. However, the hybrid algorithm produced considerably better results in fulfilling the rest of the soft constraints, with a 56% improvement in balancing the workload. Digging deeper into these schedules, we noticed a notable difference in the quality of the produced solutions at the instructors' level. While the schedules produced by the hybrid algorithm tend to limit the number of different courses and assign a consistent set of sections to each instructor, the classic BA generates more variant ones. An example to further explain this point is illustrated in Table IV below, where we took one of the instructors and compared her schedule generated by the two approaches. We can clearly see that the hybrid BA schedule is more practical and more convenient for the instructor.

TABLE IV. INSTRUCTOR ASSIGNED SECTIONS EXAMPLE

	BA Schedule			Hybrid BA Schedule		
	Course	Type	Hours	Course	Type	Hours
	CS09	Tutorial	1	CS17	Lab	6
	CS23	Lecture	3	CS23	Tutorial	2
	CS10	Tutorial	1			
	CS01	Tutorial	1			
Number of Sections	4			5		
Unique courses	4			2		
Total Workload	6			8		

TABLE V. MANUAL VS CLASSIC BA VS HYBRID BA RESULT SUMMARY

Category	Manual Schedule				BA Schedule				Hybrid BA Schedule			
	Min and max workload violation	Balance workload violation	Day off violation	Number of instructors per course violation	Min and max workload violation	Balance workload violation	Day off violation	Number of instructors per course violation	Min and max workload violation	Balance workload violation	Day off violation	Number of instructors per course violation
Students	2	0.003	0		0	0	0		0	0	0	
Professors	7	0.017	0		0	0.003	0		0	0.004	0	
Lecturers/TA	6	0.035	4		2	0.014	0		0	0.004	0	
All	15	0.055	4	1	2	0.017	0	10	0	0.008	0	6
Total Fitness	0.172				0.066				0.029			

TABLE VI. BA VS HYBRID BA AVERAGE COMPARISON

Criteria	BA	Hybrid BA
Avg. Fitness	0.076	0.062
Avg. Exec. Time	0.82	5.17

Taking the performance evaluation of the proposed algorithm a step further, we compared the two approaches in terms of the average fitness obtained along with the average time needed to generate the solutions, running each algorithm 20 times on the real dataset to get the average value.

As Table VI demonstrates, the hybrid algorithm generates better results in general. Although the improvement is not very significant on an average scale, our goal is to find the best fitted schedule which will be obtained through running the algorithm multiple times and adopting the best schedule. In other words, increasing the chance of finding a near optimal solution on a set of satisfactory solutions would be more beneficial than trying to ensure that all solutions in the set are optimal solutions. The difference in the average running time is not a concern as well, since the algorithm will only be run once each semester in practical situations. So, we can sacrifice the

increase in execution time for the sake of obtaining a much better schedule that will be adopted throughout the semester.

F. Evaluating the Results with the Assessment of the Schedule by the Scheduling Committee

Finally, our results were also assessed by members of the scheduling committee in the Computer Science Department, by answering an evaluation survey. Overall, we received a positive feedback from the scheduling committee. Generally, they strongly agreed that the algorithm fulfilled the requirement of assigning all courses to faculty members, as well as allowing each faculty member a day off per week. They also agreed that the algorithm managed to balance the workload among the faculty members. Moreover, the committee agreed that the quality of the schedule produced was satisfactory, and that they consider the resulting schedule reliable.

Finally, the scheduling committee strongly admits that our algorithm is needed and useful for the Computer Science department, and would use it if it was currently available.

VI. DISCUSSION

After applying our proposed method to solving the faculty scheduling problem, it is evident that the Bees Algorithm proved its capability and suitability for this problem. Specifically, the diversification stage played a significant role in the exploration of many different solutions. This was achieved through the greedy-randomized population creation part of the algorithm. Whereas intensification further improved the solutions obtained that being the neighboring moves' role, focusing on minimizing the violations of the soft constraints.

Moreover, hybridizing the Bees Algorithm with another meta-heuristic immensely improved our algorithm's performance, leaping to a higher level of intelligence. We used both Hill Climbing and Demon algorithms as solution acceptance algorithms. The Hill Climbing algorithm was applied on the best and elite solutions, whilst the Demon algorithm was only applied on the elite solutions, with the intension of doing more intelligent search around the elite than the other selected best solutions.

VII. CONCLUSION

In this paper, we tackled the faculty scheduling problem, which is concerned with assigning faculty members to prescheduled courses. To solve the problem, firstly, we designed the construction of the initial population that is considered a primary factor in the Bees Algorithm. We used a specially designed greedy-randomized heuristic for this purpose. Secondly, we designed the neighboring moves that will be used to improve the solutions selected by the algorithm. We hybridized the Bees algorithm with the Demon algorithm and Hill Climbing, which is considered an innovative approach in this particular problem.

We used the dataset provided by the CS department to test our algorithm and chose to use this dataset to evaluate our algorithm, because it portrays a realistic environment. We also used two theoretical datasets to further test the algorithm. The algorithm showed superior results when compared to the manually allocated one, as it managed to eliminate 'min & max

workload' constraint as well as the 'day off' constraint. Moreover, the scheduling committee in the department evaluated the schedules produced by the algorithm, and agreed that it satisfies their expectations.

Several areas of improvement arise, though, by enhancing the hybrid algorithm to solve some additional requirements, such as minimizing the number of courses assigned to each teacher. Further research could be conducted by broadening the problem and generalizing the algorithm to solve other variations of scheduling.

REFERENCES

- [1] S. Even, A. Itai, and A. Shamir, "On the complexity of time table and multi-commodity flow problems," in Foundations of Computer Science, 1975., 16th Annual Symposium on, 1975, pp. 184–193.
- [2] E.-G. Talbi, *Metaheuristics: from design to implementation*, vol. 74. John Wiley & Sons, 2009.
- [3] D. Abramson, "Constructing school timetables using simulated annealing: sequential and parallel algorithms," *Manag. Sci.*, vol. 37, no. 1, pp. 98–113, 1991.
- [4] T. Ferdoushi, P. K. Das, and M. A. H. Akhand, "Highly constrained university course scheduling using modified hybrid particle swarm optimization," in Electrical Information and Communication Technology (EICT), 2013 International Conference on, 2014, pp. 1–5.
- [5] F. Aloul, I. Zabalawi, and A. Wasfy, "A SAT-based approach to solve the faculty course scheduling problem," in AFRICON, 2013, 2013, pp. 1–5.
- [6] R. Lewis and B. Paechter, "Finding feasible timetables using group-based operators," *IEEE Trans. Evol. Comput.*, vol. 11, no. 3, pp. 397–413, 2007.
- [7] A. Gunawan and K. M. Ng, "Solving the teacher assignment problem by two metaheuristics," *Int. J. Inf. Manag. Sci.*, vol. 22, no. 2, pp. 73–86, 2011.
- [8] S. Kirkpatrick, C. D. Gelatt, M. P. Vecchi, and others, "Optimization by simulated annealing," *science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [9] R. W. Eglese, "Simulated annealing: a tool for operational research," *Eur. J. Oper. Res.*, vol. 46, no. 3, pp. 271–281, 1990.
- [10] A. G. Nikolaev and S. H. Jacobson, "Simulated annealing," in *Handbook of Metaheuristics*, Springer, 2010, pp. 1–39.
- [11] F. Glover and M. Laguna, *Tabu Search**. Springer, 2013.
- [12] A. Gunawan, K. M. Ng, and H. L. Ong, "A genetic algorithm for the teacher assignment problem for a university in Indonesia," *Inf. Manag. Sci.*, vol. 19, no. 1, pp. 1–16, 2008.
- [13] J. H. Holland, "Genetic algorithms," *Sci. Am.*, vol. 267, no. 1, pp. 66–72, 1992.
- [14] M. Mitchell, "Genetic algorithms: An overview," *Complexity*, vol. 1, no. 1, pp. 31–39, 1995.
- [15] E. Ayca and T. Ayav, "Solving the course scheduling problem using simulated annealing," in *Advance Computing Conference, 2009. IACC 2009. IEEE International*, 2009, pp. 462–466.
- [16] S. Parera, H. T. Sukmana, and L. K. Wardhani, "Application of genetic algorithm for class scheduling (Case study: Faculty of science and technology UIN Jakarta)," in *Cyber and IT Service Management, International Conference on, 2016*, pp. 1–5.
- [17] Y. OuYang and Y. Chen, "Design of automated Course Scheduling system based on hybrid genetic algorithm," in *Computer Science & Education (ICCSE), 2011 6th International Conference on, 2011*, pp. 256–259.
- [18] A. Gunawan, K. M. Ng, and K. L. Poh, "Solving the teacher assignment-course scheduling problem by a hybrid algorithm," *Int J Comput Inf. Engin.*, vol. 1, no. 2, pp. 137–142, 2007.
- [19] M. W. Carter and G. Laporte, "Recent developments in practical course timetabling," in *International Conference on the Practice and Theory of Automated Timetabling, 1997*, pp. 3–19.
- [20] M. I. Hosny, "A Heuristic Algorithm for Solving the Faculty Assignment Problem," in *Proceedings of the International Conference on Frontiers in*

- Education: Computer Science and Computer Engineering (FECS), 2012, p. 1.
- [21] D. T. Pham, A. Ghanbarzadeh, E. Koc, S. Otri, S. Rahim, and M. Zaidi, "The bees algorithm—a novel tool for complex optimisation," in *Intelligent Production Machines and Systems-2nd I* PROMS Virtual International Conference (3-14 July 2006)*, 2011.
- [22] N. ALHUWAISHEL and M. HOSNY, "A Hybrid Bees/Demon Optimization Algorithm for Solving the University Course Timetabling Problem," in *Proceedings of the 3rd NAUN International Conference on Mathematical, Computational and Statistical Sciences*. Dubai, United Arab Emirates, February, 2015.
- [23] X.-S. Yang, S. Deb, and S. Fong, "Metaheuristic algorithms: optimal balance of intensification and diversification," *Appl. Math. Inf. Sci.*, vol. 8, no. 3, p. 977, 2014.
- [24] J. E. Baker, "Reducing bias and inefficiency in the selection algorithm," in *Proceedings of the second international conference on genetic algorithms*, 1987, pp. 14–21.
- [25] A. E. Eiben, R. Hinterding, and Z. Michalewicz, "Parameter control in evolutionary algorithms," *IEEE Trans. Evol. Comput.*, vol. 3, no. 2, pp. 124–141, 1999.
- [26] A. E. Eiben and S. K. Smit, "Evolutionary algorithm parameters and methods to tune them," in *Autonomous search*, Springer, 2011, pp. 15–36.

A Lightweight Multi-Message and Multi-Receiver Heterogeneous Hybrid Signcryption Scheme based on Hyper Elliptic Curve

Abid ur Rahman

IT Department Hazara University, Mansehra

Noor-ul-Amin, Hizbullah Khattak

IT Department Hazara University, Mansehra

Insaf Ullah, Muhammad Naeem, Rehan Anwar

IT Department, Abbottabad University of Science and Technology, Abbottabad

Sultan Ullah

IT Department, University of Science and Technology, Haripur

Abstract—It is a suitable means for multi-messages to use hybrid encryption to make a safe communication. Hybrid encryption confines encryption into two parts: one part uses public key systems to scramble a one-time symmetric key, and the other part uses the symmetric key to scramble the actual message. The quick advancement of the internet technology requires distinctive message communications over the more extensive territory to upgrade the heterogeneous system security. In this paper, we present a lightweight multi-message and multi-receiver Heterogeneous hybrid signcryption scheme based on the hyper elliptic curve. We choose hyper elliptic curve for our scheme, because with 80 bits key give an equivalent level of security as contrasted and different cryptosystems like RSA and Bilinear pairing with 1024 bits key and elliptic curve with 160 bits key, respectively. Further, we validate these security requirements with our scheme, for example, confidentiality, resistance against replay attack, integrity, authenticity, non-repudiation, public verifiability, forward secrecy and unforgeability through a well-known security validation tool called Automated Validation of Internet Security Protocols and Applications (AVISPA). In addition, our approach has low computational costs, which is attractive for low resources devices and heterogeneous environment.

Keywords—Multi-receiver heterogeneous hybrid signcryption; multi-message and multi-receiver heterogeneous hybrid signcryption; hyper elliptic curve; Automated Validation of Internet Security Protocols and Applications (AVISPA)

I. INTRODUCTION

To communicate securely through a harmful network, people need the security services like authentication, integrity, confidentiality, and non-repudiation [1]. Authentication, integrity, and non-repudiation can be ensured through digital signature [2]-[7] and confidentiality can be assured through encryption [8]-[10] algorithms. In old mechanisms, the sender first signs the message and then encrypts them by using digital signature and encryption algorithms. This type of method was namely called signature-then-encryption. The approach requires more computational power, more bandwidth consumption and more machine cycle [11]. To resolve the deficiencies of old signature-then-encryption approach

signcryption was introduced [11]. Signcryption is the cryptographic primitives which combine the properties of encryption and digital signature in one logical step. After this, numbers of signcryption schemes were projected to the literature [12]-[31]. These signcryption schemes can be filled, if applications need multicast communication. Unlike unicasting, multicast communication is a proficient means to deliver a same copy of signcryptext to multicast group with less bandwidth consumption and fewer computation powers. These like of features make multicast communication an idyllic technology during if an application needs communication with group of receiver. Further, secure multicast communication attracted so many applications such as real time video conferencing, distance education and military command and control [32], respectively. For multicast communication, Zheng [33] was the pioneer to contribute a multi-receiver signcryption scheme. The proposed multi-receiver signcryption scheme enables the signcrypter to signcrypt a single message for the multi - receiver group. After, successful generation of signcryptext, then it delivered the same copy of signcryptext to multiple group. Recently, heterogeneous signcryption mechanisms have got significant attention in so many cryptographic applications [34]-[37]. It is a viable means for extensive messages to utilize hybrid encryption to create secure communication. Hybrid encryption isolates encrypted into two sections: one section utilizes public key strategies to scramble a one-time symmetric key, and the other part utilizes the symmetric key to scramble the genuine message [38], [39]. The fast advance of the internet requires different message corresponding over the more extensive territory to enhance the heterogeneous network security. To deal with these like circumstances, enhance the selection of the security prerequisites and to build the speed of data transmission for numerous messages, multi-messages signcryption were presented [40]-[43]. Recently, Shufen et al. [44] designed a Heterogeneous hybrid signcryption scheme for transmitting multi-messages to multi-receiver group. The designed approach thus suffered from replay attack and leads high computational cost due to heavy pairing operations.

Considering all the above multi-message and multi-receiver approaches, it can be suffered from high

computational cost. Because these approaches are based on RSA, Bilinear pairing and Elliptic curves, which are prominent techniques for security mechanisms. On the other hand, the Hyper-elliptic Curve Cryptosystem (HECC) with 80 bits key give an equivalent level of security as contrasted and different cryptosystems like RSA and Bilinear pairing with 1024 bits key and elliptic curve with 160 bits key, respectively. Accordingly, to reduce computational costs, we present a lightweight multi-message and multi-receiver heterogeneous hybrid signcryption scheme based on the hyper elliptic curve. Our presented scheme, give the security requirements, for example, confidentiality, integrity, authenticity, unforgeability, non-repudiations and forward secrecy. In addition, we validate these security requirements through a well-known security validation tool called Automated Validation of Internet Security Protocols and Applications (AVISPA). Furthermore, our approach has reduced computational costs, which is attracted for low resources devices and heterogeneous network environment.

II. PRELIMINARIES

The hyper elliptic curve is the short form of elliptic curves, which was initially tossed by N. Koblitz [45]-[50]. The most important factor of every cryptographic system is the discrete logarithm problem in some Abelian group. Let them choose a random number γ from the Abelian group and calculating $\gamma \cdot \mathcal{D} = \mathcal{D} + \mathcal{D} + \mathcal{D} +, \dots \dots + \mathcal{D}$ is scalar multiplication of divisors. And it is said to a hyper elliptic curve discrete logarithm problem because finding the random number γ from $\gamma \cdot \mathcal{D} = \mathcal{D} + \mathcal{D} + \mathcal{D} +, \dots \dots + \mathcal{D}$ is infeasible.

III. PROPOSED MODEL

In this sub-section, we present our newly designed a lightweight multi-message and multi-receiver Heterogeneous hybrid signcryption scheme based on the hyper elliptic curve. The security hardness and efficiency of our design scheme is based hyper elliptic curve discrete problem (\mathcal{HECDLP}). Because the hyper elliptic curve has lower known security simulation tool called Automated Validation of Internet Security Protocols and Applications (AVISPA). Our designed scheme constructed by using five phases, such as Key Generation, the Basic Notations used in the proposed scheme, Multi-Messages Signcryption Phase, Unsigncryption Phase and Signature Verification, respectively. Here in Fig. 1, we illustrate the block diagram of our designed lightweight multi-message and multi-receiver Heterogeneous hybrid signcryption scheme based on the hyper elliptic curve. In our designed scheme, before starting the communication, the signcrypter first verify the public keys each receiver, then generate the multi-message signcryptext and deliver to multi-receiver group. After receiving the signcryptext text the each unsigncrypter first confirm the public key of sender. Latter, each unsigncrypter verify the signature and decrypt the cipher text.

\mathcal{C}_i ← secret key for each receiver

\mathcal{C}_j ← encrypted messages for each receiver

enc ← encryption

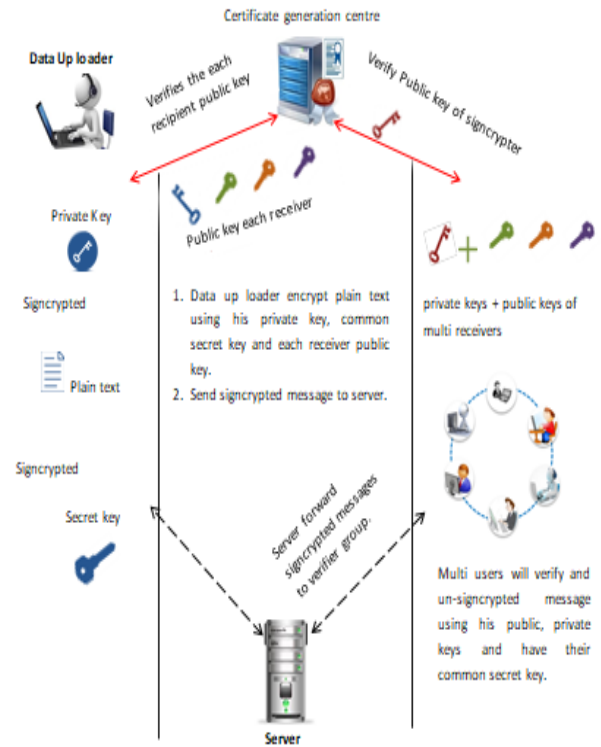


Fig. 1. Block diagram of proposed scheme.

IV. BASIC NOTATIONS

The following are the basic notations which are used in our proposed algorithm:

\mathcal{D} ← Divisor on hyper a elliptic curve

m_j ← plaintext(multi-messages) Q ← signature

$\mathcal{K}_a, \mathcal{K}_b$ ← secret keys

\mathcal{W}_s ← private key of multi-messages-signcrypter

$\mathcal{X}_s = \mathcal{W}_s \cdot \mathcal{D}$ ← public key of multi-messages-signcrypter

\mathcal{W}_i ← private key of each multi-messages-unsigncrypter

$\mathcal{X}_i = \mathcal{W}_i \cdot \mathcal{D}$ ← public key of each multi-messages-unsigncrypter

\mathcal{N}_r ← nonce

i ← receiver group

h ← one-way hash function

Dec ← decryption

\mathcal{L}, \mathcal{V} ← random numbe

A. Multi-Messages Signcryption Phase

In this first step multicast signcrypted text $(\mathcal{C}_1, \dots, \mathcal{C}_j, \mathcal{U}, \mathcal{Q}, \mathcal{C}_1, \dots, \mathcal{C}_i)$ will be generated by verifying each recipient public key by using their certificates.

- 1) First Confirms each receiver public key \mathcal{X}_i from certificate
- 2) Pick \mathcal{L} , where $0 < \mathcal{L} < n$
- 3) Split $\mathcal{L} = \mathcal{K}_a \& \mathcal{K}_b$
- 4) Compute $\mathcal{R} = \mathcal{h}(m_j, \mathcal{K}_b)$
- 5) Calculate $\mathcal{C}_j = \text{Enc}_{\mathcal{K}_a}(m_j, \mathcal{N}_r)$
- 6) Calculating the secrete key for each receiver i
 - Select \mathcal{V} where $0 < \mathcal{V} < n$
 - Computes $\mathcal{K}_i = \mathcal{V} \cdot \mathcal{X}_i$
 - Compute $\mathcal{C}_i = \mathcal{E}_{\mathcal{K}_i}(\mathcal{L})$
- 7) Computes $\mathcal{Q} = \mathcal{W}_s + \mathcal{R} \cdot \mathcal{V}$
- 8) Computes $\mathcal{U} = \mathcal{V} \cdot \mathcal{D}$
- 9) Send $(\mathcal{C}_1, \dots, \mathcal{C}_j, \mathcal{U}, \mathcal{Q}, \mathcal{C}_1, \dots, \mathcal{C}_i)$ to the group

B. Unsigncryption Phase

In the second step each recipient will receive the signcrypted text $(\mathcal{C}_1, \dots, \mathcal{C}_j, \mathcal{U}, \mathcal{Q}, \mathcal{C}_1, \dots, \mathcal{C}_i)$ through a multicast channel; and will get the plain text and will verify the sender public key \mathcal{X}_s by using his certificate.

- 1) First Confirms the public key of signcrypter \mathcal{X}_s from certificate
- 2) Calculates $\mathcal{K}_i = \mathcal{Q} \cdot \mathcal{W}_i$
- 3) Compute $\mathcal{L} = \text{Dec}_{\mathcal{K}_i}(\mathcal{C}_i)$
- 4) Split $\mathcal{L} = \mathcal{K}_a \& \mathcal{K}_b$
- 5) Calculate $m_j = \text{Dec}_{\mathcal{K}_a}(\mathcal{C}_j)$
- 6) Compute $r = \mathcal{h}(m_j, \mathcal{K}_b)$.
- 7) Computes $\mathcal{X}_s = \mathcal{Q} \cdot \mathcal{D} + r \cdot \mathcal{U}$

C. Signature Verification

The unsigncrypter verify the authenticity of received Signcrypted text as:

- Verify the public key of signcrypter \mathcal{X}_s from certificate
- Compute $\mathcal{Y} = \mathcal{Q} \cdot \mathcal{D} + r \cdot \mathcal{U}$
- Compute $\mathcal{X}_s = \mathcal{Y}$

If the last step holds, then the message is from sender otherwise the message is not sent by the sender.

V. SECURITY ANALYSIS

This phase presents the security analysis of our designed scheme. Our design scheme ensures the security requirements, for example, confidentiality, the resistance against replay attack, integrity, authenticity, non-repudiation, public verifiability, forward secrecy and unforgeability. For the validation of security requirements, we use a popular validation tool called automated validation of internet security protocols and applications (AVISPA) [51]. AVISPA is the automatic tool to validate the cryptographic schemes is either safe or un-safe. In order to find the results of developed protocol, it is essential to put in the form of HLPSSL language

according to its syntax and rules. Code written on the rules of HLPSSL language is then converted into lower level machine language through intermediate format (IF). The translation to IF is performed by the HLPSSL to IF translator. According to D. Dolev and A. Yao [52], [53], HLPSSL2IF translator checks the execution in the wisdom of given initial knowledge, every agent can construct the messages he is supposed to. AVISPA tool work with four backend [54]-[57] known as On-the-fly Model- Checker (OFMC), CL-based Attack Searcher (CL-AtSe), SAT-based Model-Checker (SATMC), and Tree-Automata-based Protocol Analyzer (TA4SP) to specify the results. Every backend have its own functionality according to their requirements. Fig. 2 shows the top down flow of AVISPA.

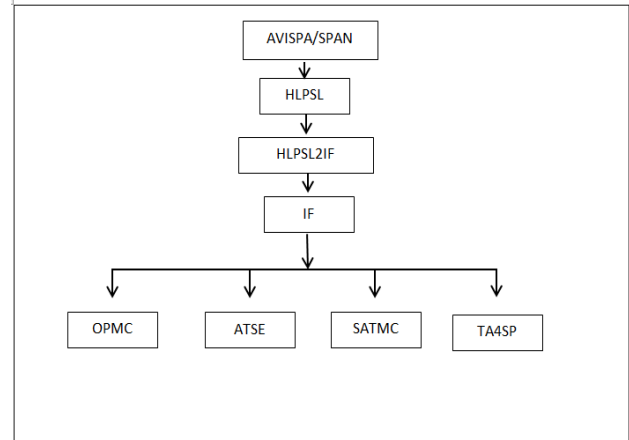


Fig. 2. Top down flow of AVISPA.

A. Confidentiality

Our method ensures the requirements of confidentiality from (1) and (2). When Alice sends message m to multi-receivers than adversary A compulsory needs secrete key \mathcal{L} to find the plain text m from cipher text \mathcal{C} . To achieve the plain text m from cipher text \mathcal{C} Adversary A needs to calculate \mathcal{C}_i from (1), to find out the \mathcal{C}_i he has to compute \mathcal{K}_i from (2). Thus, to solve \mathcal{K}_i is impossible because it is equal to calculate the hyper elliptic discrete logarithm hard problem. That's why our designed scheme ensured to obey the security requirement of confidentiality.

$$\mathcal{C}_i = \mathcal{E}(\mathcal{L}) \quad (1)$$

$$\mathcal{K}_i = \mathcal{V} \cdot \mathcal{X}_i \quad (2)$$

B. Integrity of Message

Our scheme approves that send a message is received by the original receiver and saves against any type of tampering because before sending the message hash function of the message like (3) is used. In order to achieve the integrity let us suppose that adversary A scratched the integrity by changing the cipher text \mathcal{C} as \mathcal{C}' and messages from \mathcal{C} as \mathcal{C}' then the message changes from m to m' , Therefore $m \neq m' \& r \neq \mathcal{R}$. One way hash function maintains the integrity of cipher text by denying the modification of \mathcal{C} as \mathcal{C}' and $ras\mathcal{R}$

Moreover the receiver group confirms the originality of plain text by using (4).

$$r = h(m_j, K_b) \quad (3)$$

$$Y = Q \cdot D + r \cdot U \quad (4)$$

C. Unforgeability

In order to attain the forge signature as like (5), the adversary compulsory needs W_s from (6) and V from (7). Thus to compute W_s and V is computationally hard for adversary A because it is same as to compute two time elliptic curve discrete logarithm hard problems. Hence, our scheme satisfies the security property of unforgeability

$$Q = W_s + R \cdot V \quad (5)$$

$$X_s = W_s \cdot D \quad (6)$$

$$U = V \cdot D \quad (7)$$

D. Authenticity

To achieve the authentication sender produces the signatures by using his own private key. The receiver used (6) for authentication because the sender private key associate with their public key. Furthermore, our scheme demonstrates that Authentication generated between the agents, Multi-Message-Signcrypter and Multi-Message-Unsigncrypter with the assist of nonce and encrypts the message with their secret keys K_a & K_b .

E. Non Repudiation

Our proposed scheme evidences the non-repudiation whenever a dispute occurs between sender and receiver. The Sender cannot deny what he has transmitted because third party can prove the non-repudiation using (6).

As we know that Sender sends $Q = W_s + R \cdot V$ to multi-receivers. Where W_s is the sender public key and R is already in the knowledge of the receiver. That ensures the non-repudiation property since the sender's public and private keys are associated with each other.

F. Public Verifiability

Our designed protocol provides the security property of public verifiability in case of ambiguities and disputes between agents. The designed scheme allows to verify either the message is sent by the sender or not. In case of refusal anyone can verify the message easily by performing the following steps.

- Verify the public key of signcrypter X_s from certificate
- Compute $Y = Q \cdot D + r \cdot U$
- Compute $X_s = Y$

If the last step is hold then the message from sender otherwise the message is not sent by the sender.

G. Forward Secrecy

Our designed scheme possesses the inability of an adversary A to read signcrypted messages, and recover the messages of all sessions because sender's secret key renews after every session completion. Hence, revitalization of the secret key in every session and nonce proves the goal of forward secrecy.

H. Replay attack

In our designed approach intruder may not replay old messages. Our scheme privileges the replay attack resistance by the renewal of session keys and nonce in each session. Expect that if an intruder infiltrate the message of one session, he cannot infiltrate the messages of other sessions using the same key, because the reinforcement of session key and nonce.

I. Computational Cost

In this subsection we make a comparison of our designed multi-message and multireciever with existing schemes [43], [44]. The computational cost can be computed in term most costly operations such as bilinear pairing ($\mathcal{P}\mathcal{R}$), multiplication of pairing ($\mathcal{M}\mathcal{L}$), elliptic curve multiplication ($\mathcal{H}\mathcal{m}\mathcal{L}$) and modular exponential ($e\mathcal{P}$). The Other computations such as addition, subtraction, hash and division are negligible because they need fewer computations. Table I shows the most costly operations comparison of a proposed multi-message and multireciever with existing schemes [43], [44].

TABLE I. MOST COSTLY OPERATION COMPARISON

Scheme	Multi-Signryption	Multi-Unsignryption	Total
Li [1]	$2\mathcal{P}\mathcal{R} + 3 \mathcal{M}\mathcal{L} + 1e\mathcal{P}$	$5\mathcal{P}\mathcal{R} + 3 \mathcal{M}\mathcal{L} + 2e\mathcal{P}$	$7\mathcal{P}\mathcal{R} + 6 \mathcal{M}\mathcal{L} + 3e\mathcal{P}$
Niu [2]	$2\mathcal{P}\mathcal{R} + 1 \mathcal{M}\mathcal{L} + 2e\mathcal{P}$	$4\mathcal{P}\mathcal{R} + 1\mathcal{M}\mathcal{L}$	$6\mathcal{P}\mathcal{R} + 2 \mathcal{M}\mathcal{L} + 2e\mathcal{P}$
Ours	$3 \mathcal{H}\mathcal{m}\mathcal{L}$	$3 \mathcal{H}\mathcal{m}\mathcal{L}$	$6 \mathcal{H}\mathcal{m}\mathcal{L}$

It is inspected from [58] the modular exponential consumes 1.25, pairing computation 14.31, pairing based multiplications 4.31 and elliptic curve point multiplication 0.97 milliseconds, respectively. This experiment was done by using the PC with hardware equipment's such as Intel Core i7-4510UCPU, 2.0GHz processor and 8GB of memory. The software requirement such as Windows7 Home Basic and Multi-precision Integer and Rational Arithmetic C Library (MIRACL) [59]. We assume that if elliptic curve scalar multiplication ($\mathcal{E}\mathcal{M}\mathcal{L}$) take 0.97, then hyper elliptic curve divisor multiplication ($\mathcal{H}\mathcal{m}\mathcal{L}$) take the half of elliptic curves.

Table II shows the comparisons of designing multi-message and multi-receiver with existing schemes [43], [44] in term of milliseconds. The scheme used in [43], take (129.78) milliseconds and [44] required (96.98) milliseconds for their computations. In contrast to these two schemes [43], [44], our designed multi-message and multi-receiver requires (2.88) milliseconds. Thus, it is clear from table the proposed multi-message and multi-receiver require lesser computational power.

TABLE II. COMPARISON IN MILLISECONDS

Scheme	Multi-Signryption	Multi-Unsignryption	Total
Li [43]	42.8 ms	86.98 ms	129.78 ms
Niu [44]	35.43 ms	61.55 ms	96.98 ms
Ours	1.44 ms	1.44 ms	2.88 ms

To make a reduction in computational cost among the designed multi-message and multi-receiver with existing schemes [43], [44] in term of milliseconds, we use the reduction formula [60]:

$$\frac{\text{existing approach} - \text{designed approach}}{\text{existing approach}}$$

The computational cost reduction among the designed multi-message and multi-receiver scheme from [43] is

$$\frac{129.78 - 2.88}{129.78} * 100$$

This reduces about 97.78 % and from scheme [44] is

$$\frac{96.98 - 2.88}{96.98} * 100,$$

which reduces about 97.03 %.

VI. CONCLUSION

This paper presents a lightweight multi-message and multi-receiver Heterogeneous hybrid signcryption scheme

based on the hyper elliptic curve. The proposed approach ensures the security requirements, for example, confidentiality, the resistance against replay attack, integrity, authenticity, non-repudiation, public verifiability, forward secrecy and unforgeability. Further, we validate these security requirements our scheme through a well-known security validation tool called Automated Validation of Internet Security Protocols and Applications (AVISPA). In addition, our approach has decreased in computational costs 97.03 %. To 97.78 % compare to existing schemes, this attracted the low resource devices and heterogeneous environment.

APPENDIX

In this section, we present the simulation results of our proposed scheme security requirements. We validate our proposed scheme security requirements by using a well-known security validation tool called automated validation of internet security protocols and applications (AVISPA) [51]. Fig. 3 shows that the proposed scheme is safe and Fig. 4 shows that the protocol is in working conditions.

HPLSL code

```
role
role_MultiMessageSigncrypter (MultiMessageSigncrypter:agent,MultiMessageUnsigncrypter:agent,Xs:public_key,Xi:
public_key,SND,RCV:channel(dy))
played_byMultiMessageSigncrypter
def=
  local
    State:nat,Ka:symmetric_key,Mj:text,Nr:text,Kb:symmetric_key,H1:hash_func,D:text,M1:text,V:text,Enc:ha
sh_func,L:text
  init
    State := 0
  transition
    8. State=0 /\ RCV(MultiMessageUnsigncrypter.{Nr'}_Xi) => State':=1 /\ L':=new() /\
Ka':=new() /\ Mj':=new() /\ secret(Mj',sec_2,{MultiMessageSigncrypter}) /\
witness(MultiMessageSigncrypter,MultiMessageUnsigncrypter,auth_3,Mj') /\ V':=new() /\ Kb':=new() /\
M1':=new() /\ secret(M1',sec_4,{MultiMessageUnsigncrypter}) /\ D':=new() /\
SND(MultiMessageSigncrypter.{V'.D'.{Enc(M1'.Nr')}_Ka'.{Enc(Mj'.Nr')}}_Ka'.inv(Xs).H1(Mj'.Kb').V'.{Enc(Mj'.Nr'
)}_Ka'.Enc(L')}_inv(Xs))
end role

role
role_MultiMessageUnsigncrypter (MultiMessageSigncrypter:agent,MultiMessageUnsigncrypter:agent,Xs:public_key,X
i:public_key,SND,RCV:channel(dy))
played_byMultiMessageUnsigncrypter
def=
  local
    State:nat,Ka:symmetric_key,Mj:text,Nr:text,Kb:symmetric_key,H1:hash_func,D:text,M1:text,V:text,Enc:ha
sh_func,L:text
  init
    State := 0
  transition
    8. State=0 /\ RCV(start) => State':=1 /\ Nr':=new() /\
SND(MultiMessageUnsigncrypter.{Nr'}_Xi)
    6. State=1 /\
RCV(MultiMessageSigncrypter.{V'.D'.{Enc(M1'.Nr')}_Ka'.{Enc(Mj'.Nr')}}_Ka'.inv(Xs).H1(Mj'.Kb').V'.{Enc(Mj'.Nr')}_
Ka'.Enc(L')}_inv(Xs)) => State':=2 /\ secret(Mj',sec_2,{MultiMessageSigncrypter}) /\
secret(M1',sec_4,{MultiMessageUnsigncrypter})
end role

role session1 (MultiMessageSigncrypter:agent,MultiMessageUnsigncrypter:agent,Xs:public_key,Xi:public_key)
def=
  local
    SND2,RCV2,SND1,RCV1:channel(dy)
```

```
composition
  role_MultiMessageSigncrypter (MultiMessageSigncrypter,MultiMessageUnsigncrypter,Xs,Xi,SND2,RCV2) /\
role_MultiMessageUnsigncrypter (MultiMessageSigncrypter,MultiMessageUnsigncrypter,Xs,Xi,SND1,RCV1)
end role

role session2 (MultiMessageSigncrypter:agent,MultiMessageUnsigncrypter:agent,Xs:public_key,Xi:public_key)
def=
  local
    SND1,RCV1:channel(dy)
  composition

  role_MultiMessageSigncrypter (MultiMessageSigncrypter,MultiMessageUnsigncrypter,Xs,Xi,SND1,RCV1)
end role

role environment ()
def=
  const
    hash_0:hash_func,xs:public_key,alice:agent,bob:agent,xi:public_key,const_17:agent,const_18:public_key
  ,const_16:public_key,auth_1:protocol_id,sec_2:protocol_id,auth_3:protocol_id,sec_4:protocol_id
  intruder_knowledge = {bob,alice}
  composition
    session2(i,const_17,const_18,const_16) /\ session1(alice,bob,xs,xi)
end role

goal
  authentication_on auth_1
  secrecy_of sec_2
  authentication_on auth_3
  secrecy_of sec_4
end goal

environment ()
```

VII. SIMULATION

The following Fig. 3 and 4 shows the simulation results.

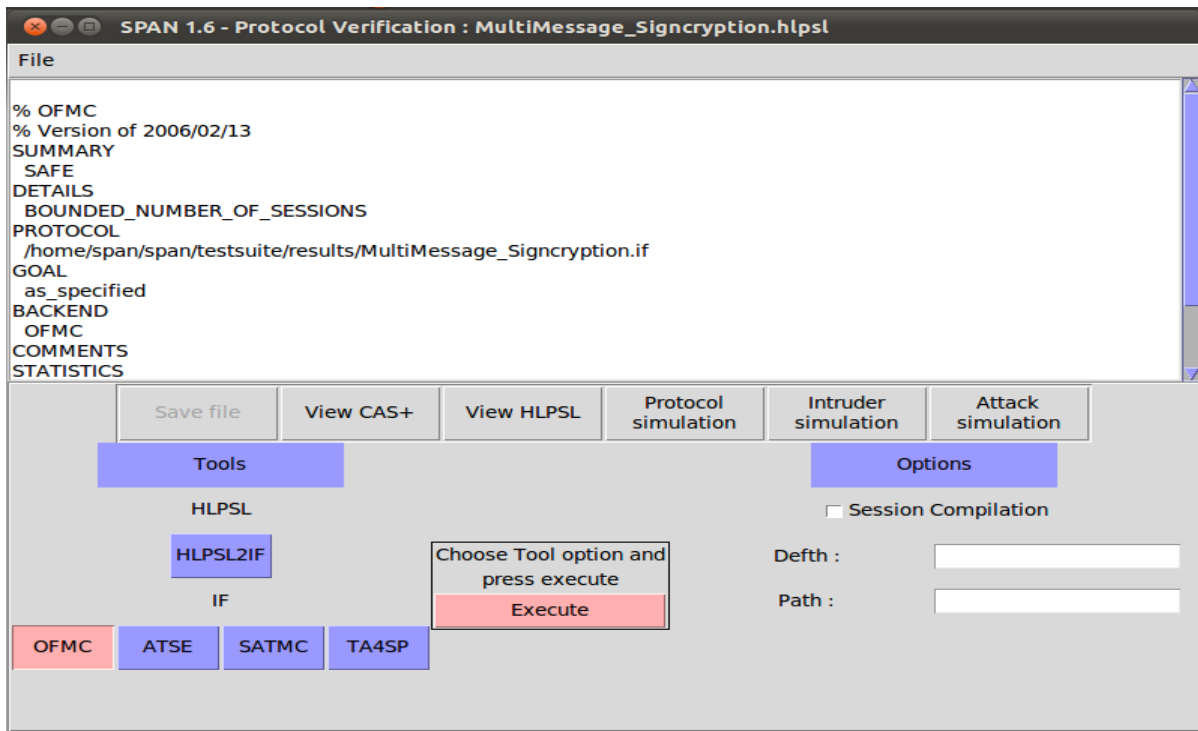


Fig. 3. Simulation results of security requirements.

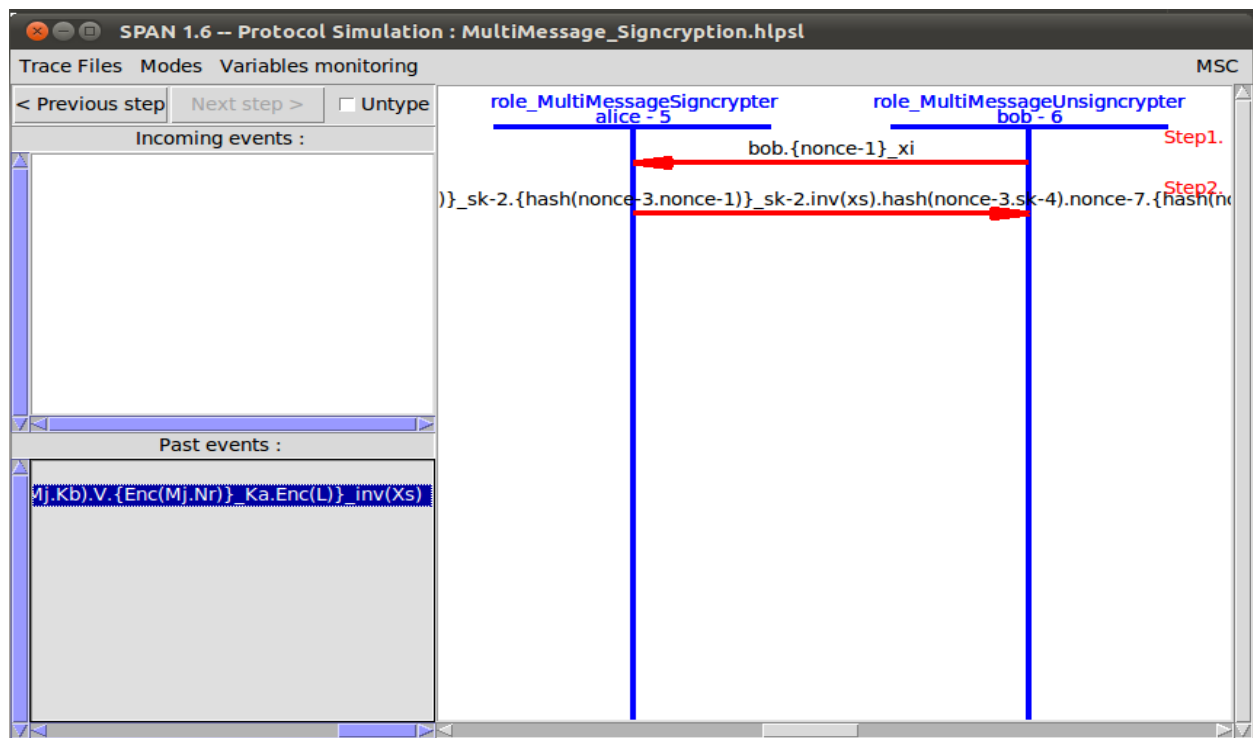


Fig. 4. Protocol in working conditions.

REFERENCES

- [1] Asad et al, "Public Verifiable Generalized Authenticated Encryption (PPG \tilde{A}) based on Hyper Elliptic Curve", *J. Appl. Environ. Biol. Sci.*, 7(12):69-73, 2017.
- [2] Arshad R, Ikram N. (2013). Elliptic curve cryptography based mutual authentication scheme for session initiation protocol. *Multimed Tools Appl* 66(2):165-178.
- [3] DegefaFB, Won D. (2013). Extended key management scheme for dynamic group in multi-cast communication. *J Conver* 4(4):7-13 7.
- [4] Diffie W, Oorschot PCV, Wiener JM (1992). Authentication and authenticated key exchanges. *Des Codes Crypt* 2:107-125
- [5] Irshad A, Sher M, Faisal MS, Ghani A, Ul Hassan M, Ashraf Ch S (2013). A secure authentication scheme for session initiation protocol by using ECC on the basis of the Tang and Liu scheme. *Security and Communication Networks* 7(8):1210-1218
- [6] Irshad A, Sher M, Rehman E, ChSA, Hassan MU, GhaniA(2013). A singleround-tripsip authentication scheme for voice over internet protocol using smart card. *Multimedia Tools Appl*:1-18
- [7] Zhang Z, Qi Q, Kumar N, Chilamkurti N, Jeong HY (2014). A secure authentication scheme with anonymity for session initiation protocol using ellipticcurve cryptography. *Multimedia Tools Appl*:1-12
- [8] Gamage C, Leiwo J, Zheng Y (1999). Encrypted message authentication by firewalls. In: *Lecture notes computer science (LNCS)*, PKC99, vol 1560. Springer-Verlag, pp 69-81
- [9] Son B, Nahm E, Kim H (2013). Voip encryption module for securing privacy. *Multimedia Tools Appl* 63(1):181-193. doi:10.1007/s11042-011-0956-1
- [10] Varalakshmi L, Florence SG (2013). An enhanced encryption algorithm for video based on multiple huffman tables. *Multimedia Tools Appl* 64(3):717-729
- [11] Yuliang Zheng.(1997). Digital signcrypton or how to achieve cost (signature & encryption) cost (signature)+ cost (encryption). In *Advances in CryptologyCRYPTO'97*, pages 165-179. Springer.
- [12] Bao F, Deng RH (1998). A signcrypton scheme with signature directly verifiable by public key. In: *Public key cryptography*. Springer, pp 55-59
- [13] Yuliang Zheng.(2001). Identification, signature and signcrypton using high order residues modulo an rsa composite. In *Public Key Cryptography*, pages 48-63. Springer.
- [14] . John Malone-Lee and Wenbo Mao.(2003). Two birds one stone: signcrypton using rsa. In *Topics in Cryptology CT-RSA* , pages 211-226. Springer.
- [15] Joonsang Baek, Ron Steinfeld, and Yuliang Zheng. (2002). Formal proofs for the security of signcrypton. In *Public Key Cryptography*, pages 80-98. Springer.
- [16] Sharma G, Bala S, Verma AK (2013). An identity-based ring signcrypton scheme. In: *IT convergence and security* . Springer, 151-157
- [17] Zheng Y, Imai H (1998). How to construct efficient signcrypton schemes on ellipticcurves. *Inf Process Lett* 68(5):227-233
- [18] Hwang RJ, Lai CH, Su FF (2005). An efficient signcrypton scheme with forward secrecy based on ellipticcurve. *Appl Math Comput* 167(2):870-881
- [19] Toorani M, Beheshti AA. (2010).An elliptic curve-based signcrypton scheme with forward secrecy. arXiv:1005.1856
- [20] Nizamuddin, Ch SA, Amin N. (2011). Signcrypton schemes with forward secrecy based on hyperelliptic curve cryptosystem. In: *High capacity optical networks and enabling technologies (HONET)*, 2011, pp 244-247. doi:10.1109/HONET.6149826
- [21] Nizamuddin, Ch SA, Nasar W, Javaid Q (2011). Efficient signcrypton schemes based on hyperelliptic curve cryptosystem. In: *7th international conference on emerging technologies (ICET)*, pp 1-4
- [22] Yiliang Han and Xiaoyuan Yang. Ecgsc. (2006). Elliptic curve based generalized signcrypton scheme. *IACR Cryptology ePrint Archive*, 2006:126.
- [23] Lal, S.; Kushwah, P. (2008). ID Based Generalized Signcrypton. *Cryptology ePrint Archive*, Report 2008/084.
- [24] Jindan Zhang and Xu an Wang.(2009). Formal security proof for generalized signcrypton. In *E-Business and Information System Security, EBISS'09*. International Conference on, pages 1-5. IEEE.
- [25] HF Ji, WB Han, and Long Zhao.(2010). Identity-based generalized signcrypton in standard model. *Appl. Res. Comput*, 27(10):3851-3854.

- [26] Zhang Chuanrong, Chi Long, and Zhang Yuqing. (2010). Secure and efficient generalized signcryption scheme based on a short ecDSA. In *Intelligent Information Hiding and Multimedia Signal Processing (IHMSP)*, 2010 Sixth International Conference on, pages 466–469. IEEE.
- [27] Yu, G.; Ma, X.; Shen, Y.; Han, W. (2010). Provable secure identity based generalized signcryption scheme. *Theor. Comput. Sci.*, 411, 3614–3624.
- [28] Gang Yu, Xiaoxiao Ma, Yong Shen, and Wenbao Han. (2010). Provable secure identity based generalized signcryption scheme. *Theoretical Computer Science*, 411(40):3614–3624.
- [29] Prashant Kushwah and Sunder Lal. (2011). An efficient identity based generalized signcryption scheme. *Theoretical Computer Science*, 412(45):6382–6389.
- [30] Shen et al. (2017). Identity Based Generalized Signcryption Scheme in the Standard Model. *Entropy*, 19, 121; doi:10.3390/e19030121.
- [31] Shehzad et al. (2014). An efficient signcryption scheme with forward secrecy and public verifiability based on hyper elliptic curve cryptography. *Multimed Tools Appl* DOI 10.1007/s11042-014-2283-9.
- [32] Anwar et al., “Multi-Receiver Signcryption Based on Hyper Elliptic Curve Crypto System”, *J. Appl. Environ. Biol. Sci.*, 7(12)194-200, 2017
- [33] Y. Zheng, H. Imai, 1998. How to construct efficient signcryption schemes on elliptic Curves: *Intl. J. Information Processing Letters* 68(5): 227-233.
- [34] Zhang Y, Zhang L, Zhang Y, Wang H, Wang C. CLPKC-to-TPKI heterogeneous signcryption scheme with anonymity. *Acta Electronica Sinica*. 2016; 44(10):2432–2439.
- [35] Li F, Han Y, Jin C. Practical access control for sensor networks in the context of the Internet of Things. *Computer Communications*. 2016; 89(9):154–164. <https://doi.org/10.1016/j.comcom.2016.03.007>
- [36] Li F, Han Y, Jin C. Practical signcryption for secure communication of wireless sensor networks. *Wireless Personal Communications*. 2016; 89(4):1391–1412. <https://doi.org/10.1007/s11277-016-3327-4>
- [37] Li Y, Wang C, Zhang Y, Niu S. Privacy-preserving multi-receiver signcryption scheme for heterogeneous systems. *Security and Communication Networks*. 2016; 9(17):4574–4584. <https://doi.org/10.1002/sec.1650>
- [38] Dent AW. Hybrid signcryption schemes with outsider security. In: *Information Security-ISC 2005*, LNCS 3650. Springer-Verlag; 2005. p. 203–217.
- [39] Dent AW. Hybrid signcryption schemes with insider security. In: *Information Security and Privacy ACISP 2005*, LNCS 3574. Springer-Verlag; 2005. p. 253–266.
- [40] H. M. Elkamchouchi, A. M. Emarah, and E. A. A. Hagra, iPublic Key Multi-Message Signcryption (PK-MMS) scheme for secure communication systems, in *Proceedings - CNSR 2007: Fifth Annual Conference on Communication Networks and Services Research*, 2007, pp. 329–334.
- [41] H. M. Elkamchouchi, A. A. M. Emarah, and E. A. Hagra, iA new efficient public key multi-message multi-recipient signcryption (PK-MM-MRS) scheme for provable secure communications, i *ICCES'07 - 2007 Int. Conf. Comput. Eng. Syst.*, pp. 89–94, 2007.
- [42] H. Elkamchouchi, M. Nasr, and R. Ismail, iA new efficient multiple broadcasters signcryption scheme (MBSS) for secure distributed networks, i *Proc. 5th Int. Conf. Netw. Serv. ICNS 2009*, pp. 204–209, 2009.
- [43] Li Y, Wang C, Zhang Y, Niu S. Privacy-preserving multi-receiver signcryption scheme for heterogeneous systems. *Security and Communication Networks*. 2016; 9(17):4574–4584. <https://doi.org/10.1002/sec.1650>
- [44] Niu S, Niu L, Yang X, Wang C, Jia X (2017) Heterogeneous hybrid signcryption for multi-message and multi-receiver. *PLoS ONE* 12(9): e0184407. <https://doi.org/10.1371/journal.pone.0184407>
- [45] N. Koblitz, Hyperelliptic cryptosystems, *Journal of Cryptology*, Vol. 1, 1989, 139-150.
- [46] N. Koblitz, Elliptic curve cryptosystems, *Mathematics of Computation*, Vol. 48, 1987, 203-209.
- [47] T. Wollinger. Software and Hardware Implementation of Hyperelliptic Curve Cryptosystem. Dissertation for the Degree of Doctor-Ingenieur. - Bochum, Germany, 2004. 201p.
- [48] Pelzl, T. Wollinger, J. Guajardo, C. Paar. Hyperelliptic Curve Cryptosystems: Closing the Performance Gap to Elliptic Curves, 2003, 15 p., <http://eprint.iacr.org/026.pdf>.
- [49] Pelzl, T. Wollinger, C. Paar. High Performance Arithmetic for Hyperelliptic Curve Cryptosystems of Genus Two., 2004, 12 p., <http://eprint.iacr.org/212.pdf>.
- [50] D. Mumford. *Tata Lectures on Theta II*. In *Prog. Math.*, volume 43. Birkhauser, 1984.
- [51] D. Dolev and A. Yao, “On the Security of Public-Key Protocols”, *IEEE Transactions on Information Theory*, 2(29), 1983. <http://ieeexplore.ieee.org/document/1056650/>
- [52] D. Basin, S. Modersheim, and L. Vigan, “An On-The-Fly Model-Checker for Security Protocol Analysis”, In *Proceedings of ESORICS'03*, LNCS 2808, pages 253–270. Springer-Verlag, 2003. https://link.springer.com/chapter/10.1007/978-3-540-39650-5_15
- [53] J. Clark and J. Jacob, “A Survey of Authentication Protocol Literature”, Version 1.0, 17. Nov. 1997. www.cs.york.ac.uk/~jac/papers/drareview.ps.gz.
- [54] B. Donovan, P. Norris, and G. Lowe, “Analyzing a Library of Security Protocols using Casper and FDR”, In *Proceedings of the FLOC'99 Workshop on Formal Methods and Security Protocols (FMSP'99)*, 1999. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.21.7256>
- [55] Y. Chevalier and L. Vigneron “Automated Unbounded Verification of Security Protocols”, In *Proc. CAV'02*, LNCS 2404. Springer, 2002. https://link.springer.com/chapter/10.1007/3-540-45657-0_24
- [56] Armando and L. Compagna, “SATMC: a SAT-based Model Checker for Security Protocols”, In *Proc. JELIA'04*, LNAI 3229. Springer, 2004. https://link.springer.com/chapter/10.1007/978-3-540-30227-8_68
- [57] Y. Boichut, P.-C. Heam, O. Kouchnarenko and F. Oehl, “Improvements on the Genet And Klay Technique to Automatically Verify Security Protocols”, In *Proc. AVIS'04*, ENTCS. https://www.researchgate.net/publication/246435265_Improvements_on_the_Genet_and_Klay_technique_to_automatically_verify_security_protocols
- [58] Caixue Zhou et al, “Certificateless Key-Insulated Generalized Signcryption Scheme without Bilinear Pairings”, *Security and Communication Networks* Volume 2017, Article ID 8405879, 17 pages <https://doi.org/10.1155/2017/8405879>
- [59] Shamus Software Ltd. Miracle library, <http://github.com/miracl/MIRACL>.
- [60] Shehzad et al. (2012). Public Verifiable Signcryption Schemes with Forward Secrecy Based on Hyper elliptic Curve Cryptosystem. *ICISTM 2012*, CCIS 285, pp. 135–142.

A Chatbot for Automatic Processing of Learner Concerns in an Online Learning Platform

Mamadou BAKOUAN¹

Laboratoire de Recherche en Informatique et
Télécommunication (LARIT)
Ecole Doctorale Polytechnique de l'INP-HB
Yamoussoukro, Côte d'Ivoire

Beman Hamidja KAMAGATE²

Laboratoire de Recherche en Informatique et
Télécommunication (LARIT)
Ecole Supérieure Africaine des TIC - ESATIC
Abidjan, Côte d'Ivoire

Tiemoman KONE³

Laboratoire d'Informatique, Signaux et Télécommunications
Institut de Recherche Mathématiques IRMA/UFHB
Abidjan, Côte d'Ivoire

Souleymane OUMTANAGA⁴, Michel BABRI⁵

Laboratoire de Recherche en Informatique et
Télécommunication (LARIT)
INP-HB
Yamoussoukro, Côte d'Ivoire

Abstract—In this article, we present a chatbot model that can automatically respond to learners' concerns on an online training platform. The proposed chatbot model is based on an adaptation of the similarity of Dice to understand the concerns of learners. The first phase of this approach allows selecting the pre-established concerns that the teacher has in a knowledge base which are closest to those posed by the learner. The second phase consists of selecting among these k most appropriate concerns based on a measure of similarity built on the concept of domain keywords. The experimentation of the prototype of this chatbot makes it possible to find the adequate answers. In the case, where the question refers to a question from the teacher, the learner is asked if the question identified is the one he was referring to. If he answers in the affirmative, the instructions associated with his request are sent to him. If not, the learner's concern is sent to the human tutor. The hybridization of this chatbot with the human agent comes to enrich the initial knowledge base of the chatbot. The results obtained with the concept based on the keywords of the domain are encouraging. The learner's comprehension rate is above 50% when applying the concept of domain keywords while the measure of Dice is below 50%.

Keywords—Metadata; ontologies; semantic similarity; natural language; semantic web; chatbot

I. INTRODUCTION

Chatbot are interactive virtual characters whose mission is to provide assistance to people in high-profile environments. Previous research has shown that this technology seems to have a positive influence on learning [1]. In addition, the presence of interactive virtual agents, also called Chatbot, taking on the role of guardian [2], seems to have positive effects on student engagement [3] and on the effectiveness of teaching [4]. In the education system in Côte d'Ivoire, the number of graduates is growing steadily, without a corresponding increase in the capacity of higher education institutions [5]. To face this situation, the government has opted for the integration of new technologies (ICT) in education through the interconnection of universities and public schools in Côte d'Ivoire [6]. This project should make it

possible to unclog university lecture halls by relying on distance learning and facilitate access to teaching resources. However, since 2015 the infrastructures of the e-Education project are not operational.

In this dynamic, the State uses e-learning through the creation of Université Virtuelle de Côte d'Ivoire (UVCI) [7]. One of UVCI's missions is to develop distance education in Côte d'Ivoire. This type of teaching is based on a set of platforms to facilitate access to learning resources for learners. In the pedagogical model of the UVCI, the human tutor plays the role of framer. It ensures the educational follow-up of the training. However, the response time of the physical tutor is low and the high number of students per physical tutor degrades the quality of the training. This sometimes gives rise to the feeling of abandonment in some students.

To remedy this, we offer a chatbot that helps to take care of students' concerns on a permanent basis. It is about lightening the task of teachers and tutors while contributing to the framing and effective management of student concerns. In the next section, we will describe the role of metadata and ontologies in how chatbot work. Then we will discuss the mechanism used by the chatbot to understand the sentences. Finally, we will see the experimentation of the prototype of the chatbot and the results.

II. LITERATURE REVIEW

Information systems have to evolve with certainty, their agility is a major requirement. Software architectures must therefore promote real flexibility and reusability to adapt to change. New software architectures have brought a real ability of an architecture to evolve in order to integrate some changes response to the complex need of integration of information systems. It is particularly in this context that the new generation of formal metadata system technologies and the semantic web, derived from the Service-Oriented Architectures paradigm, aims to respond in a relevant way to the question of interoperability related to the agility of chatbot systems.

A. Semantic Web Technology

The term semantic Web, ascribed to Tim Berners-Lee [8] in the W3C, first refers to the vision of the Web of tomorrow as a vast space of exchange of resources between human beings and machines allowing exploitation, qualitatively superior, large volumes of information and varied services. Virtual space, it should see, unlike the one we know today, the users discharged of a good part of their tasks of research, construction and combination of the results, thanks to the increased capacities of the machines to access the resources and to reason with them. The semantic web is structured in layers. These layers correspond to different categories of formalisms grouped into three levels. This is the naming / addressing level, the syntactic level and the semantic level. The semantic web respects an architecture (see Fig. 1). This figure represents the structure of semantic web components.

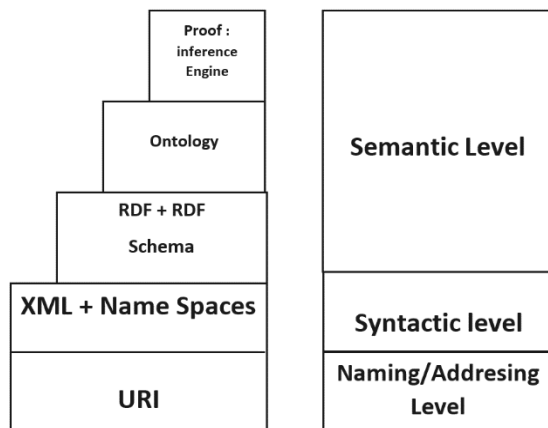


Fig. 1. Semantic web architecture [8].

Most of the languages standardized by W3C as part of the Semantic Web are XML dialects, such as RDF and RDFS. The RDFS provides basic elements for defining ontologies or vocabularies for structuring RDF resources. SPARQL is a query language for RDF. Like SQL for relational databases or Xpath and XQuery for XML documents, this language is used to retrieve information from RDF documents. The construction of ontologies and metadata requires consensus in order to avoid lexical ambiguities due to hyperonymies and polysemias. A metadata is literally a datum on a datum. That is, a structured set of information describing any resource [9]. The principle of metadata is about association a number of fields with resources for which values are assigned to each. These values can be given in a free format, as they can also conform to well-defined data formats. The operation consists of considering tags that are introduced in the files or in the appropriate programming languages. Tags have the effect of improving the efficiency of information searches compared to full-text searches. It is important to note that tagged digital resources carry with them their own metadata and those, when downloaded, copied, replicated, transmitted by email. This approach promotes interoperability for better exploitability of digital resources. Several standardization organizations have proposed and published metadata schemas that could be used by as many people as possible. We will analyze the different metadata schemas in the next section.

- LOM (Learning Object Metadata) [10]
- EAD (Encoded Archival Description) [11]
- Dublin Core [11]

The concept of metadata requires the definition of a kernel of standard and context-dependent information. This can make it difficult to exploit metadata in a learning model. To optimize this concept, the metadata schema used in learning models is enhanced by technology. Indeed, metadata is associated with domains of knowledge that can be conceptualized in ontologies.

Ontologies represent a source of very reliable and structured knowledge. For this reason and thanks to the initiatives of the Semantic Web, which brought the creation of thousands of domain ontologies, ontologies have been widely exploited in knowledge-based systems, and more specifically, for the calculation of semantic similarity. An ontology is formally defined as a pair (O, Lex) where O is an abstract ontology and Lex is a lexicon for O [12]. Let L be a logical language having formal semantics in which inference rules can be expressed. An abstract ontology is a structure $O = (C, \leq_C, R, \sigma, \leq_R, I, R)$ consists of:

Two disjoint sets C and R whose elements are respectively called Concepts and Relations;

A partial order \leq_C on C , called hierarchy of concepts or taxonomy;

A function $\sigma: C \times C$ called signature;

A partial order \leq_R on R , called hierarchy of relations where $r_1 \leq_R r_2$ implies $\sigma(r_1) \leq_C \sigma(r_2)$ with $r_1, r_2 \in R$.

A set I, R of inference rules expressed in the logical language L ;

The dom function: $R \rightarrow C$ with $\text{dom}(r) = \Pi_1(\sigma(r))$ returns the domain of r ;

The range function: $R \rightarrow C$ with $\text{rank}(r) = \Pi_2(\sigma(r))$ returns the scale of values of r ;

A lexicon for an abstract ontology $O = (C, \leq_C, R, \sigma, \leq_R, I, R)$ is a structure $Lex = (S_C, S_R, Re f_C, Re f_R)$ which consists on:

Two sets S_C and S_R whose elements are called signs, respectively for concepts and relations;

- Two relations $Re f_C \subseteq S_C \times C$ and $Re f_R \subseteq S_R \times R$, called assigning lexical references respectively for concepts and relationships;

From $Re f_C$ we define $\forall s \in S_C, Re f_C(s) = c \in C | (s, c) \in Re f_C$ and $Re f_C^{-1}(s) = s \in C | (s, c) \in Re f_C$

- From $Re f_R$ we define $\forall s \in S_R, Re f_R(s) = r \in R | (s, r) \in Re f_R$ and $Re f_R^{-1}(s) = s \in R | (s, r) \in Re f_R$

There are ontologies in different fields that support the design of learning systems including DogOnt ontology, SOUPA, CoBrA, CoDAMoS, etc. [13], [14].

- SOUPA (Standard Ontology for Ubiquitous and Pervasive Applications).
- DogOnt (Ontology Modeling for Intelligent Domestic Environments).
- CoDAMoS (Context-Driven Adaptation of Mobile Services).
- CoBrA (cobra-have overview).

In the literature several languages have been used for the description of ontologies. These languages include the eXtensible Markup Language [13], the Resource Description Framework (RDF) [15], the DAML + OIL (Darpa Modeling Language of Ontology + Ontology Inference Layer) [16] and OWL (Ontology Web Language) [14]. These languages offer different levels of expressiveness. Making yourself available to answer questions about distance learning activities related to a training module followed in a teaching platform are non-obvious tasks especially if the number of learners is important. Hence our idea, to integrate a chatbot whose role is collaboration and cooperation with the human tutor.

B. ChatBot

A Chatbot is a computer program capable of simulating a conversation with one or more users by voice or text exchange. Indeed, he plays the role of an assistant who aims to answer the questions put to him, while imitating human behavior [17]. The operating principle of a virtual guardian agent goes through three stages:

- The learner first sends questions that he would like to address to the agent.
- The agent receives the learner's question.
- He analyzes the question by consulting his knowledge base and finally provides an answer to the questions asked by the learner.

We could classify chatbot into two main categories:

- Virtual recommendation agents: This agent makes proposals to users in a virtual environment [18].
- Feedbacks chatbot: This agent makes feedbacks after performing an activity in a virtual environment [19].

Sassi researchers [18] propose a virtual recommendation agent that assists a user in his daily tasks, without any explicit request from the user. This agent aims to assist the user in his daily tasks thanks to his ability to perceive the state of the environment and to interact effectively according to the needs of the user.

Joanna's work [19] focused on the chatbot of Feedbacks. They provide a chatbot that can provide feedback to users after performing an activity in a virtual environment. Chatbot feedback and interpretation of user feedback is based on knowledge of the virtual environment. After analyzing the different works, we found that the proposed chatbot do not take into account the online learning environment. In the next section, we present some approaches for comparing texts. We will speak later of similarity between texts. The presented

approaches have been selected to best respond to the context. Thus, this document does not claim to give an exhaustive list of all the existing methods but tries to give an overview of the most used methods in the context of our study. In the next section, we will describe these different notions of similarity measure in sentences.

C. Similarity Measures Between Sentences

In automatic language processing, similarity measurement plays an important role and is one of the fundamental tasks. The automatic understanding of a sentence requires from the web agent different types of abilities: recognizing words and associating them with lexical information (morphological analysis); structure the sentence with a grammar (parsing), understand the sentence with semantic rules (semantic analysis) and take into account the context (pragmatic analysis). Huang [20] has shown that the performances of syntactic similarity based on the Jaccard index and the Dice index are very close and that they are significantly better than those of the Euclidean distance and the Levenshtein distance. The distance from Levenshtein is widely used in linguistics and bioinformatics as well as for the recognition of text blocks. Unfortunately, the computation time (complexity) is when applied to two sequences of approximately the same size. This is an obstacle in many practical applications.

In Christine's work [21], she proposes a method for measuring the semantic similarity between strings of characters. This method is based on the combination of Levenshtein's distance and Jaccard's index. This method has shortcomings when the strings correspond to names composed of several words. In addition, it requires a perfect match between each string in the two sets of strings. Thus, Hai-Hieu Vu and Jeanne Villaneau [22] proposed another method for measuring the semantic similarity between sentences that uses Wikipedia as the only linguistic resource. This method is based on a vector representation; it uses a random indexing to reduce the size of the manipulated spaces. Hai's method does not return a precise answer to the user. It returns to the user a Wikipedia article containing the elements of answer to his concern. The user is led to analyze this article in order to find an answer to his concern. Goutam Majumder and Partha Pakray [23] propose a method for calculating the semantic similarity between sentences based on the WordNet taxonomy. It allows to index, classify and put in relation the semantic and lexical contents of the English language. This method is not adapted to our context.

The similarity methods proposed in the research works are based on the TF-IDF method. TF-IDF (term frequency-inverse document frequency) is a weighting method used for finding information in the corpus. The TF-IDF method requires preprocessing of the corpus to determine the discriminating power of each word. While this pretreatment uses significant resources and lengthens the query processing time. The proposed chatbot model is an adaptation of the Dice measure based on the concept of domain keywords to understand the concerns of learners. The hybridization of the chatbot with the human agent comes to enricher the initial knowledge base of the chatbot.

III. MECHANISM USED BY THE CHATBOT TO UNDERSTAND SENTENCES

We propose a measure adaptation of Dice to calculate the similarity between sentences. This approach is based on the Dice index and the measure of similarity of the keywords of the domain. We will discuss the principle of the algorithm and the process of calculating the similarity between sentences.

A. Principle of the Algorithm

- The learner sends a question to the chatbot.
- The chatbot receives the learner's question.
- The chatbot analyzes the learner's question.
- Cleanup (Remove StopWord).
- Lemmatization (Convection of words in lemma).
- Selection of k questions (Comparison of words in common and select questions closest to the learner's question).
- Similarity based on domain words (Search among selected questions, one that is semantically close to the learner's question).
- Proposition of the question semantically close to the learner.
- The learner should confirm that the proposal corresponds to his / her concern or not.
- If the learner answers with "NO", his question is returned to a human agent.
- If the learner answers with "YES", the chatbot provides the answer to the learner's question.

B. Calculation of the Similarity between Sentences

The calculation of the similarity between sentences has been implemented by performing the following steps:

Phrase Labeling: This step deals with all of the sentences in the corpus (see Fig. 2) and converts each of their terms into lemmas. Lemmatization consists of finding the root of the bent verbs and bringing the plural and / or feminine words back to the singular masculine form (see Fig. 3).

Selection of k questions: A measure of similarity to select the k questions closest to the learners' preoccupation (see Fig. 4). This similarity approach is based on the measure of Dice. The measure of Dice calculates the similarity between two sentences Q_E and Q_S based on the number of terms common to Q_E and Q_S (see Fig. 4).

$$sim_{dice}(Q_E, Q_S) = \frac{2N_C}{S_E + S_S} \quad (1)$$

Q_E represents all the terms of the student's question.

S_E represents the number of terms after the lemmatization of the student's question.

Q_S represents all the terms of the teacher's question.

S_S is the number of terms after the lemmatization of the teacher question.

N_C is the number of terms common to Q_E et Q_S

Practical case of similarity of the Dice index between Q_E and Q_S :

Q_S "Example of question proposed by the teacher": Why I cannot read other students' posts in the forum?

Q_E "Example of student question": Unable to read messages from my fellow students in the forum.

Step 1: Cleaning the stopwords

Step 2: Converting the terms to lemma

Step 3: Analysis of terms common to Q_E and Q_S

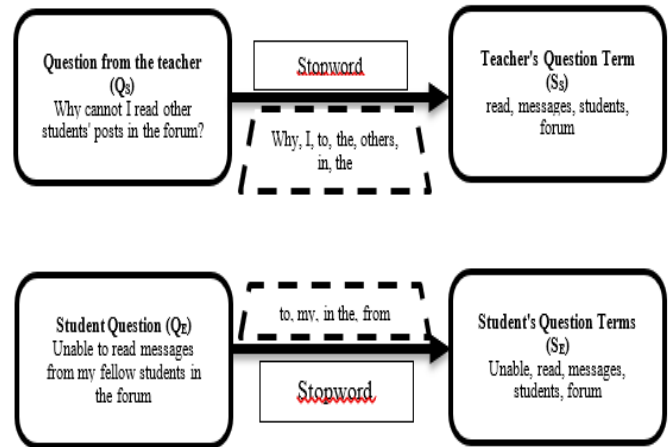


Fig. 2. Stopword cleaning process.

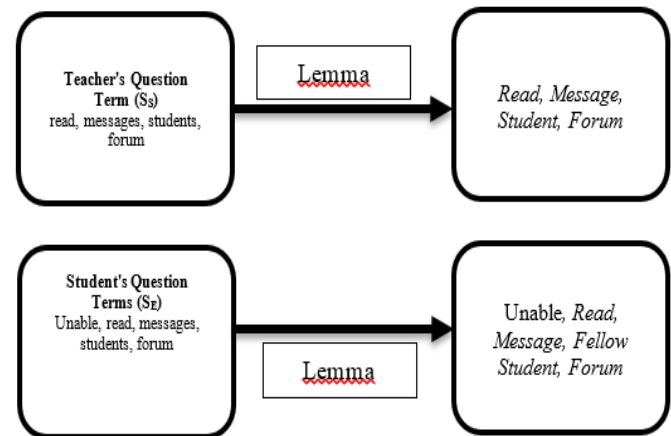


Fig. 3. Term lemmatization process.

The analysis of the terms common to Q_E and Q_S makes it possible to retain the k Q_S questions close to the Q_E questions. Then, a method of similarity based on the keywords of the domain allows to retain the Q_S question closest to the Q_E question (see Fig. 5).

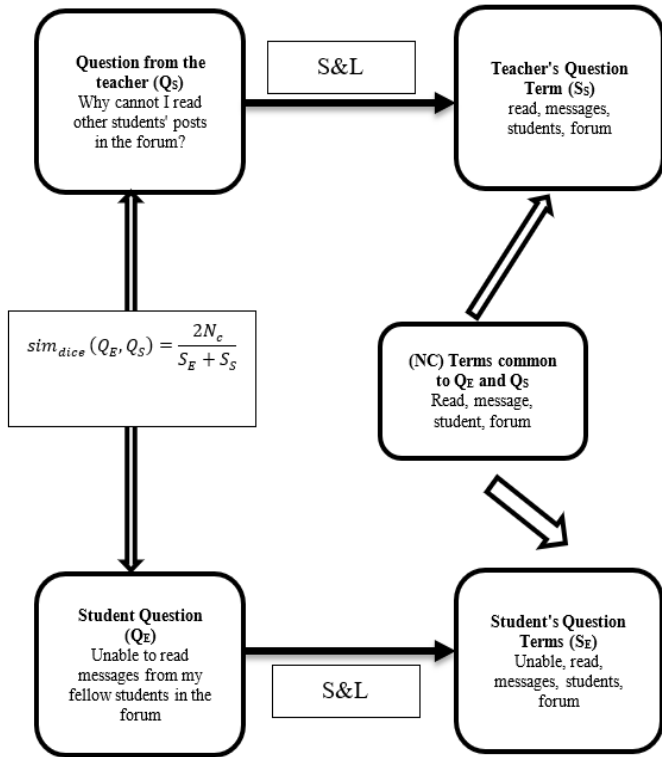


Fig. 4. Common terms analysis process.

Similarity measure built on the concept of domain keywords: this method consists in finding the Q_S question closest to the Q_E question by integrating the principle of the keywords of the domain.

MC represents all the keywords of the teacher's course.

Q_E represents all the terms of the student's question.

S_E represents the number of terms after the lemmatization of the student's question that belong to faith in Q_E and MC.

Q_S represents all the terms of the teacher's sth question with $\forall s \in \{1 \dots k\}$.

S_S represents the number of terms after the lemmatization of the teacher sth question that belong to the faith in Q_S and MC. T_S : represents all terms in common to S_S and S_E

$$S_E = \{t_i \in Q_E / t_i \in MC\} \quad (2)$$

$$S_E = Q_E \cap MC \quad (3)$$

$$S_S = \{t_i \in Q_S / t_i \in MC\} \quad (4)$$

$$S_S = Q_S \cap MC \quad (5)$$

$$T_S = |S_E \cap S_S| \quad (6)$$

$\forall s \in \{1 \dots k\} Q_E \cong Q_S$ if T_S is the maximum

Practical case of similarity between Q_E and Q_S :

Domain Keyword (s): Course, Email, Feedback, Duty, Grade, Medium, Test, Tool, Communication, Forum, Message, Discussion, Student, Publication

Question from the teacher:

Q1: Why cannot I read messages from other students in the forum?

Answer: In a Question & Answer forum, you must first contribute to the forum by submitting a contribution before having access to the messages of other students.

Q2: How can I keep up with the news of the forum (s) to which I subscribe?

Answer: In the "My classes" workspace, a message informs you of additions to the forum (s).

Q3: How to visualize all my publications in the forums?

Answer: In the tab "My page - My profile - Forum posts", you can view all contributions to forums, discussions launched, or answers given.

Question of the student:

Q1: Cannot read the messages of my fellow students in the forum

Q2: I cannot read messages from students in my group

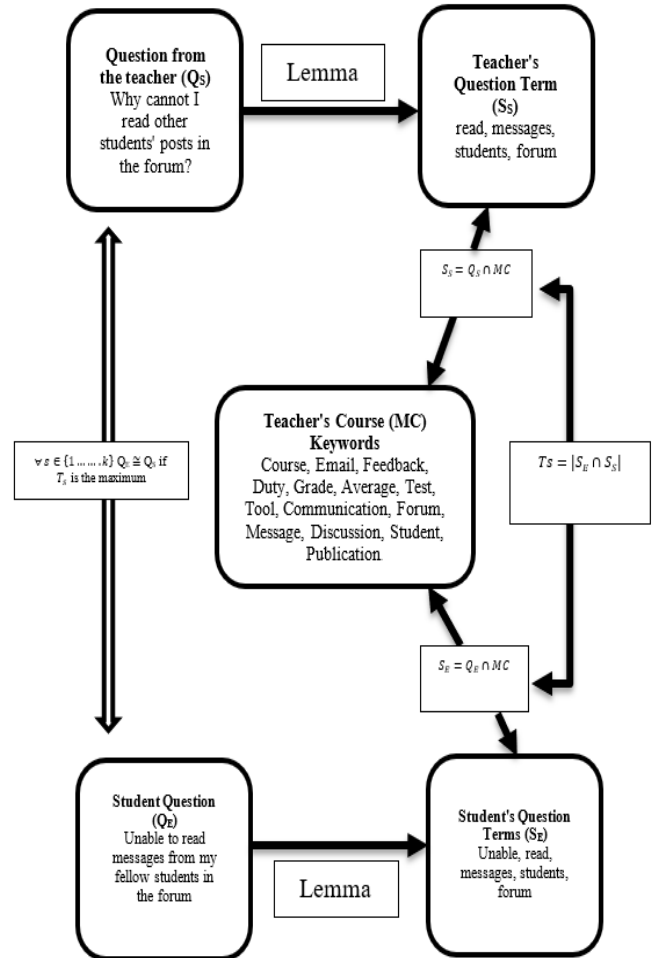


Fig. 5. Similarity measure based on domain keywords.

IV. GLOBAL OPERATION OF THE CHATBOT

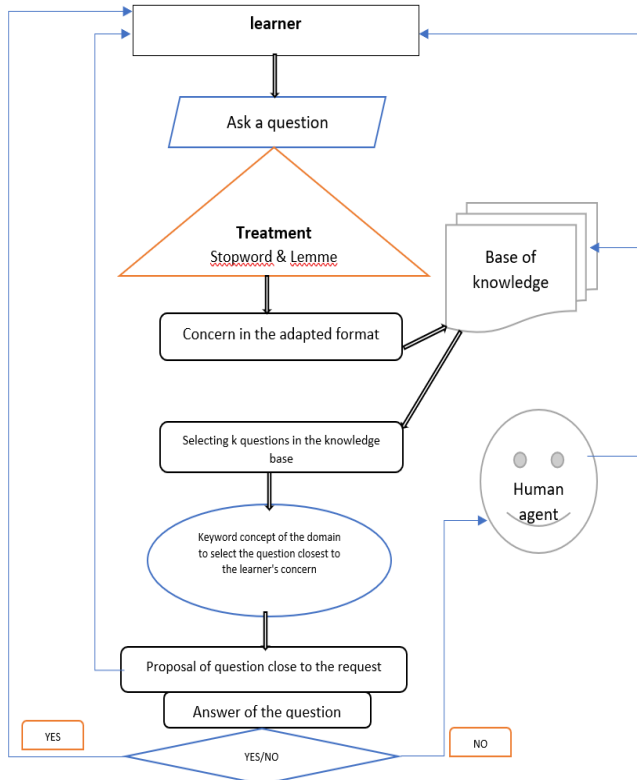


Fig. 6. Principle of function of the Chatbot.

Learner: The learner module allows to submit the concerns.

Ask a question: The module asks a question, it is the raw text of the learner's question.

Treatment: The treatment module takes into account the cleaning of the stopword and the lemmatization of the terms of the learner's preoccupation.

Lemmatization refers to the lexical analysis of the content of a text grouping the words of the same family. Each of the words of a content is thus reduced to an entity called lemma (canonical form).

Stopword: A stopword is a non-significant word in a text. It is opposed to a full word. The meaning of a word is evaluated from its distribution (in the statistical sense) in a collection of texts.

Concern in the adapted format: The learner's concern is converted into a format that allows the chatbot to understand it.

Knowledge base: The knowledge base brings together knowledge specific to the field of Université Virtuelle de Côte d'Ivoire, in a form usable by the chatbot. It contains rules that allow structuring of the data.

Selection of k questions in the knowledge base: This module allows you to browse the knowledge base in search of the teacher's questions that are close to the learner's preoccupation. We retain k questions that have high score and close to the learner's concern.

Keyword concept of the domain to select the question closest to the learner's preoccupation: Once the k closest questions are selected, we apply a concept based on the domain keyword principle. This approach selects the question of the teacher closest to the learner's concern.

Proposition of the question closest to the request: This module makes it possible to propose the question of the teacher to the learner. The learner is amenable to rewrite, if he answers by YES then the answer associated with the question of the teacher is returned to the learner. If the learner answers by NO then his concern is sent to a human agent for treatment.

Human Agent: When the chatbot does not have the answer to the learner's concern, the learner's concern is sent to the human agent who analyzes it and returns the appropriate answers. The responses of the human agent enrich the knowledge base.

Fig. 6 shows the overall operation of the chatbot and hybridization in the human agent to enrich the knowledge base.

V. EXPERIMENTATION

The experiment concerns the global operation of the prototype of the chatbot. The learner is connected to his workspace (Fig. 7) and he submits his concern to the chatbot. When he clicks the Enter key or the Submit button in the window, his concern is then converted into a language query (Fig. 8). Treatments are successively carried out as the suppression of stopwords then a lemmatization of the remaining terms.

As a result of these treatments, the query obtained is analyzed to obtain questions from the teacher close to the learner's question. Then, the treatment carried out makes it possible to find the question of the teacher closest to the question of the student. Once a question is selected, it is sent to the view to be returned to the learner (Fig. 9). Then the answer to this question will be analyzed and will return the instructions according to the following answer:

Yes: the appropriate instructions will be sent (Fig. 9)

No: the learner's concern is sent to a human tutor for analysis (Fig. 10)

The experiment is performed with the following hardware and tools: It is a Corei7 processor computer, 12GB RAM and 1TB hard drive, the object-oriented PHP programming language and a Database Management System MYSQL.



Fig. 7. Window allowing the learner to submit his concern to the chatbot.

Fig. 7 represents the window of dialogue with the learner. First a message of good is given then the student in the field seizes in the order to submit his concern to ChatBot. Finally, the learner clicks the submit button to validate his concern. As shown in Fig. 8, once the learner submits his preoccupation. This triggers the process of processing the quest. After treatment, the chatbot offers a response element to the learner. The learner has the opportunity to confirm the proposal of the chatbot. An answer item is returned to the learner based on the confirmation. If the learner answers by YES, he receives the answer adapted to his concern (Fig. 9) and if he answers by NO, the concern is sent to the human agent to have the adapted answer (Fig. 10).

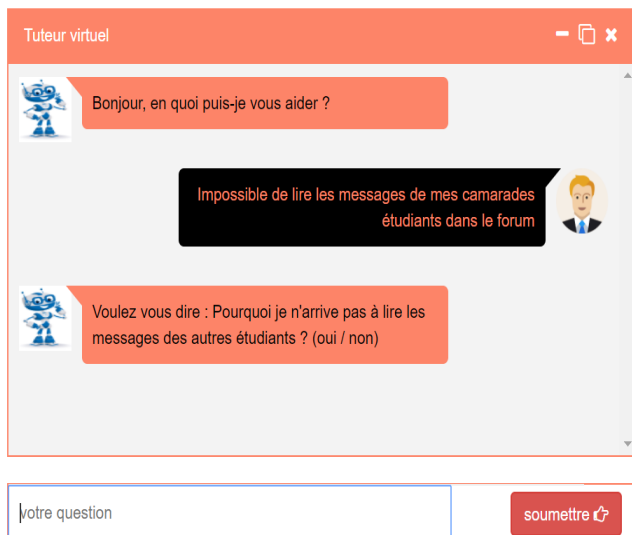


Fig. 8. Suggested question after the analysis of the learner's concern.

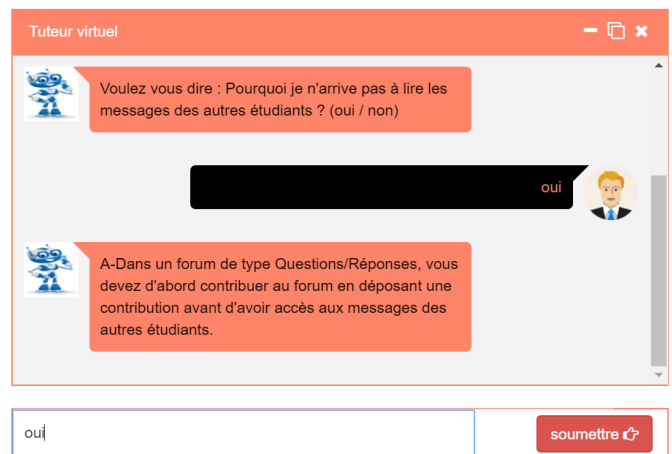


Fig. 9. The learner confirms the question proposal.

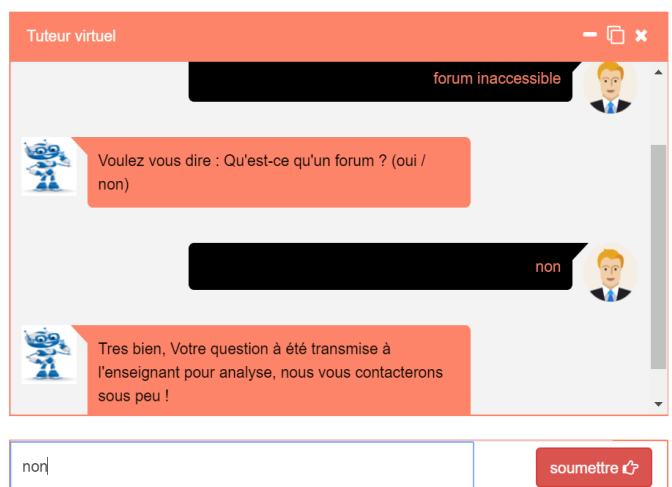


Fig. 10. The proposed question does not match the learner's concern.

The figures show the process used to respond to the learner. In the case where the question refers to a question from the teacher, the learner is asked if the question identified is the one he was referring to. If he answers yes, the instructions associated with his request are sent to him. If not, the learner's concern is sent to the human tutor.

Experimenting with the prototype of the chatbot makes it possible to find adequate answers to queries posted by the learner by applying our semantic similarity method.

A. Evaluating the Performance of the Chatbot Prototype

This assessment focuses on the learners' level of understanding of the learner's concerns. The tests are performed by applying the Dice Index method and the domain keyword-based concept (Dice Improvement) to the learner's concerns. The different tests are carried out with a series of concerns of the learners. It is a question of calculating the rate of comprehension of the concerns of the learners by the chatbot. The rate of comprehension of the questions based on the Dice index represents the ratio of the number of terms of the learner understood by the chatbot on the terms of the learner's question.

TE : The terms of the learner’s question

TEC : Learner’s terms understood

TC : Rate of understanding.

$$TC = \frac{TEC}{TE} \quad (7)$$

Calculation process of the understanding rate with the measure of Dice

TEC : 2 ; TE : 5 ; TC = 2/5

TC = 40%

The rate of understanding of the questions based on the concept of domain keywords (An improvement of the measure of Dice) represents the ratio of the number of terms of the teacher understood by the chatbot on the terms of the question of the teacher.

TS: The terms of the teacher's question that corresponds to TE

TSC: The terms of the teacher understood

TC: Rate of understanding.

$$TC = \frac{TSC}{TS} \quad (8)$$

Calculation process of the understanding rate with the concept of domain keywords

TSC: 3 ; TS : 7 ; TC = 3/7

TC = 43%

The table below represents the comprehension rate according to the number of questions asked by the learner (Table I).

NQ: The number of questions

MT: The method based on the index of Dice and concept of keywords of the domain

DICE: The measure of Dice

CMC: The concept based on the words of the domain

TC-DICE: Rate of understanding of the questions by applying the index of dice followed by the variation of the number of questions

TC-CMC: Rate of comprehension of the questions by applying the concept of words of the domain followed by the variation of the number of questions

TABLE I. RATE OF UNDERSTANDING OF THE QUESTIONS BASED ON THE DICE INDEX AND THE CONCEPT OF KEY WORDS IN THE FIELD

N Q	10		20		30		40		50	
	DIC E	CM C	DIC E	CM C	DIC E	CM C	DIC E	CM C	DIC E	CM C
T C	40 %	55 %	35 %	50 %	34 %	48 %	30 %	45 %	30 %	40 %

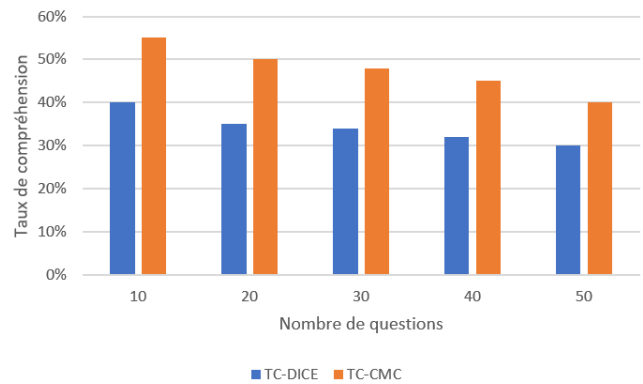


Fig. 11. Representation of the understanding rate of the questions based on the measure of Dice and the concept of keywords of the domain.

The graph above represents the rate of understanding of the questions based on the index of Dice and that based on concept of keywords of the domain (an improvement of the measure of Dice) (Fig. 11).

B. Results

Fig. 11 shows that the learner’s comprehension rate is above 50% when applying the Dice Index method. In addition, the rate of understanding of the questions is weaker and weaker as the number of questions increases while the student’s comprehension rate is above 50% when applying the concept of domain keywords. This concept is an improvement of the Dice Index. The results obtained with the concept based on the keywords of the domain are encouraging.

VI. CONCLUSION

In this paper, we presented work based on similarity measures to provide a chatbot with the ability to provide adequate responses in a learning interaction with learners. We have shown the steps to implement the prototype of the proposed chatbot that is an adaptation of the Dice Index. We also described the overall operation of the chatbot and the process used to address the learner's concern. The rest of the work consists in integrating the chatbot into the teaching of Université Virtuelle de Côte d'Ivoire and finish with the process of evaluation of learner’s satisfaction.

REFERENCES

- [1] W. Johnson, P. Rizzo, W. Bosma, S. Kole, M. Ghijsen, and H. van Welbergen. Generating socially appropriate tutorial dialog. In ISCA Workshop on Affective Dialogue Systems, pages 254-264. Berlin, Heidelberg, 2004.
- [2] M. Limniou, D. Roberts, and N. Papadopoulos. Full immersive virtual environment CAVE [TM] in chemistry education. Computers & Education, 51 (2): 584-593, 2008.
- [3] A. Kokane, H. Singhal, S. Mukherjee, and G. Reddy. Effective elearning using 3D virtual tutors and webRTC based multimedia chat. In International Conference on Recent Trends in Information Technology (ICRTIT), pages 1-6, 2014.
- [4] J. Rowe, S. McQuiggan, B. Mott, and J. Lester. Motivation in narrative-centered learning environments. Proceedings of the workshop on narrative learning environments, AIED, pages 40-49, 2007.
- [5] DECO 2017, Directorate of Examinations and Concours (DECO) - MINISTRY OF NATIONAL EDUCATION OF CÔTE D'IVOIRE - The Territorial Assembly then created on that date, the twelve (12) first departments. The decree n ° 30004 / CAB1 of May 25, 1957 fixing the

- attributions of the Ministry of National Education, creates at the same time the DECO - December 2017
- [6] EGOUV 2015 Official Portal of the Government of Ivory Coast - 2nd DAYS OF RITER: UNIVERSITY IMPORTS NEW DIGITAL SOLUTIONS FOR HIGHER EDUCATION
- [7] UVCI 2015, The Virtual University of Côte d'Ivoire abbreviated UVCI is a public university whose "decree n ° 2015-775 of December 9, 2015 establishing, attributions, organization and functioning of an administrative public institution"
- [8] Tim 2001, The term semantic web was proposed by Tim Berners Lee in 2001 ("The Semantic Web," Scientific American Magazine, May 17, 2001) an evolution of the web that would allow the available data to be usable and interpretable by software agents.
- [9] Mohamed Tahar Ileh, Imad Saleh. Virtual tutoring approach in an e-learning platform based on an SMDF multi-agent Virtual Tutor. Tahar 2009 Formal Metas Data System SMDF, Semantic Web, Foccus, Model, Remote Tutoring, Moodle and e-learning. Volume 2, Number 1, Pages 3-18
- [10] HAYNES and David. Metadata for information management and retrieval. London: Facet, 2004, pages 186.
- [11] Mirna El-Hajj Barbar, Saleh Imad, Barbar Kablan and Youssef Moncef, "An Automata Model for the Educational Evaluation Process: Generating an Automata Markup Language (AML)", Proceedings of the H2PTM'03 Conference, 24-26 September 2003, Hermès publishing, 9 Pages.
- [12] Guizzardi, G., Almeida Falbo, R., and Pereira Filho, J. G. (2001). From domain ontologies to Object Oriented Frameworks. In G. Stumme, A. Maedche, & S. Staab (Eds.), Proceedings of the ONTO 2001 Workshop on Ontologies (pp. 1-14). (THIS Workshop Proceedings, Vol 48). Germany: CEUR.
- [13] Smith, M.K., Welty, C. and McGuinness, D.L. (2004). OWL: Ontology Web Language Guide. Technical Report, W3C: World Wide Web Consortium.
- [14] Marsh, J. (2001). XML base. Technical Report, W3C: World Wide Web Consortium.
- [15] Klyne, G. and Carroll, J. J. (2004). Resource Description Framework (RDF): Concepts and Abstract Syntax. Technical Report, W3C: World Wide Web Consortium.
- [16] Connolly, D., Harmelen, F. V., Horrocks, I., McGuinness, D.L., Patel-Schneider, P.F. and Stein, L.A. (2001). DAML + OIL: Reference Description. W3C Technical Report: World Wide Web Consortium.
- [17] Messous 2014, a virtual teaching agent - an interactive environment for human learning (EIAH) featuring an educational hypermedia and an animated conversational agent (ACA) operating in a web environment. Pages 16
- [18] Sassi 2014, Assistance and advice to users in the context of ambient intelligence: a study of uses in context: X-CAMPUS eXtensible Conversational Agent for Multichannel Proactive Ubiquitous Services. Hajer Sassi, Lille 1, 2014
- [19] Joanna Taoum, Bilal Nakhil, Elisabetta Bevacqua, Ronan Querrec. A design proposal for interactive virtual tutors in an informed environment. 11th Days of the French Association of Virtual Reality (AFRV 2016), Oct 2016, Brest, France.
- [20] Huang H, et al. (2008) The crystal structure and identification of NQM1 / YGR043C, a transaldolase from *Saccharomyces cerevisiae*. *Proteins* 73 (4): 1076-81
- [21] Christine Largeron, Bernard Kaddour, Maria Fernandez. SoftJaccard: a measure of similarity between sets of strings for the unification of named entities. *Extraction and Knowledge Management (EGC 2009)*, Jan 2009, Strasbourg, France. Cépaduès-Éditions, RNTI-E-15, pp.443-444, 2009.
- [22] Measure the similarity between sentences thanks to Wikipedia using a random indexing Hai-Hieu Vu, Jeanne Villaneau, Farida Said, Pierre-François Marteau 2015
- [23] Goutam Majumder, Partha Pakray, Alexander Gelbukh, David Pinto; Semantic Textual Similar Methods, Tools, and Applications: A Survey Goutan 2016

Formalization of Behavior Change Theories to Accomplish a Health Behavior

Adnan Manzoor

Department of Information Technology
Quaid e Awam UEST
Nawabshah, Pakistan

Sohail Khokhar

Department of Electrical Engineering
Quaid e Awam UEST
Nawabshah, Pakistan

Imtiaz Ali Halepoto

Department of Computer Systems Engineering
Quaid e Awam UEST
Nawabshah, Pakistan

Nazar Hussain Phulpoto

Department of Information Technology
Quaid e Awam UEST
Nawabshah, Pakistan

Engr. Muhammad Sulleman Memon

Department of Computer Systems Engineering
Quaid e Awam UEST
Nawabshah, Pakistan

Abstract—The objective of this paper is to study theories behind behavior change and adaptation of behavior. Humans often live according to habitual behavior. Changing an existing behavior or adopting a new (healthier) behavior is not an easy task. There are a number of things which are important when considering adapting physical activity behavior. A behavior is affected by various cognitive processes, for example involving beliefs, intentions, goals, impediments. A conceptual and computational model is discussed based on state of the art theories about behavior change. The model combines different theories: the social cognitive theory, and the theory of self-regulation. In addition, health behavior interventions are discussed that may be used in a coaching system. The paper consists of two parts: the first part describes a computational model of behavior change and the second part discusses the formalization of evidence-based techniques for behavior change and questions to measure the various states of mind in order to provide tailored and personalized support.

Keywords—Behavior monitoring; healthy lifestyle; behavior change; physical activity; computational model

I. INTRODUCTION

In this paper we explore and study theories and models about behavior change, in which the particular focus is on adopting a healthy behavior. In this work we study what state of the art theories about behavior change propose and how a theory about health behavior change can be formalized in a way that it can be used with a mobile system to motivate and encourage young individuals to adopt a more physically active lifestyle. Mobile and sensor technology provide a promising way to support health behavior change interventions. Depending on what aspect of a healthy lifestyle is intended to be addressed, knowledge from different domains can be used. For instance, if the goal is to achieve a good mood or lowering daily life stress [1], [2], then theories related to emotion

regulation could be used [3]. In the first half of this paper various theories related to behavior change are discussed and one of them is used to formalize in terms of a computational model. This computational model is implemented in the virtual coaching system whose purpose is to help young adults to physically be more active. The computational model is adopted from [4]. The second part discusses the identification of various behavior change techniques and how they can be used in a mobile system; this is based on [5]. In this part we explore theories and evidence based research to understand how different theories can play a role in the design and development of an intelligent coaching system that can help motivate people to achieve an active lifestyle. One of the first steps is to identify behavioral determinants of physical activity behavior. In this respect different health behavior models have been suggested in the literature, over the years different theories proposed different kinds of focal determinants for health behavior. Some of those theories and their focal determinants are discussed below. The computational model used by the coaching system to predict health behavior determinants. Each determinant (in the computational model) is linked to a behavior change technique and a number of motivational messages. These messages help user to increase their motivation level towards physical activity health behavior.

II. THEORIES ON BEHAVIOR CHANGE

In this section, we will review the most important (psychological) theories on behavior change. Based on this discussion, a computational model of behavior change will be presented in the subsequent section.

A. The Transtheoretical Model

The transtheoretical model proposes one theory of behavior change [6]. It consists of four major components: stages of

change, processes of change, decisional balance, and self-efficacy. The model proposes that in order to change unhealthy behavior or adopt a healthy lifestyle, people go through six stages of change. These stages of change describe that a change of behavior occurs over a time period through different stages in different time slots. The first is precontemplation in which there is no awareness of the advantages and disadvantages of a behavior change and no change occurs in next six months. The Contemplation stage describes when people are better aware of the pros and cons of a behavior and in the next six months the change is likely. In the Preparation stage individuals are to change; this can take place in the next month. The Action stage is when people take the actual action to change the behavior. The Maintenance phase is when an individual avoids a situation where a relapse can happen; usually it lasts for a longer period of time which may be from six months to five years. In this stage usually individuals feel less temptation towards the unhealthy behavior. The last is the Termination stage determines the time point where a person do not feel any temptation to go back to the old habit.

B. Health Belief Model

Another model which is often referred to in health behavior literature and used in various empirical studies is the Health Belief model [7]. The Health belief model mainly consists of four constructs. Perceived susceptibility determines a person's subjective perspective on the risks of getting an illness, perceived severity is the belief which determines an individual's opinions about the seriousness of an illness. The last two are perceived benefits and perceived barriers; these are the beliefs concerning the potential benefits of adapting a healthy behavior and the likely hurdles one would face. If a person would give more weightage to potential benefits than the barriers, then he/she most likely adapt the behavior. Later the model was extended by including another construct: self-efficacy [8]. Self-efficacy is the extent to which a person believes that he/she is able to change his/her behavior; it is further discussed in the section about the Social Cognitive Theory below.

C. Theory of Planned Behavior

The Theory of planned behavior [9] is also one of the widely used models in health behavior, it is the extension of another theory by the same author, the theory of reasoned actions [10]. The model suggests that a health behavior is determined by attitudes and social influences through behavior intentions. Therefore, if an individual's attitude towards a health behavior is positive and a perceived social pressure to perform that health behavior is also present, then the individual's behavioral intention is higher to execute certain health behavior.

D. Ecological Models of Behavior Change

The proponents of ecological models propose that a population-wise change can only be achieved by taking into account various environmental factors; for example, a person's physical and social environment also play an important role in achieving behavior change. Ecological models [11] of health behavior suggest that there are more factors that could contribute to a health behavior than psychological and social factors. These may include a person's physical environment,

but also other ecological factors are taken into consideration, for example, the physical infrastructure, the community of a person, the policies, etc.

E. The Social Cognitive Theory

In the area of active physical behavior, one of the most often used theories is the Social cognitive theory [12], [13]. The most important concept in social cognitive theory is self-efficacy. As mentioned above, self-efficacy is a person's confidence in his/her capabilities to perform a certain action successfully. A high self-efficacy is associated with strong determination in the wake of difficulties and not giving up. A low self-efficacy is associated with difficulty in performing/committing to a challenging task and to give up in the wake of difficulties. Self-efficacy can be strengthened by social support and satisfaction about a difficult task done in the past. Friends can also help a person to remind him/her about the person's success on a difficult task in the past. Furthermore, self-efficacy affects the behavior through multiple paths, it can affect it directly or it can affect it indirectly through intentions or outcome expectations. Self-efficacy plays an essential role in a number of behaviors. Different kinds of interventions have been proposed based on the targeted behavior, for example, for smoking cessation [5]. Further important determinants are satisfaction intentions, impediments, long term goals, social norms. Reciprocal determinism also is an important concept in the Social cognitive theory; it proposes that a person's behavior, his/her personal characteristics and his/her environment have reciprocal interactions to determine the subsequent behavior of the person. Outcome expectation is another important concept in the theory. It determines what people expect after performing an action, outcome expectation is further based on physical outcome expectations, social outcome expectations and personal outcome expectations. Physical outcome expectations determine the changes that are felt in the body, for example, after exercising one might have a good feeling. Social outcome expectations determine people's behavior towards the action which is carried out and personal outcome expectations represent the extent to which one expects to be better or worse off.

Impediment determines the obstacle one encounters during the course of an action or before performing an action. For example an individual might avoid going to sports school by bike if the weather is not conducive or if it is raining. Another common impediment is lack of time. Social norms affect the behavior through the intentions; social norm are people's reaction towards an action. Goals are also important ingredient of the Social cognitive theory; goals can be divided into distal/long term and proximal/short term (i.e. intentions).

III. COMPUTATIONAL MODEL

Given the fact that the Social cognitive theory is one of the most often used theories in the area of active physical behavior and that we also intend to apply it in practice, choosing this theory as the basis of creating a computational model is a natural choice. This section describes a formalization of this model. The conceptualization of the model is depicted in Fig. 1.

Often in the literature, the emphasis is given to some determinants of the Social cognitive theory, but in the model that we present, all the ingredients are used. The model describes the dynamic relations among the determinants, which helps to understand how they work together to influence physical activity behavior. Below we discuss some of the concepts (with their formalization) in the model, they are adopted from [4], (for the remaining formalization please see preceding paper). Simulation result in Fig. 2 is adopted from [4], it suggest that when a person with inactive behavior encounters impediments it can cause an improvement in his/her behavior if he/she is able to surpass the potential impediments, for the rest of other results please see [4].

A. Self-efficacy

In the context of physical activity, the definition of self-efficacy is a person’s judgment about his/her capability to successfully indulge in an active lifestyle. Self-efficacy can be strengthened by social comparison, social support, accomplishing a task successfully, and psychological responses (for example a person does not feel any anxiety or fear before giving a presentation). One of the sources of self-efficacy is satisfaction about a behavior or an accomplishment. A higher satisfaction leads to an increased self-efficacy and “a lower satisfaction level results in a decreased self-efficacy”. It can be computed in the following way.

$$\text{If } SE(t) \geq Sat(t): SE(t + \Delta t) = SE(t) + \beta_{Sat,SE} \cdot (Sat(t) - SE(t)) \cdot SE(t) \cdot \Delta t$$

$$\text{If } SE(t) < Sat(t): SE(t + \Delta t) = SE(t) + \beta_{Sat,SE} \cdot (Sat(t) - SE(t)) \cdot (1 - SE(t)) \cdot \Delta t$$

B. Outcome Expectations

Outcome expectations are described in terms of three further outcomes i.e. physical outcome expectations, social outcome expectations, and personal outcome expectations. OE is computed with the following formulas:

$$OE^*(t) = (\omega_{SOE} \cdot SOE(t) + \omega_{POE} \cdot POE(t) + \omega_{PhOE} \cdot PhOE(t)) / (\omega_{SOE} + \omega_{POE} + \omega_{PhOE})$$

$$\text{If } OE^*(t) \geq SE(t): OE(t + \Delta t) = OE^*(t) + \beta_{SE,OE} \cdot (SE(t) - OE^*(t)) \cdot OE^*(t) \cdot \Delta t$$

$$\text{If } OE^*(t) < SE(t): OE(t + \Delta t) = OE^*(t) + \beta_{SE,OE} \cdot (SE(t) - OE^*(t)) \cdot (1 - OE^*(t)) \cdot \Delta t$$

C. Impediments

Confidence in one’s capability can help overcome an impediment. Other factors also play a role to avoid temptations and overcome impediments, for example self-regulation. But, if a person’s self-efficacy level is higher it helps him/her to be stronger in the face of difficulties. It is computed by the difference between self-efficacy and the impediments.

$$\text{If } SE(t) \geq Imp(t): Imp(t + \Delta t) = Imp(t) - \beta_{SE,Imp} \cdot (SE(t) - Imp(t)) \cdot Imp(t) \cdot \Delta t$$

$$\text{If } SE(t) < Imp(t): Imp(t + \Delta t) = Imp(t) - \beta_{SE,Imp} \cdot (SE(t) - Imp(t)) \cdot (1 - Imp(t)) \cdot \Delta t$$

D. Intentions

Intentions are short term goals. Outcomes expectations influence in the process of goal formation, as people aim for those action for which they expect positive consequences. Self-efficacy and outcome expectations together determine intentions. Furthermore, intentions are also affected by the impediments and the facilitators.

$$\text{Change_Int}(t) = \beta_{SE,Int} \cdot (SE(t) - Int(t)) + \beta_{SOE,Int} \cdot (SOE(t) - Int(t)) + \beta_{Fac,Int} \cdot Fac(t) - \beta_{Imp,Int} \cdot Imp(t)$$

$$\text{If } \text{Change_Int}(t) \geq 0: Int(t + \Delta t) = Int(t) + \text{Change_Int}(t) \cdot (1 - Int(t)) \cdot \Delta t$$

$$\text{If } \text{Change_Int}(t) < 0: Int(t + \Delta t) = Int(t) + \text{Change_Int}(t) \cdot Int(t) \cdot \Delta t$$

E. Satisfaction

As seen in Fig. 1, the satisfaction has incoming edges from three states i.e. behavior, intention, and facilitators/impediments. It is computed by combining these states in the following formula.

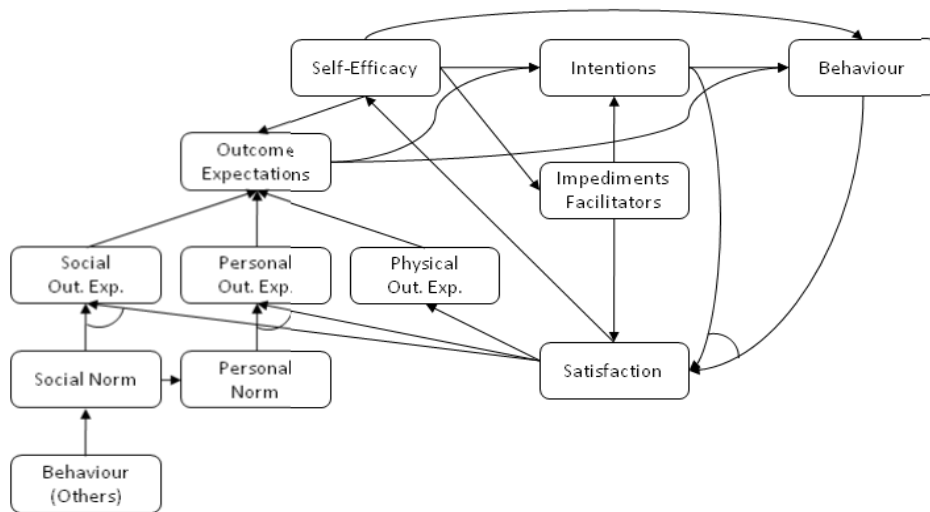


Fig. 1. A computational model for behavior change based on the social cognitive theory; the model is adopted from [4].

$$\text{Change_Sat}(t) = \beta_{\text{Int\&Beh,Sat}} \cdot (\text{Beh}(t) - \text{Int}(t)) + \beta_{\text{Imp,Sat}} \cdot \text{Imp}(t) - \beta_{\text{Fac,Sat}} \cdot \text{Fac}(t)$$

$$\text{If } \text{Change_Sat}(t) \geq 0: \text{Sat}(t + \Delta t) = \text{Sat}(t) + \text{Change_Sat}(t) \cdot (1 - \text{Sat}(t)) \cdot \Delta t$$

$$\text{If } \text{Change_Sat}(t) < 0: \text{Sat}(t + \Delta t) = \text{Sat}(t) + \text{Change_Sat}(t) \cdot \text{Sat}(t) \cdot \Delta t$$

F. Behavior

Behavior is influenced by self-efficacy, outcome expectations, and intentions. It is implemented in similar way: there are incoming edges from these states to behavior. In the current context a behavior describes a person's activity level, if value of 0 shows that person is not active and a value of 1 describes that the individual is very active.

$$\text{Change_Beh}(t) = \beta_{\text{SE,Beh}} \cdot (\text{SE}(t) - \text{Beh}(t)) + \beta_{\text{OE,Beh}} \cdot (\text{OE}(t) - \text{Beh}(t)) + \beta_{\text{Int,Beh}} \cdot (\text{Int}(t) - \text{Beh}(t)) + \beta_{\text{Fac,Beh}} \cdot \text{Fac}(t) - \beta_{\text{Imp,Beh}} \cdot \text{Imp}(t)$$

$$\text{If } \text{Change_Beh}(t) \geq 0: \text{Beh}(t + \Delta t) = \text{Beh}(t) + \text{Change_Beh}(t) \cdot (1 - \text{Beh}(t)) \cdot \Delta t$$

$$\text{If } \text{Change_Beh}(t) < 0: \text{Beh}(t + \Delta t) = \text{Beh}(t) + \text{Change_Beh}(t) \cdot \text{Beh}(t) \cdot \Delta t$$

IV. PROPOSING AND IMPLEMENTING BEHAVIOR CHANGE TECHNIQUES

The second part of this paper is concerned with the formalization of behavior change techniques (BCTs). This section is partially based on the research conducted in [5]. We first discuss the relation between behavior change and modern (mobile/sensor) technology, then we describe the concept of tailoring, and finally we identify BCTs and describe how they can be translated into actual tailored messages that can be used within a mobile system for stimulating physical activity. The BCTs described in [5] are adopted from the widely used taxonomy proposed by Michie et al. They suggested a standardized collection of BCTs which can be linked to various determinants in a theoretical framework [14].

A. Behavior Change Techniques and Modern Technology

A smartphone is an inherently a personalized device and therefore it seems to be a well suited medium to implement

BCTs and deliver tailored messages, although there is not much research done to check the effectiveness of mHealth based tailored messages [15]. Recently, new trends have emerged in an effort to help and support behavior change. Lyons et al. have conducted a study to analyze whether modern wearable devices do implement empirically tested BCTs. Based on their systematic analysis it was demonstrated that many BCTs are supported in the wearables for instance some of the most common BCTs which are supported by state of the art wearables include self-monitoring, feedback, and goal-setting. However, many of the wearables do not support action planning, behavioral practice, and problem solving which are considered essential BCTs in the physical activity domain [16].

Some individuals, who are known as Quantified-selfers are already employing modern devices for self-monitoring without the aid of an explicit intervention. Quantified-Selfers [17] are described as extreme users who form an intrinsically motivated group of people who are interested to measure various aspects of their lives, for example, to improve health or improve sleep quality [18]. Furthermore, they are also considered as proactive people who want to monitor their health (e.g., glucose levels and blood pressure) to deal with unlikely circumstances. Choe et al. conducted a study in which they interviewed people who are very much interested in quantifying different aspects of their lives, to extract what motivates them or what kind of problems they encounter during the quantification process. It is seen that in the beginning individuals are better motivated so they track many thing simultaneously but eventually some people give up because there is a large amount of data and managing and analyzing it is not trivial task. Some people's motivation decreases because no automatic feedback is provided to help them to remind them about the goal or, for example, what to track and what not to track. This is particularly in line with the study conducted by Lyon et al. [16], in which they found that one of the important features lacking in activity monitors is action planning.

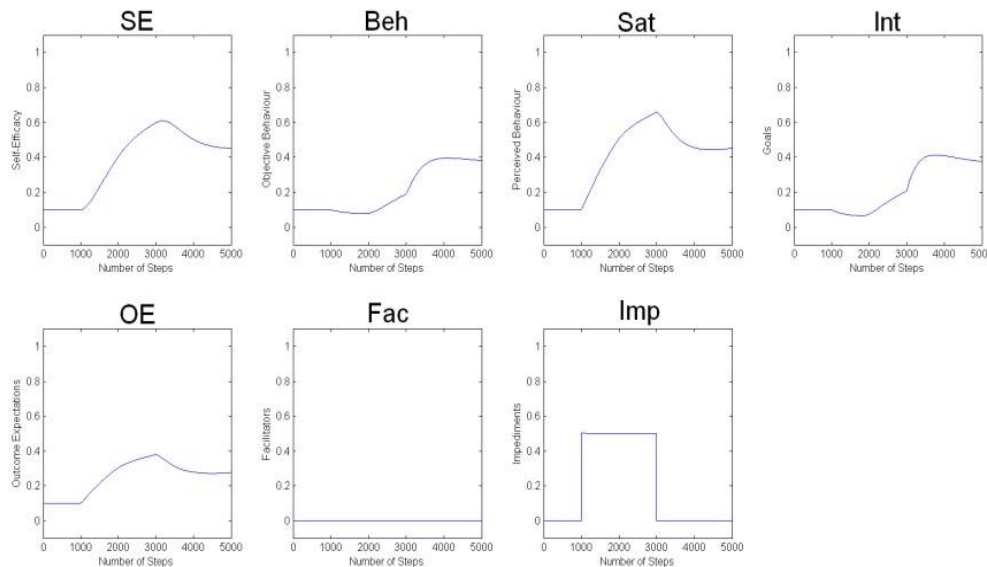


Fig. 2. An example simulation scenario of the model, as presented in (Mollee & van der Wal, 2013).

Personal informatics is another synonym term for quantified-self, Li et al. conducted a survey, and based on it they propose their five stages model to self-discovery [19]. The first stage corresponds to preparation; this stage takes place before starting a system. It consists of mainly two things: what kind of information people are motivated to record and what kind of software or wearable/tracker people would like to use or a combination of both. In this stage people face difficulties because every kind of tool has its own format and using multiple tools or switching between tools can lead to loss of information that was recorded previously. Another problem is that initially people are motivated enough to use a variety of tools, but after some time their motivation becomes less. Collection is the second stage in which people actually start recording data about different aspects of their personal life. In this stage people encounter problems that sometimes there is a lack of motivation or lack of time. The next stage is Integration; the challenge in this stage is to combine and integrate the data when it comes from multiple sources. It could lead to problems to move to next stage of Reflection which helps the person to explore and think about different aspects of the information. And Action is the last stage when people need to take certain action based on their new self. While technologies such as Personal Informatics and Wearables are promising for behavior change, their potential for behavior change tools can be leveraged by using evidence-based behavior change techniques and established behavior change models such as discussed in Section 2 of this paper.

B. Tailored Messages

There is difference between computer generated and computer tailored feedback [20], [21]. Computer generated feedback is merely a static delivery of contents without any personalization features attached to the feedback. In contrast, in computer tailored feedback an individual's psychological, social, and physical states are measured and then based on that the tailored feedback is generated by applying some algorithm(s). Moreover, there is also a difference between dynamically tailored and static tailored feedback. In the dynamic variant, timely feedback is generated based on on-going assessment, while static feedback is based on onetime assessment [20]. It has been shown that dynamically tailored feedback is significantly better. It also has been found that tailored messages have stronger effect than non-tailored messages [21]. Tailored messages require strategies/suggestions that are targeted to an individual based on his/her unique mental, physiological and environmental attributes rather than for an entire group, which are often referred as targeted generic messages [21]. In practice, software tailored messages have also shown more promising results than untailored messages [21].

C. Behavior Change Techniques in a mobile support system

To implement an automated support system for behavior change, we should identify evidence-based techniques which can be linked in an algorithmic way to the computational analysis of the determinants of behavior change. In this section,

we identify appropriate BCTs (in the physical activity domain) based on the literature (Michie et al., 2013) and link them to behavior determinants that were identified in the previous model (see Section 2). Second, it should be decided how to implement various BCTs. For example, one of the ways is to implement them in terms of feedback messages. The messages need to be tailored to each individual. Various measures can be used to personalize the messages, i.e. an individual's activity data, location data, social network and behavior determinants which can be determined by various questions.

In Table I, an overview is provided on how determinants used in the computational model described in Section 3 are related to BCTs from the taxonomy of Michie and how they can be applied in a mobile support system. The table also shows from which behavior change theory the determinants stem (see Section 2).

V. CONCLUSION

In this paper various theories related to behavior change were discussed. The Social Cognitive Theory is a widely used theory in the domain of physical activity behavior. A model based on this theory which was earlier formalized [4] is adopted in this paper. The objective of the computational model is to provide a computational means for reasoning about behavior change in a coaching system. We have discussed that modern technology such as wearables do help to achieve a behavior change to some extent but they do not yet fulfill the role of a personalized coaching system. Based on a model, a coaching system can predict behavior and generate context specific tailored messages for the users of the system. These messages are based on evidence-based strategies [5] to improve physical activity behavior, which are linked to particular BCTs and determinants in the model. Currently state of the art techniques (which are based on mobile apps) to improve health behavior do not support evidence-based behavior change techniques [15].

If the tailored interventions are combined with network interventions [22] they have a great potential for strengthening health behavior change. Social network interventions can play an important role to achieve a behavior change by utilizing the structure of the network and measuring various characteristics (mental and physical states) of individuals in the network. In addition to personalized feedback, this provides the possibility to create another type of support based on different social phenomena such as social support, social comparison, social contagion, etc. A number of strategies [22], [23] exist to find people in a social network who can work as change agents. For instance, it is possible to identify people in a social network based on the similar characteristics who can provide social support to each other. Social interventions are especially relevant as one of the important behavior change technique related to self-efficacy is social modelling or social comparison through which self-efficacy can be improved/increased [24]. A number of social interventions are proposed above, but further research can be done on implementing them in support systems.

TABLE I. LIST OF DETERMINANTS, RELATED BEHAVIOR CHANGE TECHNIQUES AND THE POSSIBLE IMPLEMENTATION IN AN AUTOMATED SYSTEM. NOTE THAT THE TABLE IS ADOPTED FROM ANOUK MIDDELWEERD ET AL. 2017

Determinant	Behavior Change Technique	Implementation
Outcome Expectations	Provide general information on consequences of behavior in general	Message in general and tailored to aspects of the questionnaire
Self-efficacy	Action planning/ time management	Message
	Social comparison	Graph tailored to preference social comparison (up-/ down wards)
	Persuasion	Persuasive messages on how to overcome barriers
	Prompt self-monitoring	Message & graph
	Plan social support	Message
Intention	Imaginary reward	Message
	Progress towards goal	Message
	Motivational messages	Message
	Modeling	Message & graph
	Prompt instruction	message
Impediments	Prompt goal setting	Message & suggestion
	Prompt Barrier identification	message
Social Norm (Descriptive and Inductive)	Social comparisons	Graph tailored to preference social comparisons (up-/ down wards)
	Information about other's approval	message
Self-regulation	self-monitoring	Graph & message
	Goalsetting	Message
	Progress towards goal	Graph & message
	Self-evaluation	message
	Imaginary reward	message
Satisfaction	Self-evaluation	message
Long-term goals	Provide general information on consequences of behavior in general	message

REFERENCES

[1] A. H. Abro, M. C. A. Klein, A. R. Manzoor, S. A. Tabatabaei, and J. Treur, "Modeling the effect of regulation of negative emotions on mood," *Biol. Inspired Cogn. Archit.*, vol. 13, pp. 35–47, Jul. 2015.

[2] L. Medeiros and T. Bosse, "Empirical Analysis of Social Support Provided via Social Media," in *Social Informatics*, 2016, pp. 439–453.

[3] J. J. Gross, "The emerging field of emotion regulation: An integrative review," *Rev. Gen. Psychol.*, vol. 2, no. 3, pp. 271–299, 1998.

[4] J. S. Mollee and C. N. van der Wal, "A computational agent model of influences on physical activity based on the social cognitive theory," in *International Conference on Principles and Practice of Multi-Agent Systems*, 2013, pp. 478–485.

[5] Anouk Middelweerd, Saskia J te Velde, Julia S Mollee, Michel CA Klein, and Johannes Brug, "Description of the development and content of Active2Gether: an app-based intervention combining evidence-based behavior change techniques with a model-based reasoning system to promote physical activity among young adults," *J. Med. Internet Res. Rev.*, 2017.

[6] J. O. Prochaska and W. F. Velicer, "The Transtheoretical Model of Health Behavior Change," *Am. J. Health Promot.*, vol. 12, no. 1, pp. 38–48, Sep. 1997.

[7] M. H. Becker, "The health belief model and sick role behavior," *Health Educ. Monogr.*, vol. 2, no. 4, pp. 409–419, 1974.

[8] I. M. Rosenstock, V. J. Strecher, and M. H. Becker, "Social learning theory and the health belief model," *Health Educ. Q.*, vol. 15, no. 2, pp. 175–183, 1988.

[9] I. Ajzen, "From intentions to actions: A theory of planned behavior," in *Action control*, Springer, 1985, pp. 11–39.

[10] M. Fishbein and I. Ajzen, "Belief, attitude, intention, and behavior: An introduction to theory and research," 1977.

[11] J. F. Sallis, R. B. Cervero, W. Ascher, K. A. Henderson, M. K. Kraft, and J. Kerr, "An ecological approach to creating active living communities," *Annu Rev Public Health*, vol. 27, pp. 297–322, 2006.

[12] A. Bandura, "Health promotion from the perspective of social cognitive theory," *Psychol. Health*, vol. 13, no. 4, pp. 623–649, 1998.

[13] A. Bandura, "Health promotion by social cognitive means," *Health Educ. Behav.*, vol. 31, no. 2, pp. 143–164, 2004.

[14] S. Michie et al., "The behavior change technique taxonomy (v1) of 93 hierarchically clustered techniques: building an international consensus for the reporting of behavior change interventions," *Ann. Behav. Med.*, vol. 46, no. 1, pp. 81–95, 2013.

[15] A. Middelweerd, J. S. Mollee, C. N. van der Wal, J. Brug, and S. J. te Velde, "Apps to promote physical activity among adults: a review and content analysis," *Int. J. Behav. Nutr. Phys. Act.*, vol. 11, no. 1, p. 97, 2014.

[16] E. J. Lyons, Z. H. Lewis, B. G. Mayrsohn, and J. L. Rowland, "Behavior change techniques implemented in electronic lifestyle activity monitors: a systematic content analysis," *J. Med. Internet Res.*, vol. 16, no. 8, p. e192, 2014.

[17] E. K. Choe, N. B. Lee, B. Lee, W. Pratt, and J. A. Kientz, "Understanding quantified-selfers' practices in collecting and exploring personal data," in *Proceedings of the 32nd annual ACM conference on Human factors in computing systems*, 2014, pp. 1143–1152.

[18] M. Swan, "Health 2050: The realization of personalized medicine through crowdsourcing, the quantified self, and the participatory biocitizen," *J. Pers. Med.*, vol. 2, no. 3, pp. 93–118, 2012.

[19] I. Li, A. Dey, and J. Forlizzi, "A stage-based model of personal informatics systems," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2010, pp. 557–566.

[20] P. Krebs, J. O. Prochaska, and J. S. Rossi, "A meta-analysis of computer-tailored interventions for health behavior change," *Prev. Med.*, vol. 51, no. 3, pp. 214–221, 2010.

[21] M. W. Kreuter, V. J. Strecher, and B. Glassman, "One size does not fit all: the case for tailoring print materials," *Ann. Behav. Med.*, vol. 21, no. 4, pp. 276–283, 1999.

[22] T. W. Valente, "Network interventions," *Science*, vol. 337, no. 6090, pp. 49–53, 2012.

[23] S. P. Borgatti, "Identifying sets of key players in a social network," *Comput. Math. Organ. Theory*, vol. 12, no. 1, pp. 21–34, 2006.

[24] N. Michinov, "Social comparison, perceived control, and occupational burnout," *Appl. Psychol.*, vol. 54, no. 1, pp. 99–118, 2005.

Implementation of Winnowing Algorithm with Dictionary English-Indonesia Technique to Detect Plagiarism

Anton Yudhana, Sunardi

Department of Electrical Engineering
Universitas Ahmad Dahlan Yogyakarta,
Indonesia

Iif Alfiatul Mukaromah

Master of Informatics Engineering
Universitas Ahmad Dahlan Yogyakarta,
Indonesia

Abstract—The ease of obtaining information that is easy, fast, and cheap from all over the world through the internet network can encourage someone to take action plagiarism. Plagiarism is an intellectual crime that often occurs in the writing world where the perpetrators take the work of others without declaring the original source; if it continues to be left it will have a negative impact on the academic community and can be a chronic disease in the progress of a nation. At this time, the process of plagiarism detection is done manually and automatically using the help of technological developments (plagiarism detection), but the automatic checks available now mostly just check every letter character contained in the document, cannot check where the plagiarist takes a quote from a foreign language and changed in plagiarist language. Detection of plagiarism in this study will use a winnowing algorithm that has a function to check every character in two samples by hashing method that can generate fingerprint from two documents. While the dictionary method English-Indonesia change the writing from English to Indonesian language. This research will produce plagiarism detection using winnowing algorithm with English-Indonesian dictionary technique.

Keywords—Plagiarism; winnowing algorithm; fingerprint; dictionary English-Indonesia

I. INTRODUCTION

Plagiarism is an intellectual crime and an unlawful act in which the offender attempts to take the work of another person either whole, in part or in small part without permission or without mentioning the original owner, so that the act of plagiarism is the same as the act of stealing [1], [2]. The act of plagiarism is not a new act but a practice that is often done (no stranger) again in the country of Indonesia and throughout the world that occurs in the academic world, the world of writing and in our society [1].

Plagiarism action occurs one of them is the ease of obtaining an information that can be accessed and taken any time example from the internet, the internet is the sophistication of technology that every year progressively rapid development [3], [4]. As the discovery of zalnur from the results of his research that there are two factors causing the occurrence of plagiarism among students that is the development of information technology is increasingly sophisticated and the burden of assignment given to the student

is heavy enough so that many students choose an instant path by doing the act of plagiarism [1].

The act of plagiarism can not be allowed to develop especially in the world of education, this action can damage the academic community and can result in the decline of a nation because its critical mindset is not honed [2]. According to [2], a nation will experience decline because someone is lazy to think, no more development of science or new discoveries produced from the children of the nation. Even the plagiarism act indicates weak character education [1]. There needs to be action that can minimize the action of plagiarism in order not to become a habit of a nation. Even plagiarism is not only done by copy-pasting and changing every word for word with the same meaning, but some people do plagiarism by taking the work of foreign writing and converting it into another language (translated into plagiarist language) [3].

Prior research has implemented several algorithms that function as document fingerprints to detect plagiarism such as rabin-karp algorithm, winnowing algorithm, smith-waterman algorithm and so on [5]-[8]. However, these algorithms can only check every word in the document file and can not check the action of plagiarism done by taking other people's work in foreign language written by another language. In this case researchers will use the winnowing algorithm to detect plagiarism and this study can only check plagiarism in the category of translated plagiarism and plagiarism ideas in English sentences translated into Indonesian. Winnowing algorithm is an algorithm that has a function to check the similarity of words using hashing techniques and will be collated with dictionary English-Indonesia technique that serves to translate the writings from english to Indonesia.

Through this research, it is hoped that the application of plagiarism detection using winnowing algorithm with dictionary-english-indonesia technique can minimize the action of plagiarism especially in the category of translated plagiarism and idea plagiarism which is done by taking or quoting the writing in English which is translated into Indonesian.

II. RESEARCH METHOD

In this research will be focused on the implementation of winnowing algorithm with dictionary English-Indonesia technique to detect the existence of an action of plagiarism.

A. Categories of Plagiarism

The act of plagiarism is an act of stealing against the work of others because it does not reveal its original source [3]. Some of the things that plagiarists do is to quote or steal other people's work by paraphrasing quotations so that their actions are unknown.

According to B. Gipp and N. Meuschke (2011) in his research categorize the act of plagiarism based on the means used, among them [3]:

- Copy & paste plagiarism, copying entirely without any changes
- Disguised plagiarism, covering the copied parts, such as shake & paste, expansive plagiarism, contractive plagiarism, and mosaic plagiarism.
- Technical disguise, summarizing quotations to be subject to automatic detection by replacing letters with foreign letters.
- Undue paraphrasing, Paraphrasing quotations or alien thinking into the plagiarist language and hiding the original owner
- Translated plagiarism, translating quotations from one language to another.
- Idea plagiarism, using foreign ideas without declaring the source.

B. Winoing Algorithm

Winoing algorithm is an algorithm that has function as document fingerprint which is used to check the similarity of words in two documents by utilizing fingerprint concept with hashing technique [7], [9].

The winnowing algorithm is the exclusion of the rabin-karp algorithm by adding a window concept. Every word in the document will first be foxed in the hash form using the hash rolling formula by changing the characters in the document into ASCII code [10].

Here is the rolling hash formula that will be shown in (1) and (2).

$$H(C_1..C_l) = C_1 \cdot b^{(l-1)} + C_2 \cdot b^{(l-2)} + \dots + C_{(l-1)} \cdot b + C_l \quad (1)$$

$$H(c_2..c_{l+1}) = (H(c_1..c_l) - c_1 \cdot b^{(l-1)}) \cdot b + c^{(l-1)} \quad (2)$$

Where:

- $H(C_1..C_l)$ = hash value
- C_l = ASCII value of Character to -1 on string
- l = string length
- b = hash base value

The winnowing algorithm uses a certain window size and each window has a fingerprint that will be used to check the similarity of words on two documents or samples. The fingerprint to be selected is the smallest fingerprint, if there are two fingerprints in one window then select the rightmost fingerprint [9]-[11]. Fig. 1 is a concept of winnowing algorithm.

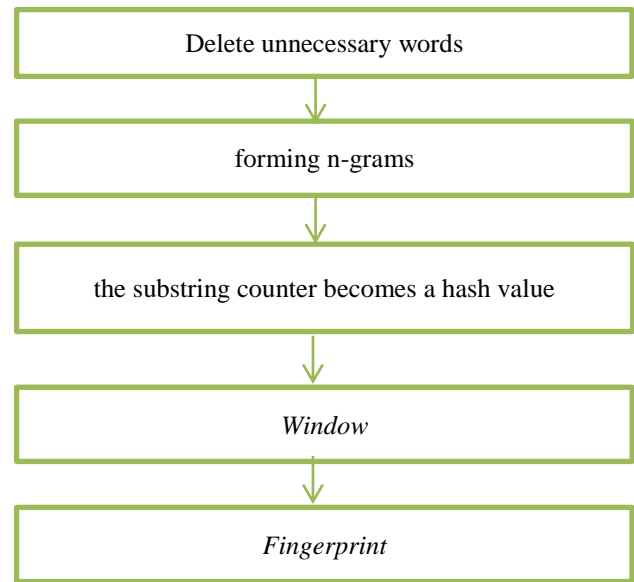


Fig. 1. The concept of Winoing Algorithm.

The concept of the winnowing algorithm in Fig. 3 removes irrelevant characters, forming n-grams of length n, computing hash values, forming window values, and selecting hash values as document fingerprint [6].

C. Dictionary English-Indonesia

Dictionary English-Indonesia is a dictionary to translate from English to Indonesian. This dictionary will be used to detect the action of plagiarism. The act of plagiarism in the writing world can not be separated from the use of foreign sentences that are translated into other languages by the actors of plagiarism to avoid the automatic detection tool.

This act of plagiarism includes the categories of Translated plagiarism and the idea of plagiarism. This category of plagiarism is very difficult to detect by means of plagiarism detection tool whose function is to check word equality on two documents by using fingerprint concept if not combined with Dictionary English-Indonesia. If in two documents there is an English vocabulary then the system will do the translation process into Indonesian before doing the process of checking the similarity in both documents. Dictionary English-Indonesia is highly dependent on databases containing English and Bahasa Indonesia vocabulary.

III. PROPOSED SYSTEM

This section will discuss some of the supporting systems for making plagiarism detection using the dictionary English-Indonesia technique. The system will be designed as follows:

A. System Design Dictionary English-Indonesia

Dictionary English-Indonesia is an important thing to check an act of plagiarism that takes a foreign scientific work in English and changed into the Indonesia language.

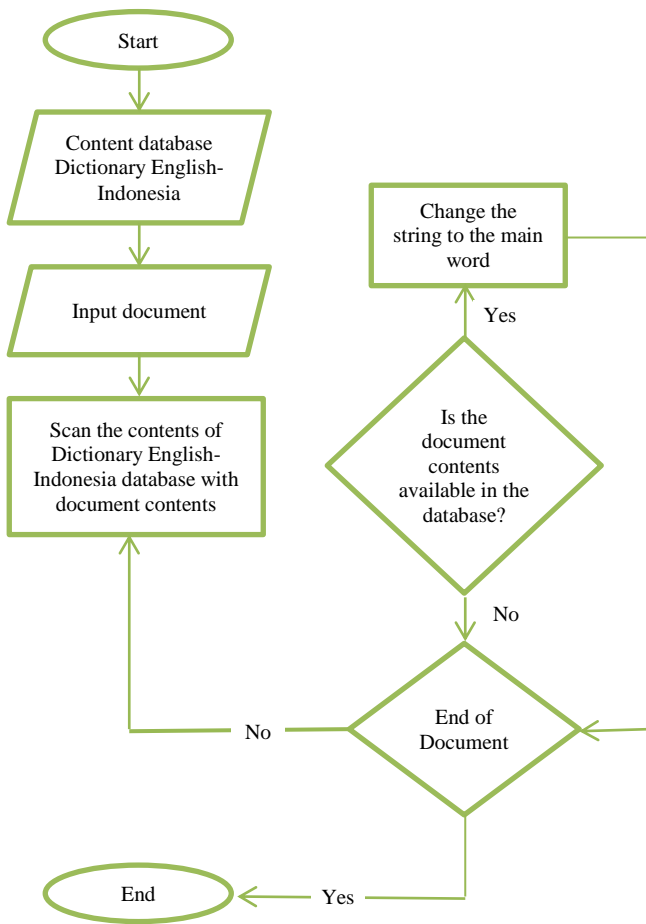


Fig. 2. Process dictionary English-Indonesia.

Based on the flowchart designated in Fig. 2 is a process of dictionary English-Indonesia, the outline is the process by which all the contents of the document will be scanned and will be matched with the existing word in the dictionary English-Indonesia database will then be modified based on the word available on the dictionary database English-Indonesia, if a string of text that match has an English word, the system will convert the text string into Indonesian language already available in the dictionary English-Indonesia database and the word will be included in the Winoxing process and will be reconciled. If a scanned and matched text string does not have an English word, then the process from dictionary English-Indonesia will not be performed. This stage will continue to repeat until the entire process of scan and string matching is complete.

Document I : saya makan apple

Document II : saya makan apel

With the dictionary English-Indonesia then the document has an English word will be changed into the Indonesia language based on dictionary English-Indonesia database.

Document I : saya makan apel

Document II : saya makan apel

After the process is complete then the sentence will go directly to the winnowing algorithm stage to checked again.

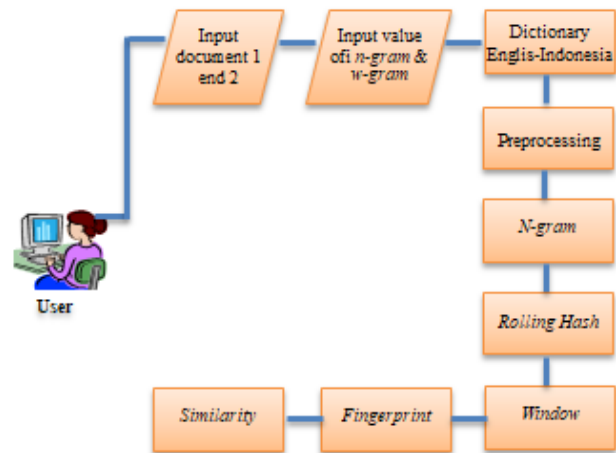


Fig. 3. Flow of Winnowing Algorithm with dictionary English-Indonesian technique for detection of plagiarism.

B. Designing Winnowing Algorithm System with Dictionary English-Indonesia Technique

In this research stages of winnowing algorithm with dictionary English-Indonesia technique. stages can be seen in Fig. 3.

Fig. 3 is a system flow for creating a tool for detecting plagiarism using a winnowing algorithm that serves to process two documents by adding the English-Bahasa Indonesia dictionary technique so that it can examine plagiarism in the category of translated plagiarism and plagiarism of ideas. At the beginning of the user interface must first enter the document to compare and include the document in comparison, then the user also need to specify and then enter the value of n-gram and w-gram (window) to be used as Fingerprint Search of both documents.

The process of generating fingerprints from two documents from the use of the winnowing algorithm with the English-Indonesian dictionary technique for detecting plagiarism shown in Fig. 3 is as follows:

1) Process dictionary English-Indonesia: The system will first perform the scanner process of the contents of two documents with the contents of the dictionary English-Indonesia database. If the user-input document has an English word available in the dictionary English-Indonesia database it will be scanned and will be converted into Indonesian, but if it does not have the English word available on the database it will proceed directly to the winnowing process.

2) Preprocessing: Removes irrelevant characters on a document [12].

a) Case Foling: The process of converting capital letters to lowercase in a document (a-z) [12].

b) Tokenizing: Removes unnecessary characters such as spaces [12].

3) N-gram: Serves to retrieve a token circuit or tangible character pieces along the length of n of a continuously inserted document (continuity) to shift according to the given offsite or end of a word or document [12], [13].

```

private function n_gram($word, $n) {
    $ngrams = array();
    $length = strlen($word);
    for($i = 0; $i < $length; $i++) {
        if($i > ($n - 2)) {
            $ng = '';
            for($j = $n-1; $j >= 0; $j--) {
                $ng .= $word[$i-$j];
            }
            $ngrams[] = $ng;
        }
    }
    return $ngrams;
}

```

Fig. 4. N-gram program.

Fig. 4 is a PHP program from N-gram, the process will take a series of characters along the n-gram value specified by the user

4) Rolling hash is a hashing method that is used to find the hash values of the grams that have been formed and gives the ability to calculate values without repeating the entire string [14]. The hash value is a numerical value formed from the ASCII code [15]. Rolling hash formula can be seen in (1) and (2) above.

5) Window is the main process of winnowing algorithm that serves to categorize hash values that have been formed to produce fingerprint [14].

6) Fingerprint selects the smallest hash value of any given group in the window stage, if there are two or more smallest hash values then select the smallest right hash value.

7) Similarity: measures the similarity in two documents or samples [16]. Similarity to be used is Jaccard Coefficient is usually used to compare documents and calculate the similarity of two objects or documents [15], [17], [18]. Jaccard Coefficient can be seen in (3) [18].

$$Similarity(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} \quad (3)$$

where:

- X = Document 1
- Y = Document 2

The size of a person performing a plagiarism act will be determined by similarity percentage level [5], [6]:

- 0% said the document has nothing in common.
- <15% say documents show plagiarism of small size / have little in common.
- 15-50% said the document is classified as moderate plagiarism.
- 50% mention the document including plagiarism with a large size.
- 100% stated that the document as a whole has something in common.

IV. RESULT

A. Implementation

Here is the implementation of winnowing algorithm with dictionary English-Indonesia technique to detect plagiarism.

Fig. 5. Input Kalimat 1, Kalimat 2, n-gram and window.

Fig. 6. Result dictionary English-Indonesia.

Fig. 5 is the first interface and the user must input the document you want to compare and must input the document that will be the comparison. The determination of n-gram and window values determines the value of similarity.

Fig. 6 shows the process of scanning on documents that have English, the document input process in the second sentence shown in Fig. 5 contains the sentence “i eat apple” and after going through the Dictionary English-Indonesia process “saya makan apel”.

Table I is an irrelevant character removal process of a document or text, whereby spaces are removed and convert capital letters into lowercase.

In Table II above can be seen that the value of n-grams 1 and 2 have similarities, so that hash values 1 and 2 also have the same overall value after doing dictionary English-Indonesia process.

TABLE I. CASE FOLDING AND TOKENIZING

Kalimat 1	sayamakanapel
Kalimat 2	sayamakanapel

TABLE II. RESULT N-GRAM AND HASH

No	N-gram 1	N-gram 2	Hash 1	Hash 2
1	say	say	775	775
2	aya	aya	727	727
3	yam	yam	787	787
4	ama	ama	703	703
5	mak	mak	737	737
6	aka	aka	699	699
7	kan	kan	732	732
8	ana	ana	705	705
9	nap	nap	746	746
10	ape	ape	713	713
11	pel	pel	758	758

TABLE III. RESULT WINDOW

Window 1	Window 2
775 727 787	775 727 787
727 787 703	727 787 703
787 703 737	787 703 737
703 737 699	703 737 699
737 699 732	737 699 732
699 732 705	699 732 705
732 705 746	732 705 746
705 746 713	705 746 713
746 713 758	746 713 758

TABLE IV. FINGERPRINT

Fingerprint 1	727 703 703 699 699 699 705 705 713
Fingerprint 2	727 703 703 699 699 699 705 705 713

Table III is the result of grouping the hash value of a number of w-gram values, from the window process generating fingerprints on each document or sentence by selecting the smallest hash value.

Table IV shows the kalimat 1 and 2 have each fingerprint 9 and have the same value, from this process already visible, the data is a word that has the same meaning after the dictionary English-Indonesia process is done. to prove what percentage of similarity documents can be seen in Fig. 7.

The results from Fig. 7 show that both documents have a 100% similarity after performing the dictionary English-Indonesia process.

B. Trials

The test results of winnowing algorithm with dictionary English-Indonesia technique to detect plagiarism with plagiarism detection test using winnowing algorithm without dictionary English-Indonesia technique with n-gram and w-gram 3 values will be displayed in Table V.

Jumlah Fingerprints kalimat 1 = 9
 Jumlah Fingerprints kalimat 2 = 9
 Union (Gabungan) Fingerprints 1 dan 2 = 18
 Intersection (fingerprints yang sama) = 9
 (Union - Intersection) = 9
 Prosentase Plagiarisme
 Koefisien Jaccard = (Intersection / (Union-Intersection)) * 100
 (9/9) * 100 = 100 %

Fig. 7. Similarity.

TABLE V. TEST RESULT OF PARAMETER

No	Input document	N-gram	W-gram	Similarity %	
				With DEI	Without DEI
1	Kesehatan itu penting	3	3	100%	0%
	Health is important				
2	Perempuan itu sangat cantik dan cerdas	3	3	66,67%	8%
	Dokter itu very beautiful and smart				
3	Plagiarism merupakan salah satu problem dalam dunia academic	3	3	49,41%	33,68%
	Plagiarisme merupakan salah satu permasalahan dalam dunia akademik dan permasalahan bagi bangsa				
4	Algoritma winnowing berfungsi sebagai document fingerprint dengan teknik hashing	3	3	100%	82,19%
	Algoritma winnowing berfungsi sebagai dokumen sidik jari dengan teknik hashing				
5	Dua dokumen tersebut memiliki nilai kesamaan yang sama	3	3	93,48%	59,65%
	Dua dokumen tersebut memiliki nilai similarity yang berbeda				

The above Table V is the result of the parametric trials of the winnowing algorithm with the dictionary English-Indonesia technique and the winnowing algorithm without dictionary English-Indonesia. Can be seen from Table V the influence of winnowing algorithm with dictionary English-Indonesia technique can give high similarity value (high accuracy) than not using Dictionary English-Indonesia technique. In Table V the number 1 data entered in sentences 1 and 2 are different data but have the same meaning. In sentence 1 data entered with the Indonesian language, while in sentence 2 data entered with English. With dictionary English-Indonesia technique then the system will change every word in English into Bahasa Indonesia, after that new system will do re-checking by using winnowing algorithm so that in Table V number 1 yields 100% similarity value. Whereas if detected using only winnowing algorithm alone without dictionary English-Indonesia technique, the system will only check every document in inputkan by user without any change, so that in Table V number 1 yields the value of similarity 0% because the document entered user in sentence 1 and 2 has a very different character though it has the same meaning.

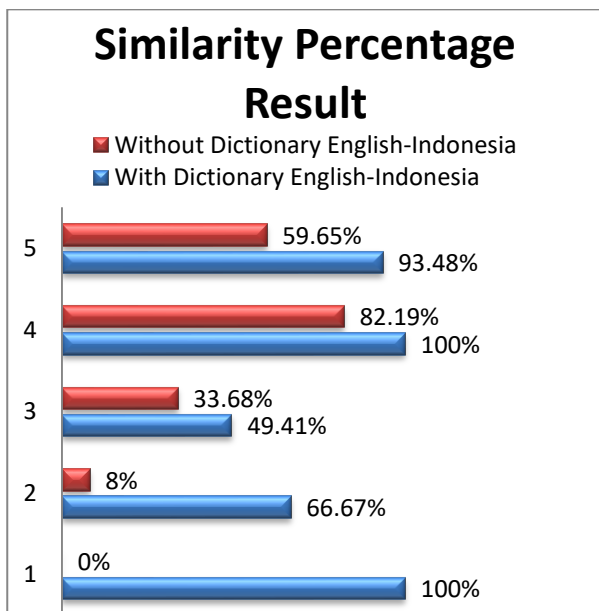


Fig. 8. Similarity percentage result.

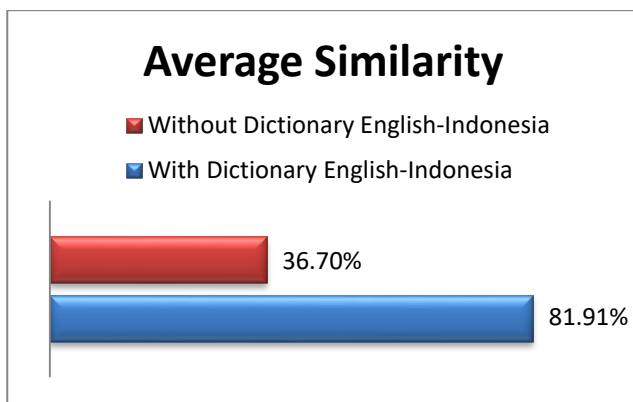


Fig. 9. Average similarity.

Fig. 8 shows the similarity percentage results of the experiments performed in Table V. It can be seen from the Fig. 8 similarity values generated using the dictionary English-Indonesia technique has a high accuracy compared to similarity value without using dictionary English-Indonesia. This is because the algorithm that serves as a document fingerprint can only check every character contained in the contents of the document entered, while the action of plagiarism is not just copy paste, but many categories in the act of plagiarism. one of which is to avoid the plagiarism detection tool usually the perpetrator hides the copied part, summarizes the copied part, paraphrases the part copied by plagiarist style and translates the part copied from the foreign work into the plagiarist language.

In Fig. 9, it can be seen that the use of dictionary English-Indonesian technique has a very big influence on the accuracy of the position of the sentence. Can be seen from Fig. 9, there is a difference of accuracy obtained up to 45.21%. This proves that winnowing algorithm with dictionary English-Indonesia technique is very useful in the accuracy of sentence position.

V. CONCLUSION

From the above research can be understood that the algorithm winnowing with dictionary English-Indonesian technique in plagiarism detection tool is very important to prevent the action of plagiarism in the category of translated plagiarism and idea plagiarism. Plagiarism detection algorithms such as winnowing, rabin-karp algorithm and so on only have a function for pattern matching according to documents or samples that have been inputkn users. The system is unable to check sentences, quotations, or paragraphs that have been paraphrased, hidden, and transransced by actors of plagiarism. Winnowing algorithm with dictionary English-Indonesia technique to detect plagiarism very well is used to minimize the action of plagiarism in the category of translated plagiarism and idea plagiarism, the dictionary English-Indonesia technique also increases the value of similarity between documents.

It is expected that the results of this study can be continued as a follow-up study by researchers themselves and by other researchers. For example, plagiarism detection tool using winnowing algorithm with dictionary Indonesia-English, Arab-English, Mandarin-English and others, to detect plagiarism in category of translated plagiarism and idea plagiarism. In addition, the development of a plagiarism detection tool using the winnowing algorithm can be further developed using the Rabin-karp algorithm, the smith-waterman algorithm, and/or the combination of some of these algorithms.

REFERENCES

- [1] M. Zalnur, "Plagiarisme Di Kalangan Mahasiswa Dalam Membuat Tugas-Tugas Perkuliahan Pada Fakultas Tarbiyah Iain Imam Bonjol Padang," *AL-Ta lim*, vol. 19, p. 55, 2012.
- [2] A. Wibowo, "Mencegah dan Menanggulangi Plagiarisme di Dunia Pendidikan," *Kesmas J. Kesehat. Masy. Nas.*, vol. 6, no. 5, pp. 195–200, 2012.
- [3] Sunardi, A. Yudhana, and I. A. Mukaromah, "Perancangan aplikasi deteksi plagiarisme karya ilmiah menggunakan algoritma winnowing," in *Prosiding SNSebatik*, 2017, vol. 1, no. 1, pp. 27–32.
- [4] N. F. Ulfa, M. Mustikasari, and I. Bastian, "Pendeteksian tingkat similaritas dokumen berbasis web menggunakan algoritma winnowing,"

- in Konferensi Nasional Teknologi Informasi dan Komunikasi (KNASTIK), 2016, pp. 194–203.
- [5] A. Yudhana and A. D. Djayali, “Implementation of Pattern Matching Algorithm for Portable Document Format,” *Int. J. Adv. Comput. Sci. Appl.*, vol. 8, no. 11, pp. 509–512, 2017.
- [6] A. Yudhana, A. D. Djayali, and Sunardi, “Sistem Deteksi Plagiarisme Dokumen Karya Ilmiah dengan Algoritma Pencocokan Pola,” *JURTI*, vol. 1, no. 2, pp. 178–187, 2017.
- [7] R. K. Wibowo and K. Hastuti, “Penerapan Algoritma Winnowing Untuk Mendeteksi Kemiripan Teks pada Tugas Akhir Manusia,” *Techno.COM*, vol. 15, no. 4, pp. 303–311, 2016.
- [8] R. V. Imbar et al., “Implementasi Cosine Similarity dan Algoritma Smith-Waterman untuk Mendeteksi Kemiripan Teks,” *J. Inform.*, vol. 10, no. 1, pp. 31–42, 2014.
- [9] A. T. Wibowo, K. W. Sudarmadi, and A. M. Barmawi, “Comparison between fingerprint and winnowing algorithm to detect plagiarism fraud on Bahasa Indonesia documents,” 2013 *Int. Conf. Inf. Commun. Technol. ICoICT 2013*, no. March, pp. 128–133, 2013.
- [10] G. Wu, M. Zhao, L. Han, and S. Li, “A Fingerprint Feature Extraction Algorithm based on optimal Decision for Text Copy Detection,” *Int. J. Secur. Its Appl.*, vol. 10, no. 11, pp. 67–78, 2016.
- [11] K. T. Tung, N. D. Hung, L. Thi, and M. Hanh, “A Comparison of Algorithms used to measure the Similarity between two documents,” *Int. J. Adv. Res. Comput. Eng. Technol.*, vol. 4, no. 4, pp. 1117–1121, 2015.
- [12] E. Nugroho, “Perancangan Sistem Deteksi Plagiarisme Dokumen Teks Dengan Menggunakan Algoritma Rabin-Karp Skripsi,” in *Fakultas Matematika dan Ilmu Pengetahuan Alam Universitas Brawijaya Malang*, 2011.
- [13] E. A. Lisangan, “Implementasi n-gram Technique dalam Deteksi Plagiarisme pada Tugas Mahasiswa,” *J. Temat.*, vol. 1, no. 2, pp. 24–30, 2013.
- [14] M. Ridho, “Rancang Bangun Aplikasi Pendeteksi Penjiplakan Dokumen Menggunakan Algoritma Biword Winnowing,” in *Teknik Informatika Universitas Islam Negeri SLTAN Syarif Kasim Pekanbaru Riau*, 2013.
- [15] K. Rinatha, “Simple Query Suggestion untuk Pencarian Artikel menggunakan Jaccard Similarity,” *J. Ilm. Rekayasa dan Manaj. Sist. Inf.*, vol. 3, no. 1, pp. 30–34, 2017.
- [16] S. A. Djayali, A. Yudhana, “Pendeteksian Plagiarisme dengan Sistem Pengukuran Similartas pada Dokumen Karya Ilmiah Menggunakan String Matching Rabin-Karp,” in *Cyber Learning & It Computer Karawang*, 2016, vol. 1, no. 1.
- [17] M. Fadelillah, I. Much Ibnu Subroto, and D. Kurniadi, “Sistem Rekomendasi Hasil Pencarian Artikel Menggunakan Metode Jaccard ’ s Coefficient,” *J. Transistor Elektro dan Inform. (TRANSISTOR EI)*, vol. 2, no. 1, pp. 1–14.
- [18] S. Sugiyanto, B. Surarso, A. Sugiharto, and S. A., “Analisa Performa Metode Cosine dan Jacard pada Pengujian Kesamaan Dokumen,” *J. Masy. Inform.*, vol. 5, no. 10, pp. 1–8, 2016.

Multi-Stage Algorithms for Solving a Generalized Capacitated P-median Location Problem

Mohammed EL AMRANI¹, Youssef BENADADA²

Smart Systems Laboratory, Rabat IT Center
ENSIAS, Mohammed V University in Rabat
Rabat, MOROCCO

Abstract—The capacitated p-median location problem is one of the famous problems widely discussed in the literature, but its generalization to a multi-capacity case has not. This generalization, called multi-capacitated location problem, is characterized by allowing facilities to use one of several capacity levels. For this purpose, a predefined list of capacity levels supported by all potential facilities is established. In this paper, we will detail the mathematical formulation and propose a new solving method. We try to construct, indeed, a multi-stage heuristic algorithm that will be called BDF (Biggest Demand First). This new method appears in two approaches: Integrated BDF (IBDF) and Hybridized BDF (HBDF) will be improved by using a local search optimization. A valid lower bound to the optimal solution value is obtained by solving a lagrangian relaxation dual of the exact formulation. Computational results are presented at the end using new instances with higher ratio between the number of customers, facilities and capacity levels or adapted from those of p-median drawn from the literature. The obtained results show that the IBDF is much faster with medium quality solution while HBDF is slower but provides very good solutions close to the optimality.

Keywords—Location; p-median; multi-capacity; heuristic; LNS; lagrangian relaxation; lower bound

I. INTRODUCTION

The location of facilities is a major problem for strategic or tactical decisions. It is much encountered in the industry as well as in the real life. Many interesting applications fields were its direct result, such as network design, telecommunications and customer distribution services. The objective is to propose an optimal assignment of customers to potential facilities subject to a number of constraints such as capacity and budget.

The CPMP (Capacitated P-Median location Problem) is a well-known variant that is characterized by the capacity constraints and the number p of medians predefined initially. It is hugely studied in the literature and constitutes several research studies in combinatorial optimization and operations research fields.

Let $G(N, M, E)$ be a bipartite graph where N represents the set of customers, M is the set of potential sites to install the medians, and E is the set of edges that connect each vertex of N to a node of M. The p-median graph is in the form of a set of clusters, each one is composed of one facility (black triangle) connected to a set of customers (points) or only a closed facility (white triangle) (Fig. 1).

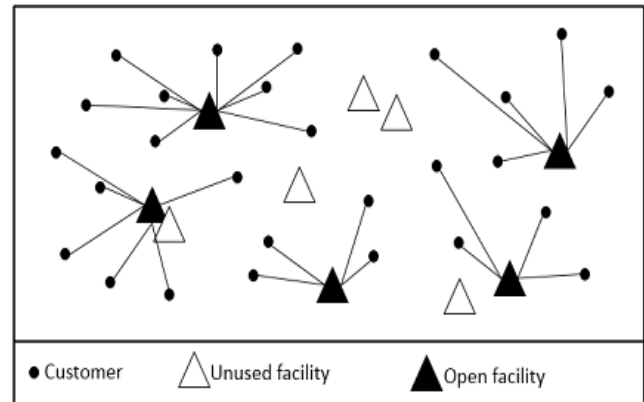


Fig. 1. Graphical representation of basic p-median.

In the industrial fields, service costs generally increase with the capacity used, and as this capacity could exceed in some cases the customers' demands, applying a basic CPMP can present a significant waste of capacity resources. In order to generalize the CPMP for more real problems and for more efficient resolution, we propose in this paper a new variant using different capacity levels. Thus, each facility is prepared to support several capacity levels and to use at the end one level at most. The total of the assigned customers' demands defines the adequate level of capacity, where each level has a corresponding cost. This new variant is called: Budget constraint Multi-Capacitated Location Problem (BMCLP).

The BMCLP's applications can appeared in many industrial areas, such as telecommunications, energy management, and several other fields. The goal of this problem is to minimize the overall cost of assigning customers to facilities. This cost is a multiplication of the unit cost measured by the distance and the demand of the served customer. The total opening costs of facilities is limited by a predefined budget.

The BMCLP is a new variant of location problems family, first time studied by the same authors in [10]. The CPMP problem is NP-complete according to the proof of [11], so its generalization BMCLP is too. Several other variants of CPMP are treated by [18], [3], [1] and [7]. Dynamic location problems are solved by [2], [8] and [6]. Network problems are appeared in [16], [9], [13] and [19]. For the resolution of the problem and its variants, several exact and approximate approaches are tested: the reference [3] applies a cutting plane algorithm based on the Fenchel cuts, references [18], [5] and [4] have chosen to

use the Branch & Price and branch & Bound methods based on Lagrangian relaxation, a resolution with column generation is applied by [4] and references [6], [15] and [12] used different approaches and techniques.

In this work, we will apply the Branch and Cut, a classical resolution method; it is an exact approach that consists of generating cuts at each node of the Branch and Bound tree. Then, we will build a heuristic, more adapted to this location problem variant called BDF (bigger demand first). To improve the solution quality, a local search LNS (Local Neighborhood Search) algorithm, used by [17], will complete the BDF approach.

The BDF, a method in the form of a multi-stage algorithm, is presented in two approaches. Firstly, by using it alone for the solution construction and it will be called IBDF for an Integrated BDF. Secondly by hybridizing it with the application of the branch and cut on a well-defined sub-problem, and it will be called HBDF for a Hybridized BDF.

In order to obtain a valid lower bound to the BMCLP, we use a heuristic procedure to solve the dual problem of our initial formulation. This heuristic procedure combines two different approaches A1 and A2, namely the relaxation of capacity constraints and the decomposition of the problem in two independent sub-problems. Procedure A1 is based on a lagrangian relaxation and sub-gradient optimization, while procedure A2 is based on an independent decomposition starting from the relaxed problem obtained in A1. Indeed, after the relaxation of the capacity constraints and putting them in the objective function, the problem becomes decomposable in two sub-problems, one with variables (x_{ij}) and the second with (y_j^k). These variables are defined in the next section below.

This paper is organized as follows. After the introduction, we discuss, in section two, the formulation of the new BMCLP. The third section is devoted to the solving methods, namely the new heuristic approaches, the Lagrange heuristic and the LNS algorithm. The computation of a valid lower bound is detailed in section four. Finally, computational results are presented in the penultimate section before the conclusion.

II. FORMULATION

The BMCLP is a new variant of capacitated location problem that is characterized by capacity levels, each facility can be used at one level at most. The concept of capacity levels appears in the mathematical formulation with new variables and additional constraints such that each facility must respect the capacity of the level used.

Let $G(N, M, E)$ be a bipartite graph where N represents the set of customers, M is the set of potential sites to install the medians, and E is the set of edges that connect each vertex of N to a node of M .

The BMCLP's graph is in the form of a number of clusters (Fig. 2); each one is composed of one median facility (colorful triangle) connected to a partition of customers set (points). It could also contain only a closed facility (white triangle). Colors represent capacity levels used (k_1, k_2, k_3) and uncolored triangle represents therefore unused facility.

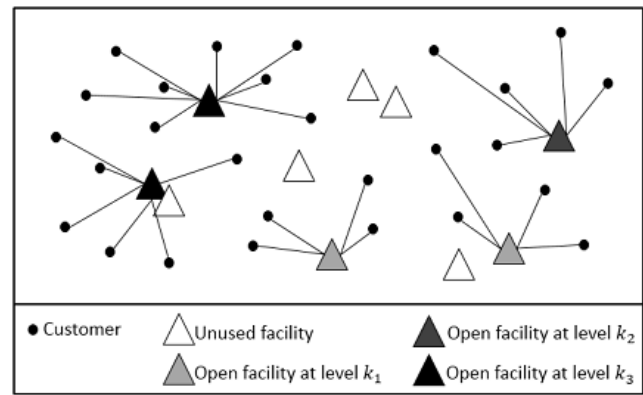


Fig. 2. Graphical representation of BMCLP.

The starting point for building the BMCLP mathematical formulation is the CPMP problem that is defined as follows:

$$\text{Min} \sum_{i \in N} \sum_{j \in M} d_i c_{ij} x_{ij} \quad (1)$$

$$\sum_{j \in M} x_{ij} = 1, \quad i \in N \quad (2)$$

$$\sum_{i \in N} d_i x_{ij} \leq u y_j, \quad j \in M \quad (3)$$

$$x_{ij} \leq y_j, \quad i \in N, j \in M \quad (4)$$

$$\sum_{j \in M} y_j = p \quad (5)$$

$$x_{ij} \in \{0,1\}, \quad i \in N, j \in M \quad (6)$$

$$y_j \in \{0,1\}, \quad j \in M \quad (7)$$

Where d_i is a demand of customer i , c_{ij} represents the assignment cost of customer i to facility j , u is the capacity of medians. We assume that all facilities have the same capacity u . p is the pre-known number of medians to use from the $|M|$ available facilities.

x_{ij} is a binary decision variable which equals to 1 if and only if the customer i is assigned to the facility j .

y_j is also a binary decision variable that is 1 if and only if the facility j is open.

The objective function (1) contains only assignment costs, which can be in the form of transportation costs or response time. In this variant, the opening costs of facilities are determined by the number p in (5). Constraints (2) ensure that each customer is assigned to one and only one median. The constraints (3) require a capacity for each facility. The constraints (4) prohibit the assignment of a customer to a closed facility. The constraints (6) and (7) are the integrality constraints.

The constraints (4) are covered by capacity constraints (3) combined with the integrality ones.

Indeed, let's prove that:

$$\sum_{i \in N} d_i x_{ij} \leq u y_j \Rightarrow x_{ij} \leq y_j \quad \forall i \in N, j \in M$$

Let $j \in M$ and $y_j \in \{0,1\}$

If $y_j = 1$ then

$$x_{ij} \leq y_j, \forall i \in N \quad \text{Because } x_{ij} \in \{0,1\}, \forall i \in N$$

Otherwise $y_j = 0$ then

$$\sum_{i \in N} d_i x_{ij} \leq 0$$

As $\forall i \in N, d_i x_{ij} \geq 0$ So $\forall i \in N, d_i x_{ij} = 0$

Hence $x_{ij} = 0$

And finally, $x_{ij} \leq y_j$

NB. We keep the redundant constraints in the formulation because it increases the efficiency of the Cplex solver. Indeed, Cplex works with branch and cut methods after the relaxation of integrality constraints, these constraints (4) then constitute valid inequalities that decrease the search area, so the algorithm converges more quickly.

The BMCLP is a generalized p-median problem. This generalization concerns facilities that can be operated at several levels of capacity. To do this we must redefine the variable y and use other additional data. We will also have a change in the constraints, the constraint requiring a number p of facilities to open is replaced by a budget constraint limiting the opening costs of factories.

Let K denotes the set of levels, u^k the capacity of level k , f^k the opening cost associated with each level k , and B the limit budget on the sum of facilities opening costs.

Let y_j^k be a binary decision variable that is 1 if and only if the facility j is open and used at the level k .

The mathematical formulation is as follows:

$$\text{Min} \sum_{i \in N} \sum_{j \in M} d_i c_{ij} x_{ij}$$

$$\sum_{j \in M} x_{ij} = 1, \quad i \in N$$

$$\sum_{i \in N} d_i x_{ij} \leq \sum_{k \in K} u^k y_j^k, \quad j \in M \quad (8)$$

$$x_{ij} \leq \sum_{k \in K} y_j^k, \quad i \in N, j \in M \quad (9)$$

$$\sum_{k \in K} y_j^k \leq 1, \quad j \in M \quad (10)$$

$$\sum_{j \in M} \sum_{k \in K} f^k y_j^k \leq B \quad (11)$$

$$x_{ij} \in \{0,1\}, \quad i \in N, j \in M$$

$$y_j^k \in \{0,1\}, \quad j \in M, k \in K \quad (12)$$

In this formulation, the constraints now take into account the multi-capacity concept. Constraints (9) represent valid inequalities that cut the feasible region. Constraints (10) force

the facility to be opened at one level at most. Constraint (11) is used to limit the opening budget of facilities.

The BMCLP problem allowed to modeling more real situations by opening facilities on several capacity levels. However, this generalization also increases the number of constraints and variables and made its resolution more difficult with a solver such as CPLEX, especially for the big size problems. It is for this reason that we seek in the next section a new heuristic approach more suitable to solve the problem. This new approach, called BDF, is in the form of multi-stage algorithm.

III. SOLVING METHODS

In this section, we describe a heuristic procedure (BDF) for finding a good feasible solution to BMCLP that is based on the following two ideas.

- The customer with the highest demand has priority in an assignment to the same facility.
- The nearest facility is favored for the assignment of any customer.

The main idea of this method is to assign the customers, iteratively, to the nearest facility while satisfying the capacity constraint, the budget constraint is a priori ignored. When assigning customers, we give priority, as implies the name of the BDF, to customers with greater demands. However, the real factor considered is not the demand value alone but its multiplication at the distance to facility (factors appeared in objective function). As the assignment of customers to the nearest facility can create a cluster with violated capacity constraints, we keep only the first customers with the biggest demands. The customers not assigned to this facility will be re-assigned to the second nearest facility.

In order to be more accurate, we will create a new priority factor for customers against the same facility.

Let i be a customer, j_1 its nearest first facility and j_2 the second nearest one.

The dissimilarity factor called DF is defined as follows.

$$DF = d_i c_{ij_2} - d_i c_{ij_1}$$

This factor calculates the additional cost when a customer transfers his assignment from a facility to next one. The objective of this method is to minimize the DF factors so as to minimize the overall cost.

The BDF, a method in the form of a multi-stage algorithm, is presented in two approaches:

- By using it alone for the solution construction and it will be called IBDF, for an Integrated BDF.
- By hybridizing it with the application of the branch and cut on a well-defined sub-problem, and it will be called HBDF for a Hybridized BDF.

We use the notation IBDF0 and HBDF0 for IBDF and HBDF respectively without applying the LH and LNS algorithms.

Both approaches provide solutions that are generally not feasible; therefore, a Lagrange heuristic (LH) will be applied for the solution feasibility. The approaches will be improved later by using LNS, (Large neighborhood search) a local optimization method.

Given that the method consists of several sub-methods, we will start by developing them before establishing both algorithms.

A. Branch and Cut (B&C)

The Branch and cut is a combination of two algorithms into one, namely the Branch and Bound and cutting planes. We will not be interested in this algorithm because we use it implicitly through the Cplex solver. It will be still used for comparison with the solution obtained by our method for smaller instances.

B. Integrated Biggest Demand First (IBDF0)

The IBDF0 is an integrated method that can find for several instances very good solutions and sometimes optimal ones. However, this heuristic method remains unreliable and does not guarantee the feasibility. For this reason, we propose to use Lagrange heuristic to make the solution feasible.

IBDF0 Algorithm

- 1) For each customer,
 - Sorting facilities by distances to this customer in ascending order.
- 2) $p \leftarrow 1$; $UC \leftarrow N$ (UC for unassigned customers set)
- 3) Assigning each customer, from the UC set, to its p^{th} nearest facility; $p \leftarrow p+1$
- 4) For each facility,
 - Sorting customers according to theirs DF in descending order.
 - Keeping the maximum customers without exceeding the higher capacity level. The remaining customers will be put in UC set.
 - Updating the capacity resources; the new capacity is the one available after satisfying the customers' demands.
- 5) If $UC = \emptyset$ Then break
 - Else go to 3-

C. Hybridized Biggest Demand First (HBDF0)

The HBDF0 is the first iteration of the IBDF0 method hybridized with the B&C, its principle is to reduce the problem size. Indeed, it proposes assignments for certain customers considered to have important demands. After the application of this method, we set the corresponding variables to the assigned customers and start the execution of the B&C on the unassigned customers' sub-problem.

- 1) For each customer,
 - Sorting facilities by distances to this customer in ascending order.
- 2) $UC \leftarrow N$ (UC for unassigned customers)
- 3) Assigning each customer to its nearest facility.
- 4) For each facility,

- Sorting customers according to theirs DF in descending order.
- Keeping the maximum customers without exceeding the higher capacity level. The remaining customers will be put in UC set.
- Updating the capacity resources; the new capacity is the one available after satisfying the customers' demands.

5) Updating the budget; the new budget is the one available after subtracting the corresponding costs at all levels used in all open facilities.

6) If $UC = \emptyset$ then break

Else applying the formulation of the problem with the new updated data (UC, capacities, budget ...), using Branch and Cut.

D. Lagrange Heuristic (LH)

The resolution, with violated constraints, gives generally unfeasible solution. To have a feasible solution, we need to apply some existing heuristics. We decided to use the Lagrange heuristic, which consists of setting some variables of the problem, and re-solve the initial problem.

Let's set the variables y_j^k to their values obtained from the unfeasible solution and reconstruct the problem.

The new formulation after setting y_j^k is as follows:

$$\begin{aligned}
 & \text{Min} \sum_{i \in N} \sum_{j \in M} d_i c_{ij} x_{ij} \\
 & \sum_{j \in M} x_{ij} = 1, \quad i \in N \\
 & \sum_{i \in N} d_i x_{ij} \leq K_j, \quad j \in M \\
 & x_{ij} \in \{0,1\}, \quad i \in N, j \in M
 \end{aligned} \tag{13}$$

Where $K_j = \sum_{k \in K} u^k y_j^k$ is a dependent constant on facility j.

The previous problem is linear and contains a reduced number of constraints and variables. It is in the form of knapsack problem with additional demand constraints. Therefore, the problem can be easily solved with Cplex. The solution of this problem is feasible but approximate.

E. Large Neighborhood Search (LNS)

To improve the solution obtained, we propose to use one of the local optimization methods. Among those that have demonstrated their effectiveness in combinatorial optimization and particularly in location and transportation problems, we mention LNS (Large Neighborhood Search). This method has the advantage of ensuring optimality for instances of small or medium size when the selected sample is also the problem studied.

LNS is a meta-heuristic used for local optimization. From a first solution, the search algorithm will try to improve it by successive samplings in its neighborhood. At each iteration, one or more clusters (facility and connected customers) are

deconstructed to obtain a sub-problem that is supposed to be easy to solve. Then we solve it using Branch and Cut algorithm to rebuild the solution. A stochastic element defines the sample used for each application of the method.

LNS Algorithm

Repeat until the stopping criterion.

1) Randomly selecting a part of the solution (i.e a set of clusters), this part will constitute a sub-problem easy to solve using Cplex.

2) Destroying the clusters of the sub-problem under consideration.

3) Executing Cplex and recovering the obtained partial solution.

4) Integrating the solution of the sub-problem found to the solution of the initial problem.

5) If during 10 iterations, no significant improvement is recorded, we declare the stopping criterion.

BDF Algorithms

After having defined all steps of the BDF method, we can then implement the following two algorithms:

IBDF:

1) *IBDF0*

2) *If the obtained solution is feasible*

- LSN
- Exit

3) *Else*

- LH
- LNS
- Exit

HBDF:

1) *HBDF0*

2) *If the obtained solution is feasible*

- LSN
- Exit

3) *Else*

- LH
- LNS
- Exit

IV. COMPUTATION OF THE LOWER BOUND

In this section, we will present a method based on lagrangian relaxation to determine a good lower bound. The calculation of this will allow us to evaluate the quality of the solution obtained. It can be noted that the formulation contains three different constraints' blocks; the first one with variables x_{ij} (constraints 2), the second with variables y_j^k (constraints 10 and 11) and the third which contains a combination of both (constraints 8 and 9). Eliminating the third block, the formulation becomes decomposable into two independent sub-

problems, one with variable x_{ij} and the second with variables y_j^k .

As we have mentioned above, constraints (9) are facultative and its violation has no impact on solution feasibility. In order to get the latter decomposition, we will forget constraints (9) and relax constraints (8).

The new relaxed problem is as follows:

$$\begin{aligned} \text{Min} \quad & \sum_{i \in N} \sum_{j \in M} d_i(c_{ij} - \lambda_j)x_{ij} + \sum_{j \in M} \sum_{k \in K} \lambda_j u^k y_j^k \quad (14) \\ & \sum_{j \in M} x_{ij} = 1, \quad i \in N \\ & \sum_{k \in K} y_j^k \leq 1, \quad j \in M \\ & \sum_{j \in M} \sum_{k \in K} f^k y_j^k \leq B \\ & x_{ij} \in \{0,1\}, \quad i \in N, j \in M \\ & y_j^k \in \{0,1\}, \quad j \in M, k \in K \end{aligned}$$

This problem is decomposable into two sub-problems:

Problem x:

$$\begin{aligned} \text{Min} \quad & \sum_{i \in N} \sum_{j \in M} d_i(c_{ij} - \lambda_j)x_{ij} \quad (15) \\ & \sum_{j \in M} x_{ij} = 1, \quad i \in N \\ & x_{ij} \in \{0,1\}, \quad i \in N, j \in M \end{aligned}$$

Problem y:

$$\begin{aligned} \text{Min} \quad & \sum_{j \in M} \sum_{k \in K} \lambda_j u^k y_j^k \quad (16) \\ & \sum_{k \in K} y_j^k \leq 1, \quad j \in M \\ & \sum_{j \in M} \sum_{k \in K} f^k y_j^k \leq B \\ & y_j^k \in \{0,1\}, \quad j \in M, k \in K \end{aligned}$$

These two problems are independent and can be solved in parallel using sub-gradient algorithm. However, we propose that at each iteration, we start with the problem x, we obtain the variables' values x which give information on the facilities that must be open then we add it as constraint in problem y before solving it. At each iteration of the sub-gradient algorithm, we can have a solution that present a lower bound for our problem; this bound will be improved from one iteration to another.

V. COMPUTATIONAL RESULTS

The BMCLP is a new problem that has not been found in the literature. Therefore, we cannot find instances for the test or for comparison. For this reason, we decide within this research to create instances using semi-random values based on justified choices and whose difficulty is measurable. We will also use p-median instances adapted to our problem to complete the calculation tests.

We turn both algorithms on an i7-2600 CPU @ 3.40 GHz machine with 8GB RAM. We use the programming language java version 7 and version 12.3 of Cplex.

The test set consists of five classes of instances representing five levels of difficulty (easy, medium, difficult, very difficult and complex). The difficulty of these instances is based on the size of the problem, which is generally measured by the number of customers, facilities, and capacity levels. However these two last numbers have a small impact on the problem's size. Each level of difficulty contains several test instances. Other difficulty factors are taken into account, namely the dispersion of customers against facilities and the available resources. Experience shows that the difficulty of the problem varies in proportion to the variance of customers' distances and their demands. At the same time, it varies inversely with the budget allocated for opening facilities and their capacity levels. Thus, by increasing the difficulty, while

keeping the feasibility, we multiply the number of iterations necessary to find the optimal solution.

These are the parameters used in the following result's table:

- LD: Level of Difficulty
- NC: Number of Customers
- NF: Number of Facilities
- NL: Number of Levels
- LB: Lower Bound
- Obj: Objective function value
- CPU: Execution time
- GAB: = Min(GAB1,GAB2)
 $GAB1 = (BDF - B\&C) / BDF$
 $GAB2 = (BDF - LB) / BDF$

The following table shows the different instances used and the execution results of the IBDH and HBDF methods as well as the lower bound (Table I):

TABLE I. COMPUTATIONAL RESULTS TABLE

Insta nce	LD	NC	NF	NL	Branch & Cut		LB	IBDF			HBDF		
					Obj	CPU(s)		Obj	CPU(s)	GAP	Obj	CPU(s)	GAP
F1	Easy	10	3	2	184	0.02	183	184	0.001	0,00%	184	1.521	0.00%
F2		10	5	2	230	1.57	228	230	0.001	0,00%	230	2.932	0.00%
F3		20	5	3	430	3.88	426	430	0.006	0,00%	430	5.875	0.00%
F4		30	8	3	372	8.75	361	386	0.011	3,63%	372	9.104	0.00%
M1	Medium	50	4	3	2120	159.45	2073	2423	1.013	12,51%	2158	12.187	1.76%
M2		50	6	4	8429	132.31	8269	8697	5.012	3,08%	8501	16.345	0.85%
M3		70	6	4	5915	234.71	5758	7079	4.004	16,44%	6002	29.297	1.45%
D1	Difficult	100	10	5	12602	5991.1	12494	13923	11.105	9,49%	12802	136.548	1.56%
D2		100	15	5	-	-	14138	15572	12.435	9,21%	14435	24.364	2.06%
D3		200	15	8	-	-	835569	863935	13.432	3,28%	854112	51.784	2.17%
V1	Very Difficult	300	25	10	-	-	318418	345808	12.726	7,92%	319853	36.273	0.45%
V2		300	30	10	-	-	194222	227616	15.762	14,67%	202615	77.238	4.14%
V3		402	30	12	-	-	369834	390885	21.253	5,39%	387483	83.684	4.55%
V4		402	40	12	-	-	405420	432915	23.932	6,35%	406778	91.105	0.33%
C1	Complex	500	50	4	-	-	259809	299221	37.317	13,17%	261644	117.634	0.70%
C2		1000	100	4	-	-	430119	471794	51.265	8,83%	451432	136.721	4.72%
C3		3038	600	10	-	-	94725	105740	56.216	10,42%	99801	294.364	5.09%
C4		3038	700	10	-	-	85082	116380	54.823	26,89%	89801	5011.784	5.25%
C5		3038	1000	10	-	-	297208	347564	47.823	14,49%	298053	3456.273	0.28%

From the numerical results, we note that the BDF method has yielded good results in most cases and across both approaches. However, we find that IBDF is much faster but less effective, while HBDF gives very good results. Although it is slower, it still works for a reasonable time. The local search LNS, improving the quality of the solution, allows in some cases and for small instances to reach the optimal solution. This is justified by the fact that the selected sub-problem to destroy coincides with global problem.

VI. CONCLUSION

In this paper we introduced the budget constraint multi-capacity location problem. We proposed the BDF method in the integrated BDF and Hybridized BDF approaches, Lagrange heuristic to ensure the solution feasibility and for the improvement step, the local LNS search. To evaluate the solution quality, a valid lower bound to the optimal solution value is obtained by solving a lagrangian relaxation dual problem.

The solution achieved might not be an optimal BMCLP solution, however the branch and cut, used for small instances, allows to estimate its maximum distance from optimality. For large instances, a valid lower bound is calculated. Computational tests on problems adapted from those proposed in the literature and on new test problems with large dimensions show the effectiveness of the proposed multi-stage algorithm. As a perspective we propose to test other solving methods that ensure the optimality and work on much large scale problems.

REFERENCES

- [1] Baldacci R., Hadjiconstantinou E., Vittorio M. & Mingozzi A., A new method for solving capacitated location problems based on a set partitioning approach, *Computers & operations research*, 29, 365-386 (2001)
- [2] Behmardi B. & Shiwoo L., Dynamic multi-commodity capacitated facility location problem in supply chain. 2008 Industrial engineering research conference, 1914-1919 (2008)
- [3] Boccia M., Sforza A., Sterle C. & Vasilyev I., A cut and branch approach for the capacitated p-median problem based on Fenchel cutting planes, *Journal of mathematical modelling and algorithms*, 7, 43-58 (2008)
- [4] Ceselli A., Two exact algorithms for capacitated p-median problem, *Springer-verlag berlin heidelberg*, 1, 319-340 (2003)
- [5] Ceselli A. & Giovanni R., A branch-and-price algorithm for the capacitated p-median problem, *Networks an international journal*, 45, 125-142 (2005)
- [6] Da Gama F. S. & Captivo M. E., A heuristic approach for the discrete dynamic location problem, *Location Science*, 6, 211-223 (1998)
- [7] Dantrakul S., Likasiri C. & Pongvuthithum R., Applied p-median and p-center algorithms for facility location problems, *Expert Systems with Applications*, 41, 3596-3604 (2014)
- [8] Dias J., Captivo M. E. & Climaco J., Capacitated dynamic location problems with opening, closure and reopening of facilities, *IMA journal of management mathematics*, 17,317-348 (2006)
- [9] Eberyab J., Krishnamoorthya M., Ernsta A. & Bolandb N, The capacitated multiple allocation hub location problem: Formulations and algorithms, *European journal of operational research*, 120, 614-631 (2000)
- [10] El Amrani M., Benadada Y. & Gendron B., Generalization of capacitated p-median location problem: modeling and resolution, *International IEEE Conference of logistics operations management*, 1-6 (2016)
- [11] Garey M. R. & Johnson D. S., *Computers and intractability, a guide to the theory of NP-completeness*, 338, The ACM digital library, Philadelphia (1979)
- [12] Klose A., A lagrangean relax-and-cut approach for the two-stage capacitated facility location problem, *European journal of operational research*, 126, 408-421 (2000)
- [13] Ling-Chieh K. & Watson-GANDY, C. D. T., An approximation algorithm for a competitive facility location problem with network effects, *European Journal of Operational Research*, 267, 176-186 (2018)
- [14] Lorena L. A. N. & Senne E. L. F., A column generation approach to capacitated p-median problems, 31, 863-876 (2003) <http://www.lac.inpe.br/lorena/instancias.html>
- [15] Mohammad S. K., Haldun S. & Cem I., A column generation approach for the location-routing problem with time windows, *Computers & Operations Research*, 90, 249-263 (2018)
- [16] Oded B., Jorg K. & Dmitry K., On covering location problems on networks with edge demand, *Computers & Operations Research*, 74, 214-227 (2016)
- [17] Ropke S., Pisinger D. An adaptive large neighborhood search heuristic for the pickup and delivery problem with time windows. (2004). Technical report, Department of Computer Science, University of Copenhagen, Denmark
- [18] Senne E. L. F. & Pereira M. A., A branch-and-price approach to p-median location problems, *Computers & operations research*, 32, 1655-1664 (2004)
- [19] Shu J., Ma Q. & Li S., Integrated location and two-echelon inventory network design under uncertainty, *Annals of operations research*, 181, 233-247 (2010)

Communicator for Hearing-Impaired Persons using Pakistan Sign Language (PSL)

Muhammad Wasim

Department of Computer Science
Usman Institute of Technology
Karachi, Pakistan

Adnan Ahmed Siddiqui

Department of Computer Science
Hamdard University
Karachi, Pakistan

Abdulbasit Shaikh

Department of Computer Science
Institute of Business Administration (IBA), Karachi,
Pakistan

Lubaid Ahmed, Syed Faisal Ali, Fauzan Saeed

Department of Computer Science
Usman Institute of Technology
Karachi, Pakistan

Abstract—Communication with a hearing-impaired individual is a big challenge for a normal person. Hearing-impaired people use hand gesture language (sign language) to communicate with each other, which is not easy to understand by a normal person because he/she is not trained to understand sign language. This communication gap between a hearing-impaired and a normal person created a big problem for hearing-impaired individuals during their shopping, hospitalization, at their schools and homes. Especially in case of emergency, it is very difficult to understand the statement of a hearing-impaired one who uses sign language. In the last few years researchers and developers from all over the world presented different ideas and works to solve this problem but no such solution is available to resolve this issue and can create two-way communication between hearing-impaired and normal persons. This paper presented a detail description about a two-way communication system based on Pakistan Sign Language (PSL). This duplex system is developed through conversion from the text in simple English into hand gestures and vice versa. However, conversion from hand gestures is available not only in text but also with voice providing more convenience to normal person. Main objective is to facilitate a large population and making hearing-impaired persons, the vital part of our civilization. A normal person can enter the text (sentence) in application, after the checking of spelling and grammar, the text is divided into tokens and sub-tokens. A token is a gesture against each word of the text while sub-tokens are the gestures of each character of the words. The combination of tokens created the gestures of text. On the other hand when gestures were input in to the application, using image processing technique, the nature of hand gesture were recognized and converted into corresponding text or voice.

Keywords—Communicator; hearing-impaired; Pakistan Sign Language (PSL); hand gesture; special person; token

I. INTRODUCTION

Deaf person (hearing-impaired) uses hand gestures as a basic language (sign language) for the purpose of communication with normal-hearing persons. Normally it is difficult to understand this sign language for hearing-persons without proper training and it creates a big gap between hearing-impaired and normal-hearing persons. The proposed

application is a dual mode application that can be used as an easy and proper communication between them. This duplex system is developed through conversion from the text in simple English into hand gestures and vice versa. However, conversion from hand gestures is available not only in text but also with voice providing more convenience to normal person. Main objective is to facilitate a large population and making special persons, the integral part of the society. The system “communicator” is based on Pakistan Sign Language (PSL). In this application a normal person can enter the text (sentences) in application, after the checking of spelling and grammar, the text is divided into tokens and sub-tokens. A token is a gesture against each word of the text while sub-tokens are the gestures of each character of the words. The combination of tokens created the gestures of text. On the other hand, when gestures were input in to the application, using image processing technique the nature of hand gesture were recognized and converted into corresponding text or voice.

According to R&D report of Pakistan (2012) the estimated population size of Pakistan is approximately 180.7 million [1]. Due to this very highly rate of population growth in Pakistan there are several issues and health is one of the most important areas of concern [2], [3]. Another fact which is published in R&D report is high number of hearing-impaired individuals in urban and rural areas of Pakistan [4], [5] as mentioned in Table I.

According to given published data the province Punjab having large number of hearing-impaired persons as compare to other provinces (Fig. 1).

TABLE I. HEARING-IMPAIRED INDIVIDUALS IN PAKISTAN

District Name	Total Population	No. of Hearing-impaired Persons
Sindh	14,32,148	89,4,11
Punjab	28,16,7,95	2,33,7,37
Khyber Pakhtunkhwa	56,02,65	42,8,94
Balochistan	21,03,91	11,1,37

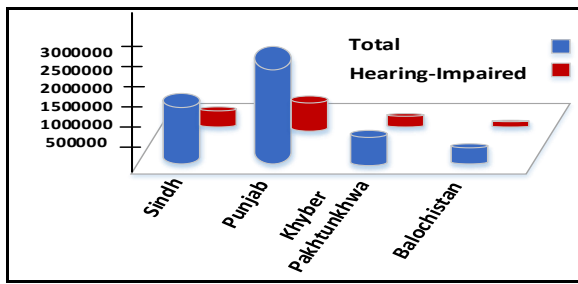


Fig. 1. Hearing-impaired people in Pakistan.

The main purpose to present this work is to facilitate hearing-impaired persons [6] and reduce the communication gap between hearing-impaired and normal individuals of the society [7]. According to published work, it was recorded that minimum duration to identify a gestures is 6-20 fps [8].

It is a fact of society that hearing-impaired persons not feeling comfort and cannot socialize among normal persons even with other persons of family as well. According to published work [9] around 300 million persons are hearing-impaired in the world and mostly are not well-trained with sign language and it creates a big gap of communication [4]. Authors in [10] also discussed the issues related to understanding of sign language. In [11], [12], some work related to electronic device, which can be used as an interpreter between normal and hearing-impaired persons.

Most of previous works were single mode that can translate hand gestures into text, but the proposed work is a dual mode, desktop based application.

II. SIGN LANGUAGE

Each country of the world having its own gesture language (Sign language) and therefore there are many development and research work have been reported in this area of research [12]-[14]. This Sign language is different in different country and based on certain gestures [15]. The gestures patterns are based on different arrangement of hand and fingers. In this proposed work these gestures were recognized first and then converted into certain text which can be understood by normal persons who cannot understand sign language.

III. LITERATURE REVIEW

In [16]-[18] authors discussed the importance and need of gesture based communication. They also discussed the way to improve the techniques. In [19]-[22], authors discussed the association of facial behavior with gestures, conversion of gestures into text for understanding, efficiency and accuracy of gesture recognition and web based application for distance learning and communication.

IV. PAKISTAN SIGN LANGUAGE (PSL) AND COMMUNICATOR

In Pakistan there are number of hearing-impaired institutions where PSL uses as a standard language for hearing-impaired persons. PSL is a combination of gestures patterns consist of alphabets, words and sentences [15]. PSL is based on single and double handed gestures. PSL is used for the purpose of communication among hearing-impaired

individuals and now it can be used as an interpreter between hearing-impaired and normal hearing persons. The presented work is a dual mode interpreter. It can convert PSL into text and for hearing-impaired persons text can be converted into gestures as well. PSL deals with both English and Urdu versions but proposed work is related with English conversion. Using communicator it might be possible that hearing-impaired and normal individuals can communicate with each other's without any hesitation. Fig. 2 showed basic alphabets symbols of Pakistan Sign Language.

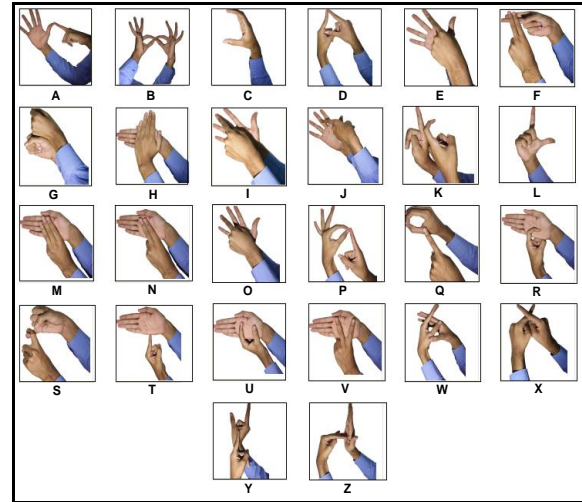


Fig. 2. Basic alphabets of PSL.

The communicator is a desktop application; a hearing-impaired person can input gestures using PSL into the system, which is converted into text/voice for normal hearing persons who did not understand language of deaf persons. On the other hand a normal hearing person can input text or voice into the system, which is translate into gestures according to PSL and easily can understand by hearing-impaired individuals.

V. METHODOLOGY

The methodology of the work is divided into two phases. In first phase authors discussed the way to translate test to gesture conversion in application and in second phase gestures to text or voice conversion has been discussed.

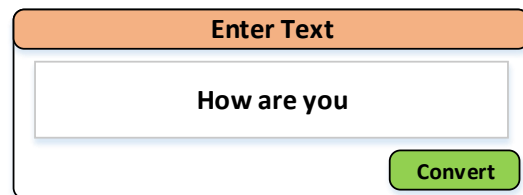
A. Text to Gestures Conversions

Text or voice can be input by the normal hearing person, which is converted into gestures using designed application.

Code:

```
var mes = document.getElementById('myTextArea').value;
```

GUI:

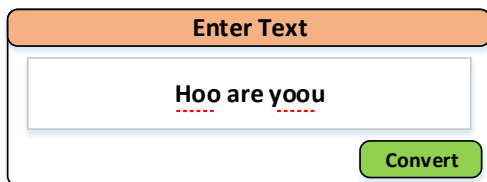


Before conversion application, first check spelling and grammar of the provided input text.

Code:

```
<script type='text/javascript' src='/JavaScriptSpellCheck/include.js'></script>
<script type='text/javascript'>$Spelling_SpellCheckAsYouType("myTextArea")</script>
```

GUI:



Once spell and grammar check application created token and sub-token against each character and words to create proper sentences.

Code:

```
var words = req.params.file.split(" ");
var Characters = req.params.file.split("");
```

Gestures "B & How" is mentioned in Fig. 3.

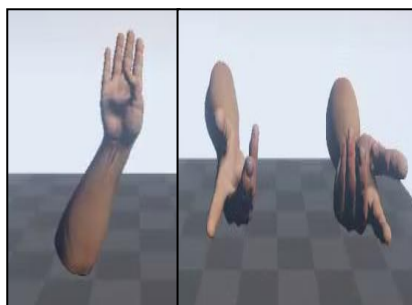


Fig. 3. Gestures of 'B' and 'How'.

Fig. 4 describes the flow-diagram of text to gesture conversion.

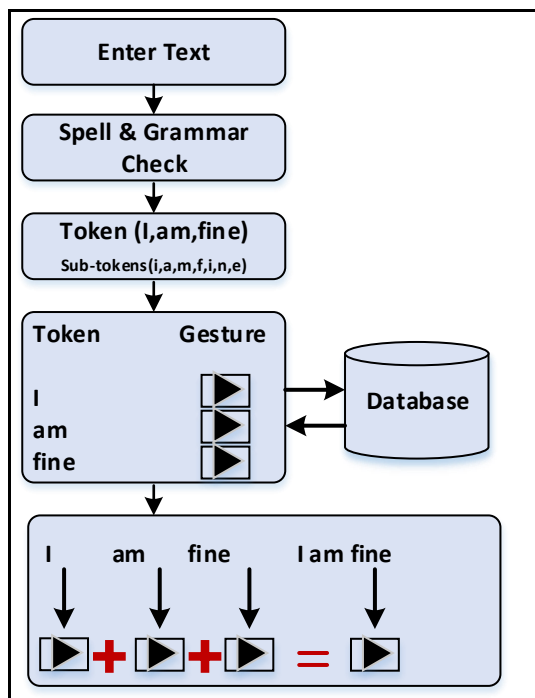


Fig. 4. Text-to-gestures conversion.

VI. GESTURES TO TEXT/VOICE COMMUNICATION

Using Leap motion controller hand gestures of hearing-impaired persons can be input into the system. Leap motion controller provide hand gesture coordinate values which helps to developed algorithm to recognized gestures and converted into related text or voice as an output. The designed system is easy to portable and more accurate to recognized hand gestures. Fig. 5 describes the way conversion between gestures (input) to text (output) and the key stages and functionalities are shown in Fig. 6.

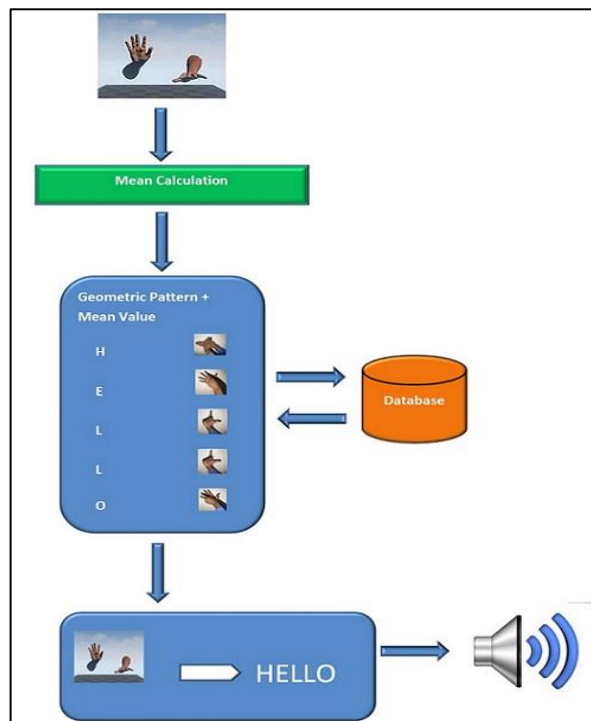


Fig. 5. Conversion of gestures into text/voice.

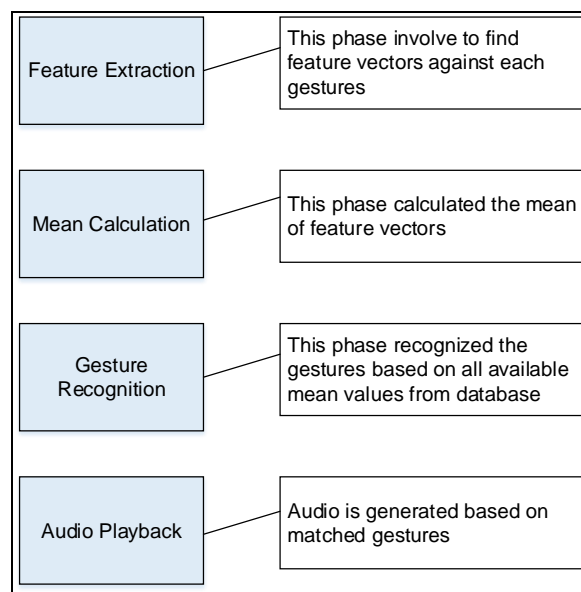


Fig. 6. Key stages and functionalities.

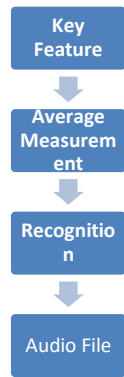


Fig. 7. System modules.

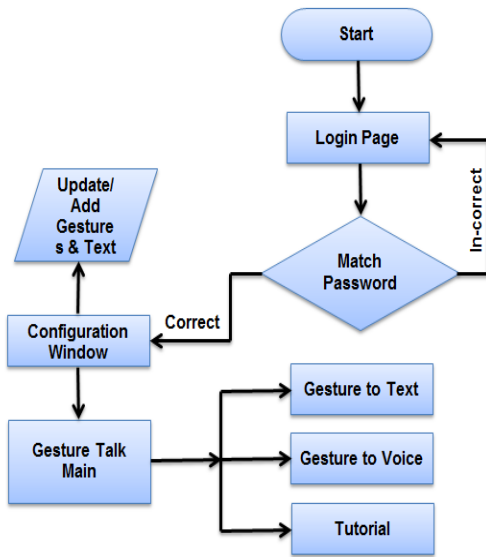


Fig. 8. System flow.

Fig. 7 and 8 describes the key phases and system flow.

VII. KEY MODULES OF SOFTWARE

In this developed system hand gestures and text/voice can be input into the system. In half portion of screen text is mentioned and in remaining part of screen corresponding gestures are shown. In this application a user need to register in the system to use and record gestures. An admin panel is used to add, update and delete ant gestures, text or audio file (Fig. 9).

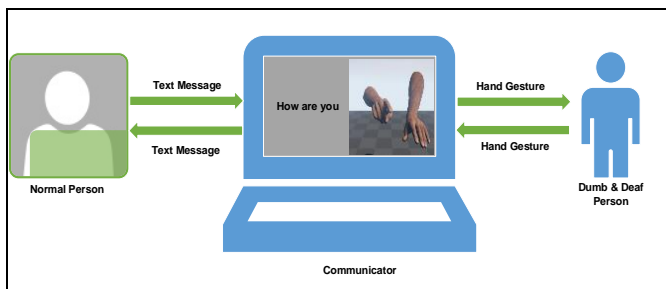


Fig. 9. Application screen shot.

VIII. RESULTS

Table I shows the results of communicator and rate of acceptance of alphabets. To achieve high accuracy all test runs were repeated 10 times for 100 individuals (deaf-persons).

Fig. 10 shows the graphical view of Table II.

Table III shows the results of gestures and text conversions.

TABLE II. ALPHABETS AND RATE OF ACCEPTANCE

Person	Alphabets	Issues with Alphabet	Acceptance Rate (%)
1.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
2.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
3.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
4.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	Issue with I & J	92.31
5.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
6.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
7.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
8.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	Issue with J	96.15
9.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100
10.	ABCDEFGHIJKLMNOPQRSTU VWXYZ	No Issue	100

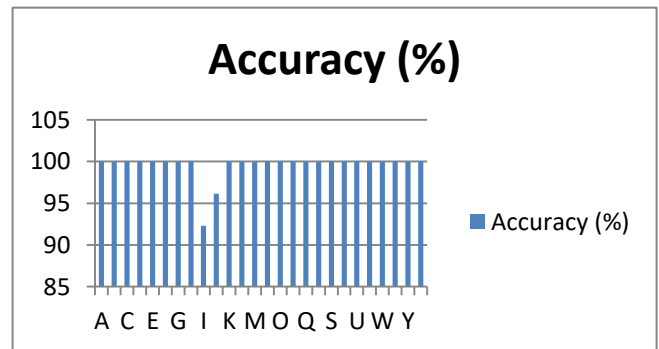


Fig. 10. Graphical representation of Table II.

TABLE III. RESULTS OF TEXT AND GESTURES CONVERSIONS

Input Text	Gestures	Persons										Acceptance Rate (%)	
		1	2	3	4	5	6	7	8	9	10		
Hello		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100
How are you		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100
I am fine		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100
No		OK	OK	OK	×	OK	OK	OK	OK	×	OK	80	
No thank you		OK	OK	OK	OK	×	OK	OK	OK	×	×	70	
Thank you		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100	
What is your name		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100	
Where are you		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100	
Where do you live		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100	
Yes		OK	OK	OK	OK	OK	OK	OK	OK	OK	OK	100	

IX. CONCLUSION

The developed application was tested on 100 deaf individuals in Image Processing Research Lab (IPRL) at Usman Institute of Technology. Each testing repeated 10 times for the purpose of accuracy. It is observed that the system is working with high accuracy in term of recognized

gestures and text/voice conversion. The proposed system is a desktop based application, which is using as an interpreter between normal hearing and hearing-impaired person. In next phase of development the authors of this research work are working to develop an android based application, which will be easy to use.

ACKNOWLEDGMENT

The authors would like to thanks Mr. Hassan Izhar, Muhammad Mubashir Khan, Hafiz Abdullah Qamar and Ashar Jamal Baig to support in programming. Special thanks to Mr. Shahroz Shamim to record gestures of alphabets of PSL in IPRL.

REFERENCES

- [1] Persons with Disabilities (PWDS) Statistics in Pakistan. Research & Development Department of HHRD, Islamabad, Pakistan, 2012.
- [2] Wasim, M., Kamal, S.A., & Shaikh, A. A Security System Employing Edge-Based Rasterstereography. *International Journal of Biology & Biotechnology*, 10(4), 483-501, 2013.
- [3] Waqar, K. Disability: Situation in Pakistan, Right to Education Pakistan Article 25A. Aga Khan University, 2014
- [4] Alvi, A. K., Azhar, M. Y. B., Usman, M., Mumtaz, S., Rafiq, S., Rehman, R. U., & Ahmed, I. Pakistan Sign Language Recognition using Statistical Template Matching. *International Journal of Information Technology*, 1(1), 1-12, 2004.
- [5] Sami, M., Ahmed, H., Wahid, A., Siraj, U., Ahmed, F., Shahid, S., & Ali Shah, S. I. Pose Recognition using Cross Correlation for Static Images of Urdu Sign Language. In *Robotics and Emerging Allied Technologies in Engineering (iCREATE) International Conference, IEEE. 200-204, 2014.*
- [6] Kuroki, K., Zhou, Y., Cheng, Z., Lu, Z., Zhou, Y., & Jing, L. A Remote Conversation Support System for Hearing-impaired-mute Persons Based on Bimanual Gestures Recognition Using Finger-worn Devices. *Workshop on Sensing Systems and Applications Using Wrist Worn Smart Devices, IEEE. 574-578, 2015.*
- [7] Ahire, P. G., Tilekar, K. B., Jawake, T.A., & Warale, P.B. Two Way Communicator Between Hearing-impaired and Dumb People and Normal People. *International Conference on Computing Communication Control and Automation, IEEE. 641-644, 2015.*
- [8] Setiawardhana, Hakkun, R. Y., Baharuddin, A. Sign Language Learning based on Android for Hearing-impaired and Speech Impaired People. *International Electronics Symposium (IES). 114-117, 2015.*
- [9] Rastogi, R., Mittal, S., & Agarwal, S. A Novel Approach for Communication among Blind, Hearing-impaired and Dumb People. *2nd International Conference on Computing for Sustainable Global Development (INDIA Com). 605-610, 2015.*
- [10] Neha Baranwal ; Kumud Tripathi ; G.C. Nandi. Possibility Theory Based Continuous Indian Sign Language Gesture Recognition. *TENCON 2015 - 2015 IEEE Region 10 Conference, 2015.*
- [11] Ahmed, S., Islam, R., Zishan, Md. S. R., Hasan, M. R., & Islam, Md. N., Electronic Speaking System for Speech Impaired People: Speak Up. *2nd Int'l Conf. on Electrical Engineering and Information & Communication Technology (ICEEICT) Jahangirnagar University, Dhaka-1342, Bangladesh, 2015.*
- [12] Bueno, J., Requisitos para um ambiente de comunicação como ferramenta de apoio na alfabetização belíngue de crianças surdas. *Universidade Estadual do Paraná, Setor de Ciências Exatas, Curitiba, 2009*
- [13] Verma, V. K., Shrivastava, S. & Kumar, N. A. Comprehensive Review on Automation of Indian Sign Language” in *Computer Engineering and Applications (ICACEA), International Conference on Advances. 138-142, 2015.*
- [14] P. Parmar, A. P., Chitaliya, N. G., Gesture Recognition System for Indian Sign Language on Smart Phone. *International Journal of Advanced Research in Computer and Communication Engineering. 5(2), 376-379, 2016.*
- [15] Sulman, N., & Zuberi, S. ().Pakistan Sign Language – A Synopsis. 1-32.
- [16] Debevc, M., Kosec, P., Rotovnik, M., & Holzinger, M., Accessible Multimodal Web Pages with Sign Language Translations for Hearing-impaired and hard of hearing 72 users. *20th International Workshop on Database and Expert Systems Application. 279–283, 2009.*
- [17] Kitunen, S. Designing a Hearing-impaired Culture Specific Web Site. *Participatory Design Research for Knack. University of Art and Design Helsinki, Finlândia, 2009.*
- [18] de Queiroz, M. A., Acessibilidade web: Tudo tem sua Primeira Vez. <http://www.bengalalegal.com/capitulomaq.php>, 2015.
- [19] Kosec, P., Debevc, M., & Holzinger, A., Sign Language Interpreter Module: accessible video retrieval with subtitles. *Miesenberger, K., Klaus, J., Zagler, W., Karshmer, A. (eds.). ICCHP 2010, Part II. LNCS, Springer, Heidelberg , Vol (6180), 221–228, 2010.*
- [20] e-MAG., Modelo de acessibilidade em governo eletrônico. <http://www.governoeletronico.gov.br/acoes-e-projetos/e-MAG>, 2011.
- [21] WCAG, 2.0., Web Content Accessibility Guidelines. *Web Accessibility Initiative (WAI)*, <http://www.w3.org/TR/WCAG20/>, 2013.
- [22] Cynthia Says. <http://www.cynthiasays.com>, 2014.

TPACK Adaptation among Faculty Members of Education and ICT Departments in University of Sindh, Pakistan

Saira Soomro

Department of Distance Continuing and Computer Education, University of Sindh, Elsa Kazi Campus, Sindh, Pakistan

Tariq Bhatti

Faculty of Education
University of Sindh, Elsa Kazi Campus,
Hyderabad-76080, Sindh, Pakistan

Arjumand Bano Soomro

Institute of Information & Communication Technology
University of Sindh, Jamshoro, Pakistan

Nazish Basir

Institute of Information & Communication Technology
University of Sindh, Jamshoro, Pakistan

Najma Imtiaz Ali

Department of Information Systems,
International Islamic University, Kuala Lumpur, Malaysia/
Institute of Mathematics and Computer Science,
University of Sindh, Jamshoro, Pakistan

Nazia Parveen Gill

Department of Statistics
University of Sindh, Jamshoro,
Pakistan

Abstract—Technological Pedagogical Content Knowledge (TPACK) framework has been to investigate the technological and instructive knowledge of teachers. Many researchers have found this framework a useful tool to explore teachers' awareness regarding TPACK and how do they are relating it in learning and teaching process in different educational settings. During its first generation time period which was from year 2006 to year 2016, TPACK constructs took a decade to get explained and interpreted by researchers. Now, it has entered in its second generation but still contextual aspect yet not being explored in detail. This study addresses two areas; firstly, to measure the TPACK of faculty members of ICT and Education departments of University of Sindh; and secondly, to unfold the impact of four circumstantial/contextual factors (Technological, Culture of Institute, Interpersonal, and Intrapersonal) on the selected faculty members in using TPACK into their own subject domains. The results showed that both faculties are already taking in technology along with their teaching practices instead of limited technological resources. Besides this, they were found collaborative in teaching and open to the technology. This study reports the TPACK framework adaptation among higher education faculty members at University of Sindh. It also helped in understanding the intrapersonal beliefs of faculty members regarding technology integration with pedagogical and content knowledge.

Keywords—TPACK; teaching-learning; circumstantial and contextual factors

I. INTRODUCTION

Effective teaching-learning rely on the subject matter transformation and transfer to the learner in an understandable manner, this refer to the concept of pedagogical content knowledge (PCK). From the beginning of 21st century,

Information and Communication Technology (ICT) provides new ways to access and process knowledge in every field. In domain of education teachers also started using ICT for transferring their PCK to individual learners in their specific contexts. Higher Education institutions are one of those hubs where this transition is occurring very rapidly. For this the faculty members of higher education institutions should have to meet up with challenges caused by ICT integration into pedagogy and content in their specific subject domains. Technological Pedagogical Content Knowledge (TPACK) framework defines how ICT can be blended with pedagogical and content knowledge. The purpose of this research was to observe the TPACK awareness and adaptation among the faculty members at higher education institutions. The researchers have tried to explore and compare the TPACK knowledge of IT and Education Faculty of University of Sindh, Jamshoro Pakistan. The findings of this study will become a fact-finding analysis for teachers to improve their TPACK knowledge in respective subject domain.

II. LITERATURE REVIEW

An effective teaching encompasses the continuous improvement in teaching methods, in subject content and effective use of ICT (Information and Communication Technology) in teaching [1]. The TPACK framework gives a baseline to teachers about the integration of knowledge, content, pedagogy and ICT.

The first model of PCK (Pedagogical content knowledge) was suggested by Shulman [2] as shown in Fig. 1, it was further derived, rearranged and represented as TPACK by Mishra and Koehler [3]. TPACK is an extension of pedagogical content knowledge (PCK) concept. It has seven

complex constructs known as: 1) knowledge of technology (TK), 2) knowledge of content (CK), 3) knowledge of pedagogy (PK), 4) knowledge of pedagogy content (PCK), 5) knowledge of technology content (TCK), 6) knowledge of technology pedagogy (TPK), and 7) technological pedagogical content knowledge (TPACK) [4]–[8].

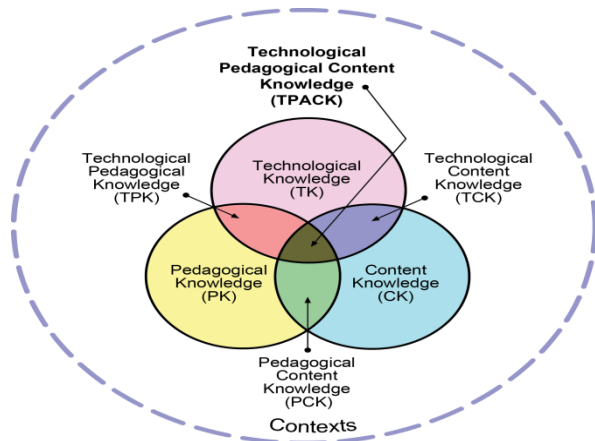


Fig. 1. Schema of Technological Pedagogical Content Knowledge (TPACK) source: [26].

The technological knowledge (TK), refers to the use of management and maintenance of ICT technologies such as wireless broadband, dial-up internet connection, digital photos and videos creation, software programs and hardware, and the use of interactive whiteboards, blackboards [3], [8]–[10].

Content knowledge (CK) is the knowledge of subject content which students expect to learn from teachers. The subject teacher is required to get equipped with sufficient content knowledge so he/she would be able to provide explanations to the students' queries [9], [11]. For the effective delivery of subject content, the pedagogy knowledge is required for teachers. It comprises of the knowledge of classroom instructions, classroom management skills, strategies for effective teaching, session planning and different assessment methods in alignment of set learning objectives of that particular session [9], [12]. The pedagogical content knowledge (PCK) refers to the particular knowledge of pedagogy that is suitable for specific subject content [6]. For better results of teaching instructions, the combination of pedagogy and content knowledge is essential [3], [9]. Technological content knowledge (TCK) helps in understanding that how different technologies can be used with subject contents to make teaching-learning process more effective [3]. Technological pedagogical knowledge (TPK) is an effective application of technology into teaching practices by the teachers. It can improve the instructions given to the students and help them in learning more effective way [4], [10], [13]. Finally the technological pedagogical content knowledge (TPACK) is fitting together the knowledge of content, technology, and pedagogy for delivery of subject content in efficient and effective manner [3].

Several studies have been conducted to authenticate the TPACK framework in different subject domains but most of them were conducted only on teacher educators or on pre-service teachers [5], [6], [14]–[19].

A 30-item questionnaire was developed by Lee and Tsai [16] to measure the World Wide Web (WWW) knowledge among teachers. The instrument used TPACK model as a framework. In total 558 teachers were selected from different Taiwan's schools were selected. The TPACK framework has been used to measure their Web knowledge at two levels, the Web knowledge in general and Web knowledge in communication.

Another study done by Koh, Chai and Tsai [20] to observe the TPACK knowledge of 1185 pre-service teachers who were enrolled in the Postgraduate Diploma/Diploma in Education programme at a higher education institute in Singapore. They performed exploratory factor analysis and found five out of seven TPACK constructs distinctive, namely technological knowledge, content knowledge, knowledge of pedagogy, knowledge of teaching with technology and knowledge from critical reflection. They reported that the participants of the study were unaware of differences between the TPACK constructs, particularly technological content knowledge and technological pedagogical knowledge.

Doukakis et al. [21] have adopted the TPACK framework and its instrument to measure the TPACK knowledge among upper secondary in-service teachers in Greece. The results showed that the sample of 1032 computer science teachers were high in content and technology knowledge. The teachers were found below average in their pedagogical content knowledge and technological content knowledge. This shows that they were unaware or unable to apply suitable teaching technique and technology both together in their subject area for teaching. In conclusion they still need to be guided to improve their PCK and TCK.

In another study conducted by Liang et al. [22], the in-service pre-school teachers were assessed by using an instrument based on TPACK in Taiwan. Their study explored 336 educational technology teachers by using the 42-items TPACK survey originally including seven scales (CK, PK, PCK, TK, TPK, TCK and TPCK) as discussed earlier in this paper. They performed exploratory factor analysis (EFA) to test reliability and validity of instrument which was adequate. The EFA also produced six scales, out of which five were same as per the original TPACK framework (i.e. CK, PK, PCK, TK and TPCK) and the sixth one was combination of technological pedagogical knowledge and technological content knowledge (TPTCK). Further results showed that pre-school teacher with more seniority possessed a certain level of resistance against technology-based teaching. It was also reported that teachers higher education qualification were more equipped with technological knowledge and comfortable with technology integrated teaching environment.

In a recent study of Mahmud [23] the TPACK has been used to assess English subject teachers in Indonesian context. The researcher through random sampling selected 74 in-service senior high school teachers. The 45-items based instrument validation and reliability reported as good. The results showed the English subject teachers at senior high school Pekanbaru, Indonesia were capable in integrating technology with content and pedagogy, as most of the teachers were experienced with good qualification. The factors with technology e.g. TCK, TK,

TPK and TPCK reported were reported low in their mean score. The researcher related this low score with the teachers' age and the English subject which they teach as it has no direct association with the technology. To bring the technological competency in teachers the author recommended the authorities and technical experts to facilitate teachers in acquire technological knowledge.

The TPACK model is a knowledge triad of content, pedagogy and technology and this intersected framework is compulsory for teachers to acquire for effective teaching-learning process. Before this TPACK framework was only used to measure the triad knowledge of teacher educators. In this paper this triad knowledge of other higher education teachers was measured and different insights will have explored due to different subject domains.

III. METHODOLOGY

A. Research Method

Within the framework of mixed method approach, the present study was based on survey procedures [24]. By using this procedure, both quantitative and qualitative data have been collected and analyzed. In beginning the quantitative study has been conducted by using a survey questionnaire. Then followed by a qualitative method in which researchers planned interviews for qualitative data collection. Mixed method allows the researcher to gain more insight from the combination of both qualitative and quantitative research [24]. This particular design is selected for two main reasons. Firstly, it offers an opportunity to counterbalance the weaknesses embedded within one method with the strengths of the other. Secondly, it helps researchers to perform exploration with a few cases or individuals (see Fig. 2).



Fig. 2. Conceptual framework of survey procedure.

B. Population

The target population for the study included all faculty members of departments of ICT and Education at University of Sindh, Pakistan. This population would have consisted of 38 faculty members of department of ICT and 26 faculty members of Department of Education.

C. Sample

Multistage sampling procedure was used because of study mixed method design. At the first stage, all faculty members were employed were selected as a sample of the study for the collection of quantitative data which include (26) Department of Education and (38) ICT faculty members. At the second stage for making sample highly representative two participants were selected from each department were selected for interviews.

D. Response Rate

The overall response rate was (67%), 43 questionnaires have been returned out of 64, which were distributed among faculty members of ICT (38) and Department of Education

(26). From Department of Education (21) were returned (80.7%) and from ICT (22) were returned (57.89%).

E. Research Instruments

In this study, the questionnaire, and interview were selected as instruments for collecting data from the participants. In this study, the instruments to collect data from the participants were questionnaire, in-depth interview and focus group discussion.

F. Questionnaire

For the first research question the questionnaire was used. The questionnaire consisted on seven constructs of TPACK and adopted from the study "DEVELOPMENT OF SURVEY OF TECHNOLOGICAL PEDAGOGICAL AND CONTENT KNOWLEDGE (TPACK)" already conducted by I. Sahin [15]. As the instrument's reliability and validity was already tested, so pilot testing was not directed.

Interviews: The interview was consisted of eight questions. From the all interviews four themes were emerged to address the second research question. All interviews were audio recorded. The qualitative data produced through interviews were analyzed for major themes to address the research questions that were posted at the outset of this study.

IV. RESULTS

We tried to answer our first research question in this part of the paper. The research question is given below.

Q#1. What difference can be found in measuring TPACK knowledge among university faculty of different subject domains (I.T and Education)?

A. Questionnaire Results

Subjects: In total the subject of this research were 43 University of Sindh teachers, out of which 21 affiliated with Department of Education and 22 with department of ICT. The respondents belong to different job levels and working experience (see Tables I and II). Most of the participants (32) were female (74.4%) out of total 43 and (11) were male (25.5%). In faculty, wise gender distribution there were (14) females (66.6%) and (7) male (33.3%) in Department of Education. In ICT, the female participants were (12) that was (55%) and (10) were male (45.45%). Define abbreviations and acronyms the first time they are used in the text, even after they have been defined in the abstract. Abbreviations such as IEEE, SI, MKS, CGS, sc, dc, and rms do not have to be defined. Do not use abbreviations in the title or heads unless they are unavoidable.

TABLE I. WORK EXPERIENCE VICE DISTRIBUTION

Work Experience (in years)	No. of Participants
1-5	17
6-10	4
11-15	10
16-20	7
21 & above	5

TABLE II. JOB VICE DISTRIBUTION

Job Title	No. of Participants
Research Associate	8
Instructor/Lecturer	18
Assistant Professor	16
Associate Professor	1

B. TPACK score distribution

In Fig. 3 the distribution shows the Technological Knowledge (TK) has considerably has higher score among all. It also shows that department of education is lower in Technological Pedagogical Knowledge (TPK).

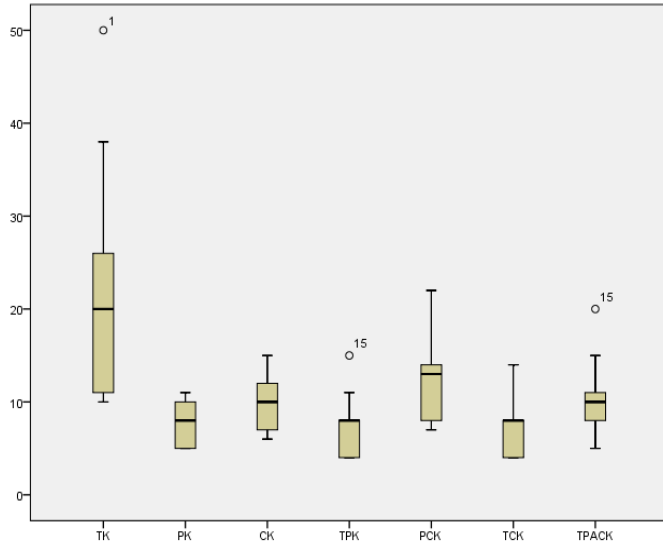


Fig. 3. TPACK adaptation in Department of Education.

In Fig. 4, the Pedagogical Content Knowledge (PCK) scores distribution much higher than other TPACK factors. And Technological Content Knowledge (TCK) has lowest score distribution among all within faculty members of ICT.

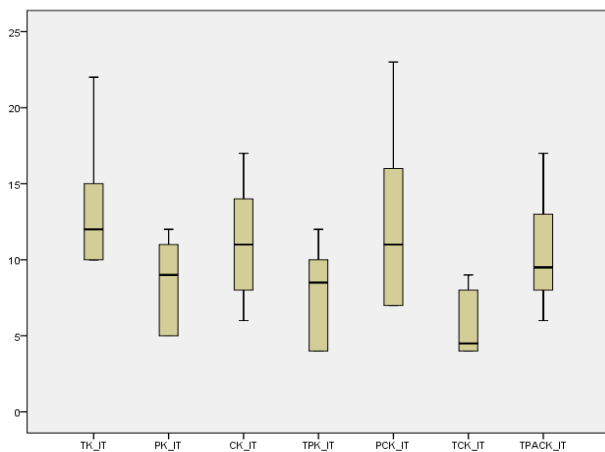


Fig. 4. TPACK adaptation in ICT.

In Fig. 5, the overall TPACK score distribution for ICT and Education department can be seen together.

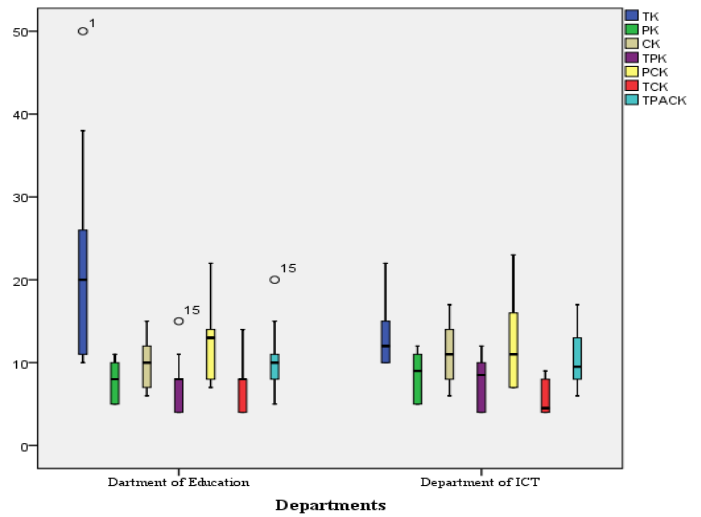


Fig. 5. Overall factor vice score distribution.

Fig. 6 shows the overall score distribution for TPACK adapted by both faculty members of ICT and Education. The figures show there is overlapping in boxes, therefore there is likely to be a difference between both groups. Although the score distribution for ICT faculty members seem higher but the median values for both (10) department of education and (9) ICT are nearly same.

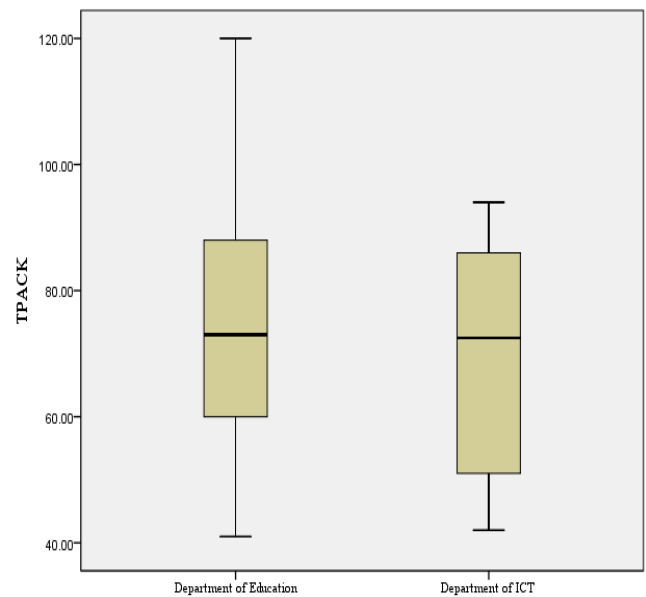


Fig. 6. TPACK adaptation in departments of Education and ICT.

In Fig. 7, the histogram for department of Education can be seen, it contains outlier data. In the graphical check for outliers in data we have already generated the box plots (see Fig. 2 and 4). The reason for these outliers could be the unusual or unrealistic response towards the TPACK adaptation, particularly factors TK, TPK and TPACK.

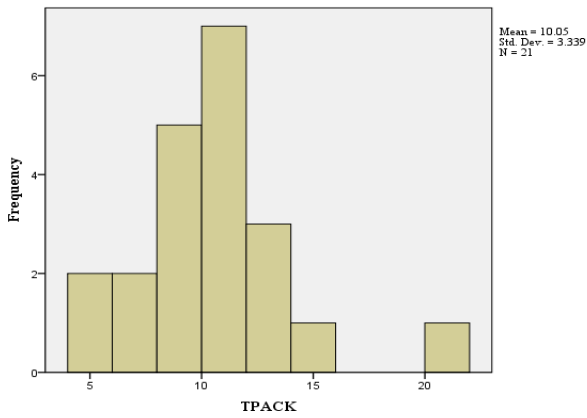


Fig. 7. Histogram for Department of Education.

In the similar way the Fig. 8 also showing outliers in the histogram diagram for department of ICT. But for this the box plot (see Fig. 2 & 4) there are no outliers shown.

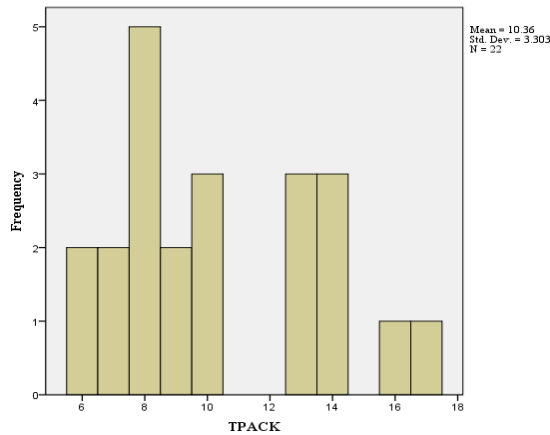


Fig. 8. Histogram for department of ICT.

For further check the Q-Q plots also generated (see Fig. 9 and 10). In Fig. 8, an outlier can be noticed prominently on the other hand Fig. 8 has no outliers in it. Beside this both Q-Q plots show the normality of data upto some level. Although in case of ICT (see Fig. 10) the points are bit far from normal line and it also show there are two groups within data among which, one is outlying from the other part.

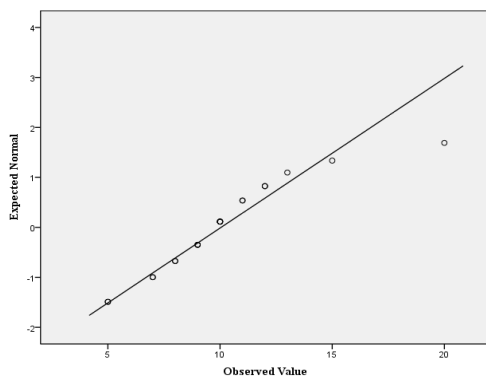


Fig. 9. Q-Q plot for Department of Education.

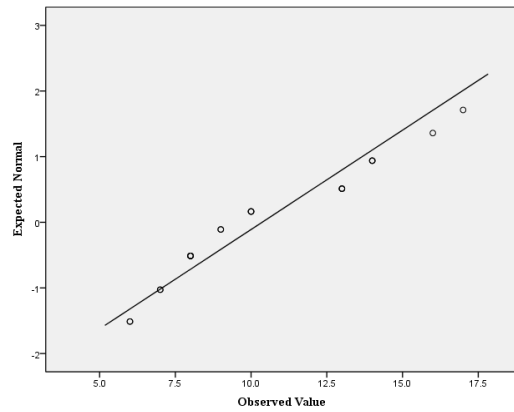


Fig. 10. Q-Q plot for Department of ICT.

C. Interview Results

Q#2. How the circumstantial/contextual factors affect the TPACK of individual faculty?

This question brought four themes to the forefront of the study, (a) Technological (b) Culture of institute (c) Intrapersonal and (d) Interpersonal. Under each theme different questions were from participants of each research site.

1) THEME 01: TECHNOLOGICAL

Q1 Is computer available at your office for session planning?

Participant#1 “Yes, computers are available in every office.”

Participant#2 “Each faculty member has his/her personal laptops.”

Participant#3 “Yes, computers are available for use.”

Participant#4 “Almost every faculty member has their own laptop bought personally by them”.

Q2 Does your department provide you fully furnished computer lab with updated hardware and software in executing your sessions?

Participant#1 “No, computer lab is not updated for classes because funds are properly given.”

Participant#2 “Yes, as I belong to ICT department so there are 5 to 6 computer labs.”

Participant#3 “Yes, fully furnished computer labs are available for use.”

Participant#4 “We do not take classes in lab, each faculty bring their laptop to conduct practical if they wish to use computer in class.”

2) THEME 02: CULTURE OF INSTITUTE

Q3 Is your nationwide higher education policies promote technology use with pedagogy and content?

Participant#1 “Yes, there are various but technology with pedagogy is a bit difficult.”

Participant#2 “Yes, we use technology with content but as we are from ICT department so most of the faculty members do not have B.Ed. / M.Ed degrees.”

Participant#3 “Yes, few of the teachers of education department use technology with content and pedagogy.”

Participant#4 “Using technology with pedagogy is our strength, but computer lab setup restricts us for extensive use.”

Q4 Do you integrate technology into course outlines of various subjects of your domain?

Participant#1 “Yes, as technology is our subject so we use it frequently.”

Participant#2 “Yes, we technology with content subjects.”

Participant#3 “Yes, few of the teachers to integrate technology with their subjects.”

Participant#4 “We have devices like multimedia projector and other technology tools.”

3) THEME 03: INTRAPERSONAL

Q5 What are your beliefs about using technology into teaching-learning process?

Participant#1 “Technology boosts teaching-learning, and can help teachers to achieve their session objectives.”

Participant#2 “Technology helps us as a teacher to make our students clear about the topics of the course.”

Participant#3 “Yes, I believe that technology helps teachers to integrate technology with their subjects.”

Participant#4 “We use different devices like multimedia projector and other technology tools.”

4) THEME 04: INTERPERSONAL

Q6 Do you think collaboration with colleagues increases your motivation to use technology into teaching-learning process?

Participant#1 “Yes, with help and collaboration of colleagues increases our motivation”

Participant#2 “It depends on department environment to use technology with collaboration into teaching-learning process.”

Participant#3 “Yes, I believe that technology helps teachers to do team teaching.”

Participant#4 “We use different collaboration tools in collaboration with other teachers.”

The second research question “How the circumstantial/contextual factors affect the TPACK of individual faculty?” uncovered the effect of circumstantial/ contextual factors on teachers’ TPACK. The factor that influenced teachers’ TPACK was the Technological devices availability. The result of interview showed positive trend towards technology availability. Most of the faculty members have personal laptops. On the fully furnished computer labs the researcher got mixed responses. As in ICT department there were 5 to 6 labs but in comparison to education department they possess only one lab. The second factor was culture of institute in

which researcher got responses on nationwide policies for technology. One participant answered that “Yes, there are various policies but using technology with pedagogy is bit difficult”. Another said that “We believe in technology with content but as we are from ICT department so most of the faculty does not have B.Ed./M.Ed”. On enquiry of one question about technology integration researcher received the answers that “We use technology is our subject frequently”. Another participant said that, “Yes few of the teachers integrate technology in their subjects”. From the result, it is obvious teachers need more time and effort to use technology into subjects as a pedagogy but still students’ skills regarding ICT skills need to be studied. The interview further explored the teacher’s intrapersonal beliefs in which they said that technology helps us as a teacher to make our students clear about the topics of the course. Another said they believe that technology helps teachers to teach subjects with more clarity. This same belief was also found in a study beliefs influence their practice of ICT integration [25].

V. DISCUSSION

Pakistan is a developing country although the technology trends are not unfamiliar in this region but still the adaptation is slow. Especially in the field of education there is still no proper infrastructure which can pursue teacher and learner to adapt technology in their teaching-learning process.

University of Sindh situated in the Province of Sindh at the distance of around 155km from the capital city Karachi. Although this University is not rich in ICT resources but faculty members are well equipped with the knowledge of using the technological instruments and try to implement in their teaching-learning process.

Being the part of this University faculty we conducted this study to get the clear picture about other faculty members’ perceptions and approach to adapt TPACK in their particular domain of teaching. For this we chose two departments ICT and education, one of the reasons to choose these was the easy access to the faculty members and to get their response on time as two of the authors belong to these departments, this helped us in data collection through questionnaires and interviews.

The results produced from the data collected through questionnaires were unable to provide us the clear picture therefore we have conducted interviews. While comparing the two groups of data or in other words to do analysis of variance we needed to do some checks. The main requirement for this parametric technique to compare the groups was the continuous scale instead of discrete, the TPACK instrument contained 5-point Likert scale, the other check was the random sampling, the number of faculty members of our selected departments was small therefore we employed all the faculty members which made us to violate this second check. The third check was to keep the observations independent, during questionnaire distribution we visited the individual faculty members in their offices instead of a particular venue where they gather and do discussions. For the normality check we ran explore in SPSS 22 the results are given in the data analysis part. The outliers in results and the violation of second that was random sampling made us to conduct the interviews. Besides all these issues the questionnaire part was not totally worthless

it helped us to answer our first research question and up to some level we figured out about the TPACK factors which were highly adapted in both faculty members of ICT and Education department.

From interviews, it revealed that both faculties are already persuaded to incorporate technology in their teaching and learning process, although the technological resources are not sufficient but they are compensating this by having their personal gadgets. They have no issue in collaborative teaching and they are also open to the technology. Only if there would be a technological infrastructure, then they can employ TPACK in more effective manner.

ACKNOWLEDGMENT

We acknowledge faculty members of I.T. and Education Faculty University of Sindh, Jamshoro, who voluntarily participated in this study and gave their precious time to this insight sharing activity, which became a research study ultimately.

REFERENCES

- [1] G. Alayyar, P. Fisser, and J. Voogt, "Developing technological pedagogical content knowledge in pre-service science teachers: Support from blended learning," *Australas. J. Educ.*, 2012.
- [2] L. Shulman, "Those who understand: Knowledge growth in teaching," *Educ. Res.*, 1986.
- [3] P. Mishra and M. Koehler, "Technological pedagogical content knowledge: A framework for teacher knowledge," *Teach. Coll. Rec.*, 2006.
- [4] E. Baran, H. H.-H. Chuang, and A. Thompson, "tpack: an emerging research and development tool for teacher educators," *Turkish online J. Educ. Technol.*, vol. 10, no. 4, pp. 370–377, 2011.
- [5] D. Schmidt, E. Baran, and A. Thompson, "Technological pedagogical content knowledge (TPACK) the development and validation of an assessment instrument for preservice teachers," *J. Res.*, 2009.
- [6] M. Koehler and P. Mishra, "What happens when teachers design educational technology? The development of technological pedagogical content knowledge," *J. Educ. Comput. Res.*, 2005.
- [7] M. Koehler, P. Mishra, and K. Yahya, "Tracing the development of teacher knowledge in a design seminar: Integrating content, pedagogy and technology," *Comput. Educ.*, 2007.
- [8] M. Koehler and P. Mishra, "What is technological pedagogical content knowledge," *issues Technol. Teach. Educ.*, 2009.
- [9] "TPACK: An emerging research and development tool for teacher educators," *TOJET: The Turkish*, 2011.
- [10] M. Koehler and P. Mishra, "Introducing TPCK. AACTE Committee on Innovation and Technology (Ed.), The handbook of technological pedagogical content knowledge (TPCK) for," 2008.
- [11] P. Nilsson, "Teaching for understanding: The complex nature of pedagogical content knowledge in pre - service education," *Int. J. Sci. Educ.*, 2008.
- [12] J. Hinojosa, C. Labbé, L. López, and H. Iost, "Traditional and emerging IT applications for learning," *Int. Handb.*, 2008.
- [13] J. Harris, P. Mishra, and M. Koehler, "Teachers' technological pedagogical content knowledge and learning activity types: Curriculum-based technology integration reframed," *J. Res. Technol.*, 2009.
- [14] S. Jang and M. Tsai, "Exploring the TPACK of Taiwanese elementary mathematics and science teachers with respect to use of interactive whiteboards," *Comput. Educ.*, 2012.
- [15] I. Sahin, "Development of survey of technological pedagogical and content knowledge (TPACK)," *TOJET Turkish Online J. Educ.*, 2011.
- [16] [M. Lee and C. Tsai, "Exploring teachers' perceived self efficacy and technological pedagogical content knowledge with respect to educational use of the World Wide Web," *Instr. Sci.*, 2010.
- [17] C. Angeli and N. Valanides, "Epistemological and methodological issues for the conceptualization, development, and assessment of ICT-TPCK: Advances in technological pedagogical content," *Comput. Educ.*, 2009.
- [18] [18] L. Archambault and K. Crippen, "Examining TPACK among K-12 online distance educators in the United States," *Contemp. issues Technol.*, 2009.
- [19] R. Graham, N. Burgoyne, P. Cantrell, L. Smith, and L. S. Clair, "Measuring the TPACK confidence of inservice science teachers," *TechTrends*, 2009.
- [20] J. Koh, C. Chai, and C. Tsai, "Examining the technological pedagogical content knowledge of Singapore pre - service teachers with a large - scale survey," *J. Comput. Assist.*, 2010.
- [21] S. Doukakis, A. Psaltidou, and A. Stavragi, "Measuring the technological pedagogical content knowledge (TPACK) of in-service teachers of computer science who teach algorithms and programming in," *Technol.*, 2010.
- [22] J. Liang, C. Chai, J. Koh, C. Yang, and C. Tsai, "Surveying in-service preschool teachers' technological pedagogical content knowledge," 2013.
- [23] M. Mahdum, "Technological Pedagogical and Content Knowledge (TPACK) of English Teachers in Pekanbaru, Riau, Indonesia," *Mediterr. J. Soc. Sci.*, 2015.
- [24] J. W. Creswell, *Research Design: Qualitative, Quantitative and Mixed Method Approaches*. Sage Publications Inc., 2013.
- [25] P. Ertmer, A. Ottenbreit-Leftwich, and O. Sadik, "Teacher beliefs and technology integration practices: A critical relationship," *Comput.*, 2012.
- [26] http://www.tpack.org/tpck/index.php?title=TPCK_ _Technological_Pedagogical_Content_Knowledge

Quality Aspects of Continuous Delivery in Practice

Maryam Shahzeydi

Computer Engineering Department
Dolatabad Branch, Islamic Azad University,
Isfahan, Iran

Taghi Javdani Gandomani*

Computer Engineering Department
Boroujen Branch, Islamic Azad University, Boroujen, Iran
Computer Engineering Department
Dolatabad Branch, Islamic Azad University, Isfahan, Iran

Rasool Sadeghi

Computer Engineering Department
Dolatabad Branch, Islamic Azad University, Isfahan, Iran

Abstract—Continuous Delivery is recently used in software projects to facilitate the process of product delivery in Agile software development. As an Agile practice, this practice is mainly used to achieve better quality of software development process and higher customer satisfaction. However, less attention has been paid on exploring the quality factors related to Continuous Delivery as well as quality model. The main aim of this paper is to figure out the quality aspects and factors of Continuous Delivery. Initial data analysis showed that this practice is impressed by people related factors, organizational issues, tools and process related factors as well.

Keywords—Continuous delivery; quality model; agile software development; agile methods; agile practice

I. INTRODUCTION

Agile methods are widely using in software development projects since the last decade. These methods promote a different style of software development which distinguishes the development from traditional or disciplined methods in software engineering. Focusing on Agile values and principles, defined in Agile manifesto [1], these methods promote early and frequent delivery, higher quality level, better customer collaboration, embracing required changes in customer's requirements and so on [2]. It's why many software companies are looking for the best way to adopt these methods in their software product lines [3]. However, they are faced with various challenges [4].

Agile software development includes various methods such as Scrum, Extreme Programming (XP), Crystal family, Test Driven Development (TDD), Feature Driven Development (FDD), etc. [2], each defines its own particular practices, roles and artifacts. However, usually, Agile software teams use various practices that can be commonly used in all Agile methods [4], [5]. Continuous Delivery (CD) is one of the popular practices which recently has gained special importance for Agile projects.

CD focuses on releasing reliable software product through software development, test and deployment [6]. This practice was introduced in 2010, as the ability to release every time [7]. However, the core concept of CD is not really continuous code development; it is the ability of release at any time [6], [8]. Indeed, the recently developed code should have the ability to

be added by new features and functionalities as easy as possible.

Since the ultimate goal of software development is achieving customer satisfaction through increasing quality of both development process and product, quality of all the development practices is important. Quality of CD also play a great role in customer satisfaction. Better conduction of this practice may lead to higher customer trust directly and satisfaction indirectly. However, the literature review shows less effort on exploring the quality related aspects of CD, proposing a clear quality model, or even providing guidelines to increase quality of this practice in real environments.

This article tries to explain the concept of CD from the lens of quality and address the most related previous studies and finally describe the outline of a required quality model dedicated for this practice. So, the rest of this article is organized as follows: Section 2 describe the underlying concept of CD briefly. Section 3 addresses the most related works followed by Section 4 which outlines a quality model associated for CD. Section 5, finally, concludes the paper.

II. CONTINUOUS DELIVERY

Agile software development defines an underpinning concept, short cycles, which its focus is on early and frequent delivery. To establish such concept, Agile approach defines proper practices, among them CD plays a critical role. As mentioned before, CD focuses on the ability of software release whenever customer needs [6], [8]. CD is a really a practice to help the software stakeholders (i.e. business and technical parties) to collaborate in development and deployment of a software product in short cycles while focusing on the quality factors.

Technically speaking, CD is considered as an Agile practice which facilitates the process of delivery of product increments upon the customer request. However, CD focuses more on commitment to ensuring the recently developed code is able to be released at any time rather than the delivery process [9]. The promised advantages of CD temp both business and technical practices of software development to adopt it in their product line [8]. Accelerating time-to-value, quick user feedback, achieving clear and visible believable

progress, reducing the risks of delivery, providing innovations in the release process, better quality and data-driven decision making are the most addressed advantages and benefits of CD in practice [10], [11].

The above advantages have root in the concepts and goals of CD. For instance, frequent delivery and release provides the ability to get customer feedback timely and faster. Also, short cycles and frequent delivery increase the chance of risk discovery and avoid them in the next delivery. So, better quality will be expected indirectly. Recently and in the

competitive software industry, lots of the reputed companies such as Facebook, Google, IBM and Microsoft are trying to use CD as a compulsory development practice in their project [6], [11].

CD process includes a series of activities all together are known as “Continuous Delivery Pipeline”. As shown in Fig. 1, this pipeline involves some automatic and manual tasks. Although, literature review shows different steps for this pipeline, all are almost the same in tasks and activities in which Build, Staging and Production are constant [7], [11], [12].

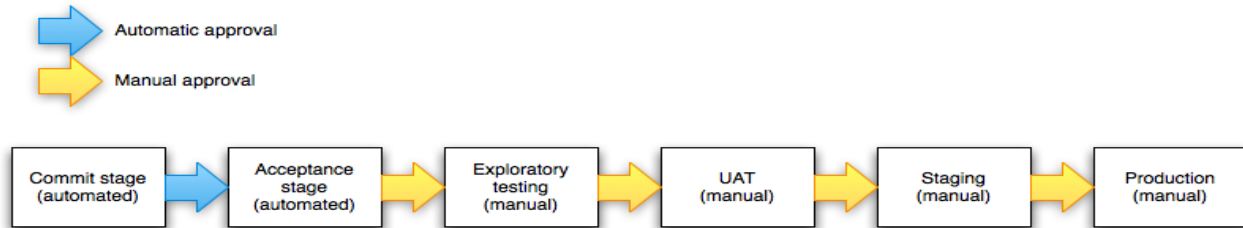


Fig. 1. Typical continuous delivery pipeline.

In the Build stage software teams use source repository as input and store an artifact in the artifact repository. The main goal of this stage is software development, software test, packaging and archiving. Unit tests are mostly used in this stage. The second stage, Staging, software teams install and deploy the recently build artifact in a staging environment and simultaneously perform regression, performance, integration and functional tests. Production stage, finally, focuses on deployment of the recently testes software into the production environment [7].

Despite of its simple concept, employing of CD in practice needs proper conditions. For instance, software development process needs to support iterative development in advance [13]. Indeed, without defining several iterations, CD cannot be considered. This would be a serious limitation for small projects where number of iterations are limited. Furthermore, extensive of positive team climate and also positive atmosphere between customer and development team is necessary [13].

III. RESEARCH BACKGROUND

Most of the previous studies paid attention to introduction and employment of CD only. Indeed, less attempts have been made to determine and highlight the quality factors and aspects of CD in practice. However, a few studies have referred to this issue.

Some studies focused on the barriers and challenges of employing and quality of CD in real environments. “organizational challenges” was reported as a serious challenge in the CD pipeline [5]. Another study [13] technical, procedural and customer-related challenges of CD have been addressed and the details of each were explained. For instance, issues with CD downtime, problems and limitations of automatic test process and configuration related problems are listed as technical challenges.

Another study tried to create a trade-off between risk of lower release quality and time-to-market while adoption of CD

[7]. Agile practices and their impacts on employing of CD were investigated in another study. This study showed that while some Agile practices like TDD, Pair testing, and customer involvement and collaboration have a positive and significant impact on the CD, some others like Pair programming have not such impact.

In another study, a new eco-system, Rugby, was proposed to support the CD life cycle and facilitate its pipeline [12]. The main focus of this study was to indicate the impact of Agile approach on the CD pipeline. The proposed eco-system defined some new roles such as team leader, project leader, customer, and developer to support and facilitate the CD adoption in real environments. The results of this showed the increase of frequency and quality of interactions between development team and customer party.

In another study, some of the adaptable quality metrics of CD were addressed. These metrics have been categorized as project level, product level, and pipeline level. The addressed metrics are suitable to be used in evaluating quality of CD.

In sum up, literature review shows that only a few studies focused on quality aspects of CD. This indicates a research gap that can be fulfilled by conducting proper research studies in practice. Focusing on this gap, the next section provides some quality factors that may affect the process of CD in practice.

IV. OUTLINE OF CD QUALITY

Conducting a qualitative research study led to collection of proper data related the topic under study, CD quality aspect. Data collection and analysis are ongoing at the time of this writing. However, some aspects of the results can be showed in this article. This section provides the main findings of this study. However, the details of each aspect and evaluation of the findings will be provided in another article in future.

Data analysis showed that quality of CD is impressed by four different aspects including People, Process, Organization, and Tools, as shown in Fig. 2. These aspects are the high level abstract for various quality factors. Indeed, each of them

consists of several quality factors which together impress the quality of CD in practice.

'People' category mainly indicates that people related issues are important factors that impress quality of CD. People relationship is so critical in performing CD since this practices connect both technical and business parties. Furthermore, the relationship between development team members also is important, because it seems that collaborative teams conduct CD in a better quality.

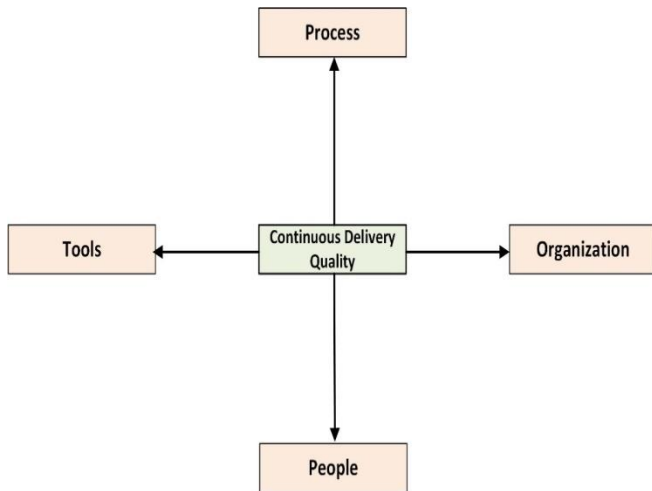


Fig. 2. The outline of quality aspects of CD.

Usually the people involved the CD process having different level of experience and so this aspect can impress the quality of CD too.

'Process' of CD has a great impact on quality of this practice. Various activities included in this process such as TDD, frequent testing, mechanisms used for requirement prioritization, and daily continuous integration seriously needed to be perform in the professional manner. Therefore, any weakness in doing such activities results in low quality of CD directly.

In 'organization' category the main focus in on the organizational culture and its related issues. Existence of culture of CD in organizational processes is compulsory to achieve the desired quality of CD. Also, providing mechanisms to manage the potential technical and human related risks greatly can lead to better quality of CD. Moreover, quality control and assurance and its process positively impress the quality of all the involved practices generally and CD particularly.

'Tools' category deals with tools related issues. For instance, automatic facilitates directly accelerate the process of CD and avoid the human related errors in this practice. Also, existence of mechanisms for version controlling leads to reduce configuration related defects.

In sum up, it seems that quality of CD depends of various technical and human related activities. However, more data analysis is necessary to explore the details of the above mentioned aspects, as noted earlier.

V. CONCLUSION AND FUTURE WORK

CD is one the most important practices which recently is widely used in software projects. This practice focuses on the ability of software release at any time. CD defines a sequential set of activities to facilitate the release process. Quality of CD directly impresses the quality of development process. To explore the quality factors and aspects of CD, a qualitative study has been conducted. Initial data analysis showed that quality of CD is impressed by four aspects including People, Organization, Process, and Tools related factors. Each of these aspects by involving some quality factors may lead to better quality of CD in practice and real environments.

For the future work, the authors intend to employ the proposed model in two case studies to evaluate its usefulness and applicability in an empirical study.

ACKNOWLEDGMENT

The authors would like to express their gratitude to all the participants of the research.

REFERENCES

- [1] K. Beck et al. (2001, May 2014). Agile Manifesto. Available: www.agilemanifesto.org
- [2] D. Cohen, M. Lindvall, and P. Costa, "An introduction to Agile methods," *Advances in computers*, vol. 62, pp. 1-66, 2004.
- [3] T. J. Gandomani and M. Z. Nafchi, "An empirically-developed framework for Agile transition and adoption: A Grounded Theory approach," *Journal of Systems and Software*, vol. 107, pp. 204-219, 2015.
- [4] T. J. Gandomani and M. Z. Nafchi, "Agile transition and adoption human-related challenges and issues: A Grounded Theory approach," *Computers in Human Behavior*, vol. 62, pp. 257-266, 2016.
- [5] H. C. Esfahani, "Transitioning to Agile: A Framework for Pre-Adoption Analysis using Empirical Knowledge and Strategic Modeling," PhD, Graduate Department of Computer Science, University of Toronto, Canada, 2012.
- [6] L. Chen, "Continuous delivery: Huge benefits, but challenges too," *IEEE Software*, vol. 32, no. 2, pp. 50-54, 2015.
- [7] J. Humble and D. Farley, *Continuous Delivery: Reliable Software Releases through Build, Test, and Deployment Automation* (Adobe Reader). Pearson Education, 2010.
- [8] G. Schermann, J. Cito, P. Leitner, and H. C. Gall, "Towards quality gates in continuous delivery and deployment," in *Program Comprehension (ICPC), 2016 IEEE 24th International Conference on*, 2016, pp. 1-4: IEEE.
- [9] Atlassian. (2016, Dec. 2017). Continuous Delivery. Available: <https://www.atlassian.com/continuous-delivery>
- [10] T. Dingsøyr and C. Lassenius, "Emerging themes in agile software development: Introduction to the special section on continuous value delivery," *Information and Software Technology*, vol. 77, pp. 56-60, 2016.
- [11] G. G. Claps, R. B. Svensson, and A. Aurum, "On the journey to continuous deployment: Technical and social challenges along the way," *Information and Software technology*, vol. 57, pp. 21-31, 2015.
- [12] S. Neely and S. Stolt, "Continuous delivery? easy! just change everything (well, maybe it is not that easy)," in *Agile Conference (AGILE), 2013*, 2013, pp. 121-128: IEEE.
- [13] S. Krusche, L. Alperowitz, B. Bruegge, and M. O. Wagner, "Rugby: an agile process model based on continuous delivery," in *Proceedings of the 1st International Workshop on Rapid Continuous Software Engineering*, 2014, pp. 42-50: ACM.

Integration of Heterogeneous Requirements using Ontologies

Ahmad Mustafa, Wan M.N. Wan-Kadir,
Noraini Ibrahim
Software Engineering Department
Universiti Teknologi Malaysia
Johor Bahru, Malaysia

Muhammad Arif Shah
Software Engineering Department
Universiti Teknologi Malaysia, Johor Bahru, Malaysia
City University of Sciences and Information Technology,
Peshawar Pakistan

Muhammad Younas
Software Engineering Department
Universiti Teknologi Malaysia, Johor Bahru, Malaysia
Department of Computer Science
Government College University, Faisalabad, Pakistan

Abstract—Ontology-driven approaches are used to sustain the requirement engineering process. Ontologies can be used to define information and knowledge semantics during the requirements engineering phases, such as analysis, specification, validation and management of requirements. However, requirement analysts face difficulties in using ontologies for requirement engineering. In this study, a framework has been proposed to integrate heterogeneous requirements by using local and global ontologies.

Keywords—Heterogeneous requirements; requirement engineering; local ontologies; global ontologies

I. INTRODUCTION

The success of software system can be quantified by the degree to which it meets the proposed envision. Software Requirements Engineering (RE) is a well-defined process to identify stakeholders and their needs. It is a way to document the requirements for agreeable analysis, communication, and further implementation [1]. Moreover, all the efforts and resources applied during RE process must be reflected in the developed system in term of quality product and delivery of product time to market. However, Standish group report illustrates that 31.1% of projects are canceled before completion, and 52.7% of projects cost increased up to 189% to its original estimated cost. Moreover, many projects are failed due to lack of user input, incomplete requirement specifications and change in requirement specifications [2].

In the early phases of software development process, RE emphasizes on elicitation analysis, specification, validation, and management of requirements [3]. It is recognized that RE highlights on the enhancement of the quality of system under development. Moreover, RE focuses on reasons that may

propagate the risks such as budget overrun, time delay and project failures [4], [5].

In software development industry and academia, ontologies are being used in requirements engineering phase. Ontology is defined as “explicit specification of a conceptualization” [6], [7]. Ontology has explicit classes and properties and used as a standard form for knowledge representation of concepts inside a domain. Furthermore, it establishes an association in such a way that is allowable for automated reasoning [6], [7].

Furthermore, the ontological concepts can be used to address or resolve various kinds of issues in RE. It is used to write complete, unambiguous and consistent requirements statements. Furthermore, ontologies can be used to manage heterogeneous requirements, accomplish consistency analysis, represents domain knowledge model and requirements changes [8]-[10].

Rest of the paper is organized as follows. Section II describes the overview of heterogeneous requirements. Section III describes the related studies. The proposed approach is mentioned in Section IV which have following subsections, ontology as shared vocabulary using Multilanguage WordNet, generic requirements formalism using Pivot Model and transformation of requirements into Pivot Model. Section V describes the limitation of the study, and the last section the conclusion and future work of the study.

II. OVERVIEW OF HETEROGENEOUS REQUIREMENTS

A Multinational Corporation (MNC)¹, has several departments situated in different countries. For instance, as shown in Fig. 1 requirements are described in Malay, Arabic and German languages.

¹ https://en.wikipedia.org/wiki/Multinational_corporation

III. RELATED MATERIAL

A. Requirements Engineering

Requirements engineering activities aim to manage requirements-related knowledge such as natural language documents, storyboards, use cases, and business process specifications. These artifacts are called Requirements Document. A requirement document is an initial process in software development process. The development of these documents is taken as one of the challenging task [11].

Requirements engineering is concerned with (a) elicitation: actors and requirements identification, (b) modeling of identified requirements, (c) analysis of requirements to detect inconsistency and ambiguity and (d) validation of requirements [12]. According to Sommerville [3], the requirement statements are the descriptions of what services the system should provide and the restrictions on their operations. These requirements imitate the needs of customers for a system that serves a specific purpose such as placing an order, finding information, and controlling a device. The process of discovering, analyzing, recording requirements documents and verifying these services and constraints is known as requirements engineering (Fig. 2).

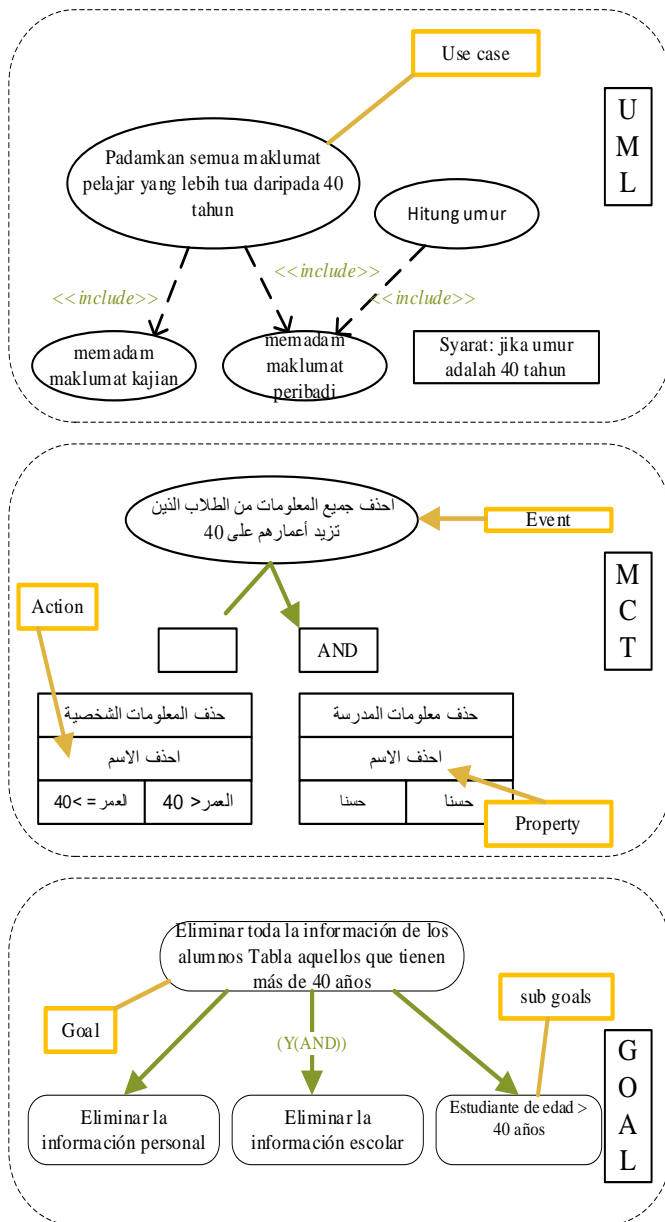


Fig. 1. Example heterogeneous requirements.

Moreover, each department is not only using different languages but also using different formalism such as department from Malaysia is using Unified Modeling Language (UML), Gulf Countries are using MERISE Conceptual Model of Treatment (MCT), and German version of requirements formalism is a goal-oriented approach. The example showed in the Fig. 1 mainly focuses on the following:

- Integration of heterogeneous requirements from a different partner into one language which is understood by all partners.
- Reasoning on requirements for identification of relationship. i.e., in Fig. 1 all partners have the same need.

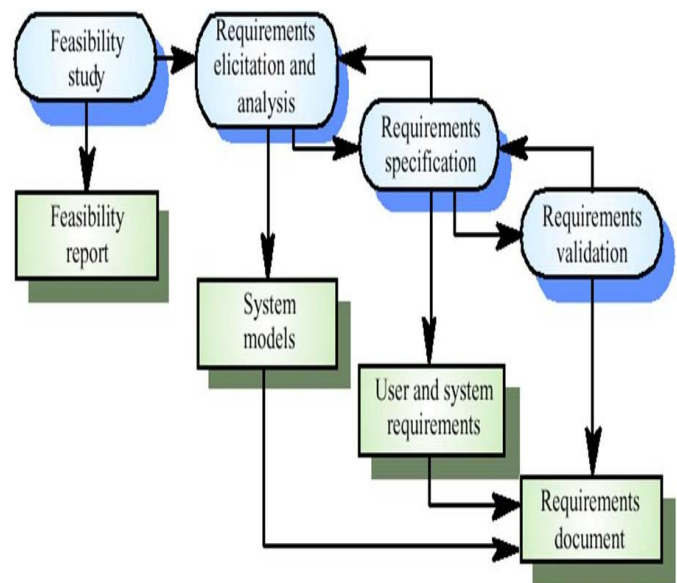


Fig. 2. Process of requirement engineering [3].

B. Ontology Engineering

This section presents ontology definition and how ontology can be used to store information captured from functional requirements.

A standard definition of ontology is coined by Gruber [13], "An ontology is an explicit, formal specification of a shared conceptualization." This definition emphasizes on two main points: formal conceptualization and automated reasoning, and designing the domain-specific ontology.

The word ontology has borrowed initially from philosophy within less than twenty years, and it means the philosophical study of nature of existence. In other words, ontology comprises recognizing the essential categories of things such as

ontology may be used to group objects as abstract or concrete, essential or existential. Although, term ontology is borrowed from philosophy though it gains substantial popularity in computer science and information science [14], [15].

Likewise, the Web Ontology Language (OWL)² is a language for modeling ontologies for the semantic web, and it is recommended standard proposed by the World Wide Web Consortium (W3C). OWL design is a more expressive language which defines the classes axiomatically and with the support consistency reasoning [16]. OWL is evolved from languages those were developed by joining two standards: European standard Ontology Interface Layer (OIL) and American standard DARPA³ Agent Markup Language (DAML). OWL allows specifying various ontology modeling paradigms such as hierarchical relationships, restrictions, modeling attributes, and associations under a well-defined semantic [17].

Numerous definitions have been proposed for ontology [13], [18], [19], domain ontology is considered as domain conceptualization of formal, consensual and referenceable regarding classes and properties. This definition focuses on three criteria that differentiate ontologies from other models used in computer science:

- **Formal:** In the context of requirements engineering, ontologies offer reasoning capabilities to diagnose inconsistency and incompleteness from requirements specification. The formal ontology conceptualization is grounded in a formal theory which checks the level of consistency and automated reasoning over the ontological concepts and individuals.
- **Consensual:** The consensual aspect allows the designer to share their models using global ontologies.
- **Referenceable:** It is the capability of ontology that each concept either class or property can be referred through a unique identifier. It helps in defining requirements semantics.

Another definition of ontology in computer science is coined by Gruber Gruber [20], "An ontology is a representation of specification of a conceptualization." Consequently, an ontology represents the conceptual model of the specific domain of interest, describing it in a declarative fashion [21].

C. *Ontology-driven Approaches Contribution in RE*

This section briefly presents state of the art studies of ontology-driven in RE. With the globalization of world in the case of the complex system, many technical experts from different fields, department, research lab and from a different part of the world participate in the project. Most of the time, they use their favorite languages to formalize the requirements of their assign part given to them [18], [22], [23].

To address the issues of heterogeneity in requirements, the study [18] proposed a pivot model to assess the user requirements role in data design repository. In the study, firstly, ontological concepts are presented relating to the formalisms of user requirement. Secondly, a proposed model is intended to integrate with different semi-formal models. For the validation of approach, a case study is an implementation using model-driven approach.

Similarly, another study [22] presents a multi-perspective framework to manage requirements traceability using ontology as a knowledge management mechanism. To generate traceability relation ontology matching is applied as a reasoning mechanism. The associations are recognized by originating semantic similarity of ontological concepts representing of requirements elements. The precision and recall are used to compare the results of traceability relations identified by the framework and manually identified by the user, as validation of approach.

Furthermore, ontology in the context of software requirements is used to store information derived from requirements. Ontology is a structured way to organize information and data stores in ontologies can be accessed via queries. Software requirements specification ontology helps in capturing domain knowledge and knowledge of software under development. As software requirements specifications documents are used in all stages of development, an ontology is developed from software requirements specifications can support different development activities during the development process. These ontologies can be used to present domain knowledge in a processable format that may be helpful in test case generation to support software testing process [24], [25].

Likewise, the ontologies can be used, to reduce the adverse effects of factors such as ambiguous statements, insufficient specification, and changing requirements on requirements engineering process. The possible application of ontologies in RE process: To develop the requirements model, a paradigmatic way to write requirements, and acquisition of domain knowledge in a structured way [8], [26]. Sommerville defined that requirements specifications document is a description of the desired software characteristics specified by the customers [3].

Similarly, ontologies can be used to reduce the barrier of understanding that if machines are not recognizing the knowledge, ontology formalized that knowledge in a computer and human-readable form. It allows the user to find information based on purpose rather than syntax. A significant issue is the definition of standards to represent the underlying structures, the ontologies [17], [27]. The Resource Description Framework (RDF)⁴ allows defining taxonomies and relations between concepts. The RDF has three object types: resources, properties, and statements.

² <https://www.w3.org/TR/owl-features/>

³ <https://en.wikipedia.org/wiki/DARPA>

⁴ <https://www.w3.org/TR/rdf-concepts/>

Requirement Models	Concepts		Pivot Model				Relationships
			Actor	Requirement			
				Action	Criteria	Result	
Goal-oriented	Actor		✓				
	Goal	Task		✓			
		Metric					
		Result					
	AND/OR						✓
Combined Relationships						✓	
Use case Model	Actor						✓
	Use case	Action	✓				
		Result					
		Extension points					
		Condition					
	Generalization						
	Include						
Extend							
MERISE	Actor		✓				
	Event						
	Treatment	Action			✓		
		Result					
		Emission rules					
Synchronization						✓	

Fig. 3. The Core notations of three requirements formalisms [12], [18], [28].

IV. PROPOSED FRAMEWORK

In this study, we present a framework to combine different user requirements in MNC as shown in Fig. 4. This study is based on the use of ontologies to amalgamate the vocabularies used to software specification in different languages. Also, this study proposed a generic model to express the software requirements statements defined in different languages. At the final stage, framework checks the consistency of the integrated requirements. As shown in Fig. 3, the framework consists of subsequent phases: (1) Each partner from different countries defines its own Local Ontology (LO), and LO will be derived from Global Ontology (GO). These ontologies dealt with the heterogeneity of different languages vocabulary. (2) Different requirements formalism, are integrated using Pivot Model, (3) Local requirements transformation into Pivot Model.

A. Ontology as Shared vocabulary using Multilanguage WordNet

Domain ontology is defined as a formal, consensual, and referenceable dictionary for classes and properties of entities. The glossary in the framework highlights that objects or

association in the ontology domain is referred to language-independent identifier. To develop a global ontology that can coordinate with multiple languages and vocabularies, we suggest using Multilanguage WordNet. The approach relies on the following assumptions.

- A global ontology stores standards ontology such as, International Electrotechnical Commission⁵ (IEC) standard, ISO standards on the targeted domain by the application to be developed.
- WordNet is a sizeable English database, which helps to identify the part of speech tags such as nouns, verbs, adjectives, and adverbs. In framework Fig. 5, WordNet acts as lexical ontology. Multilanguage WordNet⁶ module is used to access different languages.
- Each designer develops an ontology in the local language, and global ontology is also inherited into the local ontology. Requirements defined in LO later exported regarding the GO.

B. Generic Requirements Formalism using Pivot Model

We combine different languages into formal requirements we apply pivot model. The designers are free to define his concepts in LO. Though each designer's formalism may be different; The analyst used three different DESIGNS such as use case of UML, goal oriented and MCT model of MERISE method. Fig. 3 shows the formalism identified by core notations used in the requirement statements defined as set of actions, results, criteria, and the relation between requirements. By merging the core notions and metamodeling, general requirements are formalized [18], [28].

A generic model to formally defines the Pivot model:

(Actors, Requirements, Relationships) where:

- The requirements can be defined as is a set of conditions specified by an actor. It can be described as {A, R, C, T}, where:
 - A is a set of activities to satisfy a requirement.
 - R is the results obtained if the requirements are satisfied.
 - C is set of criteria or conditions to quantified results.
 - T is the type of requirements.
- Relationships can be defined as is a set of relationships among requirement: Relationship = {Requires, constrains Refines, conflicts}

Formalization detail is adopted from study [18], [28].

C. Transformation of Formalized Requirements into Pivot Model

We used transformation model techniques [29], for mapping between local formalism such use case, goal, MCT Pivot model (Fig. 5).

⁵ <https://webstore.iec.ch/publication/21869>

⁶ <http://compling.hss.ntu.edu.sg/omw/>

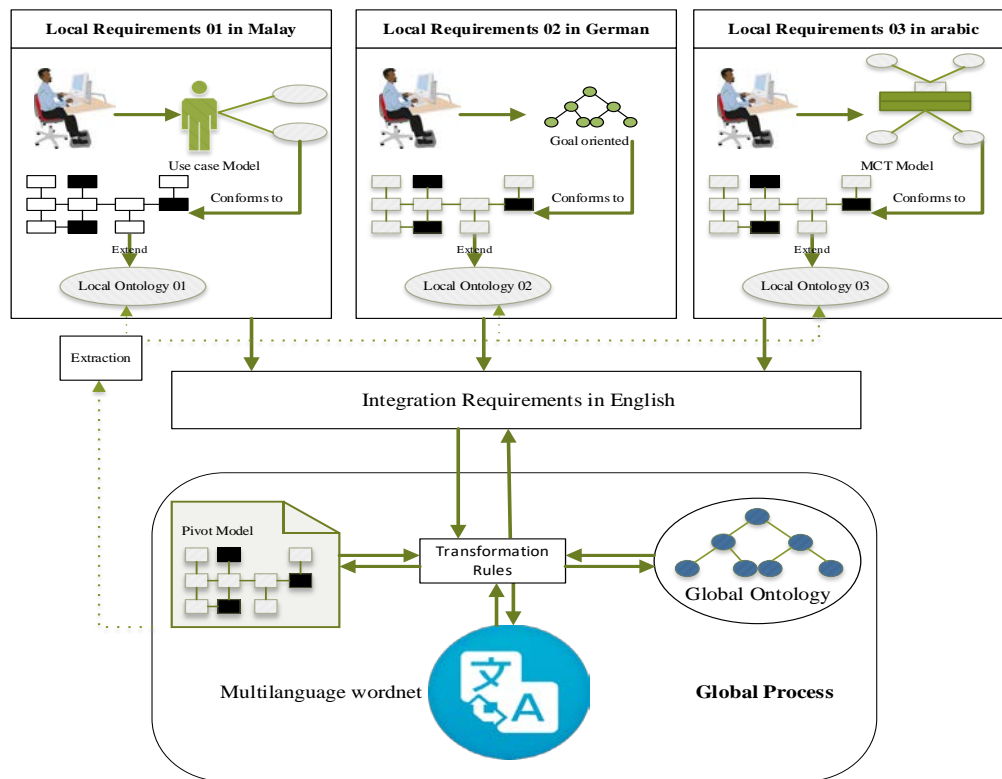


Fig. 4. Overview of proposed framework to integrate requirements.

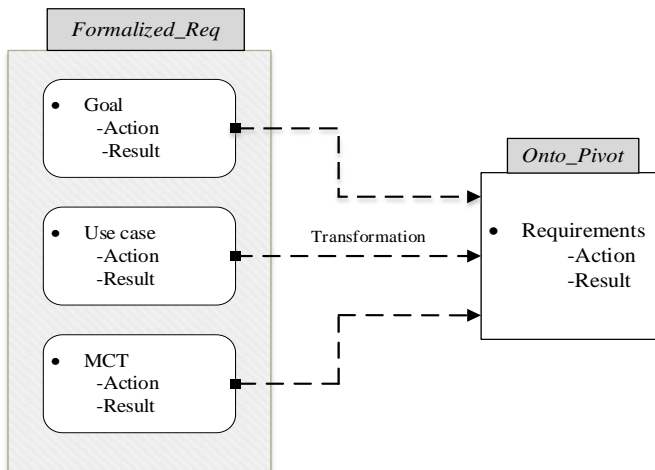


Fig. 5. Transformation of requirements.

As there are many tools and formats used to design specifications, we use ATL⁷ transformation language. ATL is a set of the model to model transformation which can integrate into Eclipse framework. It provides ways to produce a set of target models from a set of source models.

Requirements formalism is composed of three sub-formalisms as shown in Fig. 5. These instances are related to the local ontologies as described in the prior sections. The transformation rules are then implemented in these instances to translate them into Pivot model.

V. LIMITATIONS

This study does not conduct detail experiments on problem stated in Section 2. However, only presents a proposed framework to address the issue to integrate heterogenous requirement in the context of MNC.

VI. CONCLUSION AND FUTURE WORK

In this study, we discuss the problem that how heterogeneous requirements may be combined using local and global ontologies. The approach is based on general user requirement model that utilized three existing semi-formal languages: use cases of UML, MCT model of the MERISE method, and goal-oriented languages. This model is linked with local ontology and global ontology. These ontologies play the part of an amalgamated dictionary of Multilanguage Wordnet. Each designer from different locations develops the ontology using concepts and properties. By combining different user requirements in the context of MNC, our method focuses on minimizing the effect of heterogeneity of the quality of software product. In future, we will develop a system to integrate other languages.

ACKNOWLEDGMENT

The authors would like to express their deepest gratitude to Research Management Center (RMC), Universiti Teknologi Malaysia (UTM) and Ministry of Higher Education Malaysia (MOHE) for their financial support under Research University Grant Scheme (Vot number Q.J130000.2516.19H64).

⁷ <http://www.eclipse.org/at/>

REFERENCES

- [1] B. Nuseibeh and S. Easterbrook, "Requirements engineering: a roadmap," in Proceedings of the Conference on the Future of Software Engineering, 2000, pp. 35-46.
- [2] S. Hastie and S. Wojewoda, "Standish group 2015 chaos report-q&a with Jennifer Lynch," Retrieved, vol. 1, p. 2016, 2015.
- [3] I. Sommerville, "Software Engineering, Boston, Massachusetts: Pearson Education," ed: Inc, 2011.
- [4] R. W. Selby, Software engineering: Barry W. Boehm's lifetime contributions to software development, management, and research vol. 69: John Wiley & Sons, 2007.
- [5] M. Younas, D. Jawawi, I. Ghani, and R. Kazmi, "Non-Functional Requirements Elicitation Guideline for Agile Methods," Journal of Telecommunication, Electronic and Computer Engineering (JTEC), vol. 9, pp. 137-142, 2017.
- [6] D. Gašević, N. Kaviani, and M. Milanović, "Ontologies and software engineering," in Handbook on Ontologies, ed: Springer, 2009, pp. 593-615.
- [7] J. Z. Pan, S. Staab, U. Almann, J. Ebert, and Y. Zhao, Ontology-driven software development: Springer Science & Business Media, 2012.
- [8] V. Castañeda, L. Ballejos, M. L. Caliusco, and M. R. Galli, "The use of ontologies in requirements engineering," Global journal of researches in engineering, vol. 10, pp. 2-8, 2010.
- [9] K. Siegemund, E. J. Thomas, Y. Zhao, J. Pan, and U. Assmann, "Towards ontology-driven requirements engineering," in Workshop semantic web enabled software engineering at 10th international semantic web conference (ISWC), Bonn, 2011.
- [10] A. Mustafa, W. M. W. Kadir, and N. Ibrahim, "Automated Natural Language Requirements Analysis using General Architecture for Text Engineering (GATE) Framework," Journal of Telecommunication, Electronic and Computer Engineering (JTEC), vol. 9, pp. 97-101, 2017.
- [11] B. González-Baixauli, M. A. Laguna, and Y. Crespo, "Product lines, features, and MDD," in EWMT 2005 workshop, 2005.
- [12] I. Boukhari, S. Jean, I. Ait-Sadoune, and L. Bellatreche, "The role of user requirements in data repository design," International journal on software tools for technology transfer, pp. 1-16, 2016.
- [13] T. R. Gruber, "A translation approach to portable ontology specifications," Knowledge acquisition, vol. 5, pp. 199-220, 1993.
- [14] F. S. Parreiras, Semantic Web and model-driven engineering: John Wiley & Sons, 2012.
- [15] M. Hepp, "Ontologies: State of the art, business potential, and grand challenges," in Ontology Management, ed: Springer, 2008, pp. 3-22.
- [16] T. Diamantopoulos, M. Roth, A. Symeonidis, and E. Klein, "Software requirements as an application domain for natural language processing," Language Resources and Evaluation, vol. 51, pp. 495-524, 2017.
- [17] M. Ehrig, Ontology alignment: bridging the semantic gap vol. 4: Springer Science & Business Media, 2006.
- [18] I. Boukhari, L. Bellatreche, and S. Jean, "An ontological pivot model to interoperate heterogeneous user requirements," in International Symposium On Leveraging Applications of Formal Methods, Verification and Validation, 2012, pp. 344-358.
- [19] G. Pierra, "Context representation in domain ontologies and its use for semantic integration of data," in Journal on data semantics X, ed: Springer, 2008, pp. 174-211.
- [20] T. R. Gruber, "Toward principles for the design of ontologies used for knowledge sharing?," International journal of human-computer studies, vol. 43, pp. 907-928, 1995.
- [21] D. Sonntag, Ontologies and adaptivity in dialogue for question answering vol. 4: IOS Press, 2010.
- [22] N. Assawamekin, T. Sunetnanta, and C. Pluempitiwiriwajewj, "Ontology-based multiperspective requirements traceability framework," Knowledge and Information Systems, vol. 25, pp. 493-522, 2010.
- [23] S. J. Muhammad Irfan Marwat, Muhammad Arif Shah, Syed Zafar Ali Shah, "Towards optimization of software engineering ontologies," in Computer Applications and Information Systems (WCCAIS), 2014 World Congress on, 2014, pp. 1-6.
- [24] É. F. Souza, R. A. Falbo, and N. Vijaykumar, "Ontologies in software testing: a systematic," in VI Seminar on Ontology Research in Brazil, 2013, p. 71.
- [25] V. Tarasov, H. Tan, M. Ismail, A. Adlemo, and M. Johansson, "Application of inference rules to a software requirements ontology to generate software test cases," in International Experiences and Directions Workshop on OWL, 2016, pp. 82-94.
- [26] M. I. Marwat, S. Jan, M. A. Shah, and S. Z. A. Shah, "Towards optimization of software engineering ontologies," in Computer Applications and Information Systems (WCCAIS), 2014 World Congress on, 2014, pp. 1-6.
- [27] D. Brickley and R. V. Guha, "RDF vocabulary description language 1.0: RDF schema," 2004.
- [28] L. Bellatreche, N. X. Dung, G. Pierra, and D. Hondjack, "Contribution of ontology-based data modeling to automatic integration of electronic catalogues within engineering databases," Computers in Industry, vol. 57, pp. 711-724, 2006.
- [29] A. G. Kleppe, J. B. Warner, and W. Bast, MDA explained: the model driven architecture: practice and promise: Addison-Wesley Professional, 2003.

Effect of Service Broker Policies and Load Balancing Algorithms on the Performance of Large Scale Internet Applications in Cloud Datacenters

Ali Meftah

College of Computer and Information Sciences
King Saud University
Riyadh, Saudi Arabia

Ahmed E. Youssef^{1,2} and Mohammad Zakariah¹

¹College of Computer and Information Sciences
King Saud University, Riyadh, Saudi Arabia

²Faculty of Engineering, Helwan University, Cairo, Egypt

Abstract—Cloud computing is advancing rapidly. With such advancement, it has become possible to develop and host large scale distributed applications on the Internet more economically and more flexibly. However, the geographical distribution of user bases, the available Internet infrastructure within those geographical areas, and the dynamic nature of usage patterns of the user bases are critical factors that affect the performance of these applications. Therefore, it is necessary to compromise between datacenters, service broker policies, and load balancing algorithms to optimize the performance of the application and the cost to the owners. This paper aims at studying the effect of service broker policies and load balancing algorithms on the performance of large-scale Internet applications under different configurations of datacenters. To achieve this goal, we modeled the behavior of the popular Facebook application with the most recent worldwide users' statistics. Then, we evaluated the performance of this application under different configurations of datacenters using: 1) two different service broker policies, namely, closest datacenter and optimum response time; and 2) three load-balancing algorithms, namely, round robin, equally spread current execution, and throttled load balancer. The overall average response time of the application and the overall average time spent for processing a user request by a datacenter are measured and the results are discussed. This study would help service providers generate valuable insights on coordination between datacenters, service policies, and load balancing algorithms when designing Cloud infrastructure services in geographically distributed areas. In addition, application designers would benefit greatly from this study in identifying the optimal arrangement for their applications.

Keywords—Cloud computing; datacenters; load balancing algorithms; service broker policies; CloudAnalyst

I. INTRODUCTION

Cloud computing (CC) has become a prevalent technology in recent years. It provides a flexible and straightforward approach for maintaining and recovering information. Furthermore, it facilitates the collection of extensive information and the dissemination of records to various clients around the globe. Dealing with these vast information collections requires several strategies to enhance and streamline operations and provide attractive levels of execution to clients. CC incorporates computational and capacity benefits through a pay-per-use business model. Thus, it is exceptionally

desirable to business holders as it eliminates the provisioning overhead and enables organizations to begin very small and extend their assets only when needed [1], [2]. It is likely the main innovation that totally supplements the web, “cloud computing alludes to registering on the Internet, instead of processing on a desktop” [3].

With the advancement of the Cloud, it is now possible to build and host large scale applications such as social networking sites and e-commerce on the Internet more economically and more flexibly. Cloud Service Providers (CSP) are willing to provide large scaled computing infrastructure at a cheaper price (i.e., on pay per use basis) and in a very flexible manner (i.e., the users can scale up or down at will). However, several issues must be addressed to optimize the performance of applications such as the geographic distribution of the user bases, the available Internet infrastructure within those geographic areas, the dynamic nature of the usage patterns of the user base and how well the cloud services can adapt itself [21].

In practice, cloud computing clients request specific services and require that their demands be fulfilled ahead of schedule at limited costs. As shown in Fig. 1, the traffic routing between user bases (UB) and datacenters is controlled by a service broker that decides which datacenter should provide the service to the requests coming from each user base. Thus, the service broker controls traffic routing between user bases and datacenters. The load balancing algorithm determines which VM should be assigned to the next client request for processing according to different policies. For example, the round robin algorithm processes in a circular order by handling the process without priority, but equally spread current execution algorithm processes using priority. With the throttled algorithm, the client first requests the load balancer to find a suitable VM to perform the required operation. VMs should be assigned in a way that guarantees low response time and minimum transfer delay. Accordingly, service broker policies, load balance algorithms, and datacenters configuration are critical factors that influence the performance of applications.

This paper aims at studying the effect of service broker policies and load balancing algorithms on the performance of distributed large scale applications in cloud computing environments.

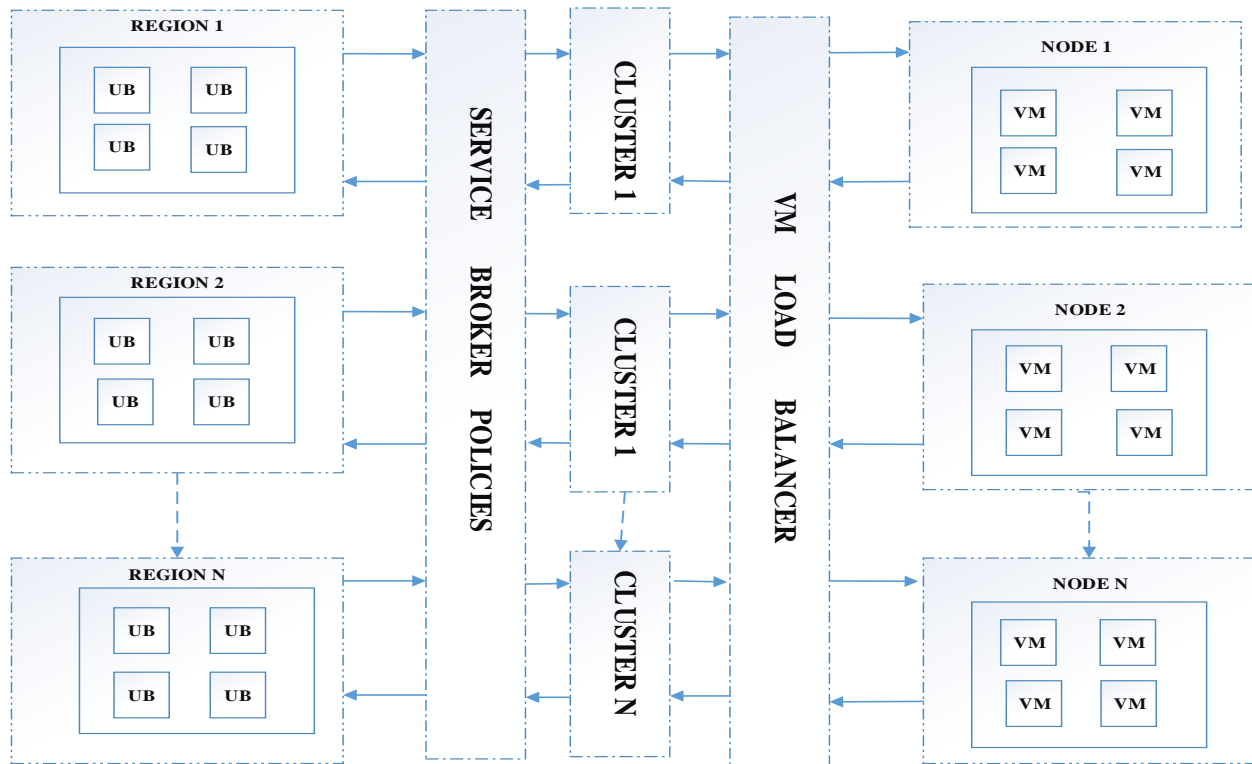


Fig. 1. Example cloud computing architecture [4].

To achieve this goal, we simulated the behavior of the popular Facebook application with the most recent worldwide users' statistics using CloudAnalyst tool. We examined the performance of this application in different scenarios (i.e., with different configurations of datacenters) and using: 1) two different service broker policies; namely, closest datacenter and optimum response time; and 2) three load balancing algorithms, namely, round robin, equally spread current execution, and throttled load balancer. The overall average response time of the application and the overall average time spent for processing a user request by a datacenter are recorded and the results are discussed.

This study helps CSP generate valuable insights on coordination between datacenters, service broker policies, and load balancing algorithms when designing Cloud infrastructure services in geographically distributed areas to optimize the application performance and the cost to the owners. In addition, application designers may use this study in identifying the optimal arrangement for their applications.

The rest of this paper is organized as follows: Section 2 provides background information on cloud computing and the most popular service broker policies and load balance algorithms; Section 3 describes the CloudAnalyst simulation tool; Section 4 shows and discusses our experimental results; Section 5 lists some related work; and finally, in Section 6, we give our conclusion and future work.

II. ESSENTIAL CONCEPTS

In this section, we review cloud computing and the most popular service broker policies and load balancing algorithms.

A. Cloud Computing

A cloud is a collection of IT resources, including hardware and software, deployed in a datacenter. The data center is built, operated, and managed by a Cloud Service Provider (CSP) which is an organization that provides cloud services. The provider may be an external (i.e., off-premise) provider to the consumer organization such as Amazon and Microsoft, or internal to the consumer organization (i.e., on premise), for example, the IT department. Cloud computing service models enable consumers to conveniently provision IT assets from a CSP in a way similar to a utility service such as electricity, wherein a consumer simply plugs in an electrical appliance to a socket, turns it on, and only pays for the amount of electricity used. The services provided by the cloud ranged from network-accessible data storage and processing, software development and deployment tools, to fully-featured software applications. CC paradigm brings uncountable benefits to users such as high scalability, rapid elasticity, and excellent availability of computing resources.

The National Institute of Standards and Technology (NIST) [23] defines different cloud service models. The first type of cloud service model is called Infrastructure-as-a-Service (IaaS) where the capability provided by CSP to the consumer is to provision processing, storage and networks. The underlying infrastructure (hardware) is managed solely by CSP. However,

the consumer has control over operating systems and is able to run and deploy software applications. The second type of service model is called Platform-as-a-Service (PaaS) where the consumer is able to develop and deploy onto the cloud infrastructure applications using programming languages, libraries, services, and tools provided by CSP. The consumer does not manage or control the underlying cloud infrastructure (i.e., network, storage, compute system) or operating systems, but has control over the deployed applications. The last cloud service model is called Software-as-a-Service (SaaS) where the capability provided to the consumer is to use the CSP's applications running on a cloud infrastructure. The applications are accessible from various client devices through either a thin client interface, such as a web browser (for example, web-based email), or a program interface. The consumer does not manage or control the underlying cloud infrastructure, operating systems, or even individual application capabilities.

B. Service Broker Policies

Service brokers handle traffic routing between user bases and datacenters by employing different policies such as:

Closest datacenter: The default routing policy routes traffic to the closest datacenter in terms of network latency from the source user base. This policy utilizes the concept of region proximity in selecting the datacenter to which the user request has to be directed. A region proximity list is maintained using the "lowest network latency first" criterion to set the order of occurrence of datacenters in the list. The datacenter that occurs first in the list, i.e., the closest data center, is selected to fulfill the request using this policy. In cases where more than one datacenter with the same latency are available, a random selection of the datacenters is made. This policy is, therefore, beneficial in cases where the request can be satisfied by a datacenter that is quite close or within the same region.

Optimal response time policy: This policy calls for the service broker to first identify the closest datacenter by making use of the network latency parameter, as in the previous policy. Then, the current response time is estimated for each datacenter. If the estimated response time is the one for the closest datacenter, then the closest datacenter is selected. Otherwise, the closest datacenter or the datacenter with the lowest response time is selected with a 50:50 chance.

C. Load Balancing Algorithms

Load balancing algorithms distribute the workload among server hubs within a datacenter as shown in Fig. 2. The primary difference between load balancing algorithms lies in the manner in which they choose the server hubs and direct new demands to those particular hubs. In fact, load balancing consists of distributing the workload to individual centers that may have fewer loads than others. Load balancing is applied to enable successful utilization of assets, to improve the reaction time of the assignment, and to eliminate situations in which some VMs are heavily loaded while others are only marginally loaded [5], [6]. Load balancing strategies are utilized by different server farms to adjust the workload among available VMs. Currently, existing load balancing algorithms include the Round Robin (RR) algorithm, equally spread current execution algorithm and the throttled algorithm.

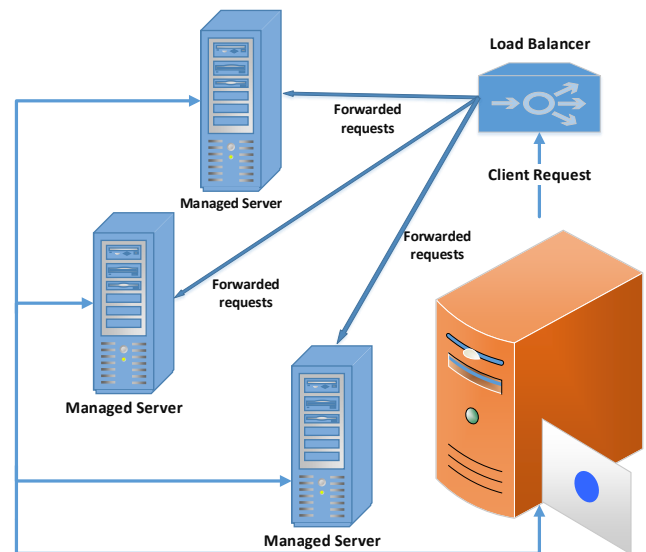


Fig. 2. Load balancing scenario [26].

Round Robin (RR) Algorithm: The RR algorithm, shown in Fig. 3, is one of the conventionally utilized algorithms that arbitrarily selects VMs. In the RR approach, time cuts are allocated to each assignment in a roundabout manner. Each undertaking is apportioned to available VMs in a roundabout request. However, there may be instances in which a few hubs are significantly loaded while others are only somewhat loaded. Thus, there are instances where the framework stack becomes irregular [7]. The RR algorithm is also a clear and static planning system that uses the guideline of time cuts, in which time is divided into various interims, and each VM is given a particular time cut or time interim [8], [9].

Equally spread current execution algorithm: The name of this algorithm suggests that it operates by equally spreading the execution load across different VMs [10]. It distributes the load randomly by first checking the size of the process and transferring the load to VMs that are only lightly loaded or which can handle the task easily in a small amount of time while maximizing throughput. This algorithm is a dynamic load balancing algorithm that determines the priority by checking the size of the process. It requires a load balancer that monitors the jobs to be executed.

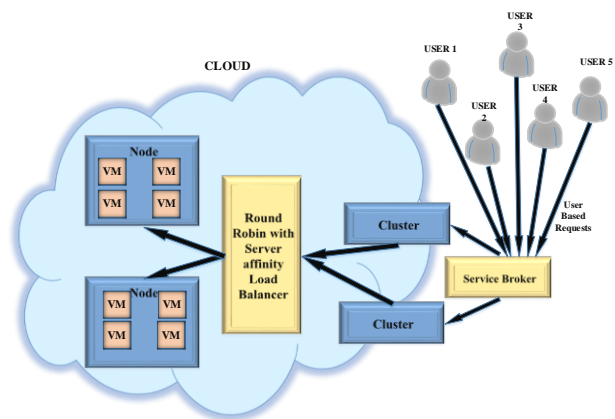


Fig. 3. Round robin load balancing [24].

The task of the load balancer is to queue the jobs and hand them to different VMs. The load balancer analyzes the queue frequently for new jobs and then allots them to the list of free virtual servers. The load balancer also maintains a list of tasks allotted to the virtual servers, which helps them to identify which VMs are free and need to be allotted new jobs.

Throttled Algorithm: The throttled algorithm, shown in Fig. 4, begins by allocating suitable VMs when the client sends a request to the load balancer. The load balancer maintains a file table of all the VMs together with their states, occupied or open mode. At the beginning, every VM is set to accessible mode. Then, the server controller guides the balancer for the next VM allotment, before it gets another demand. The balancer checks the table altogether until an important match of the VM was found. On the off chance that the ideal VM is discovered, then the balancer returns the ID of that particular VM to the server controller. Immediately, the server controller sends a demand to the VM with the specified ID. From that point onwards, the server controller sends a warning to balance the new distribution with the goal of refreshing the table. On the off chance there exists a case, when the VM is present, the balancer returns with that server information. When the VM is finished handling the doled out demand, the server controller gets a reaction cloudlet and it sends a message to the load balancer for VM de-assignment [11], [12].

III. CLOUDANALYST SIMULATOR

CloudAnalyst [12], [21] is a GUI-simulator developed to study the behavior of large scaled Internet applications in CC environment. Using CloudAnalyst, application developers or designers can determine the best strategies for selecting datacenters to serve specific requests, allocating resources among available datacenters, and the costs related to such operations. It is an open source simulation tool [19] that enables the recreation and assessment of the execution of different cloud administrations.

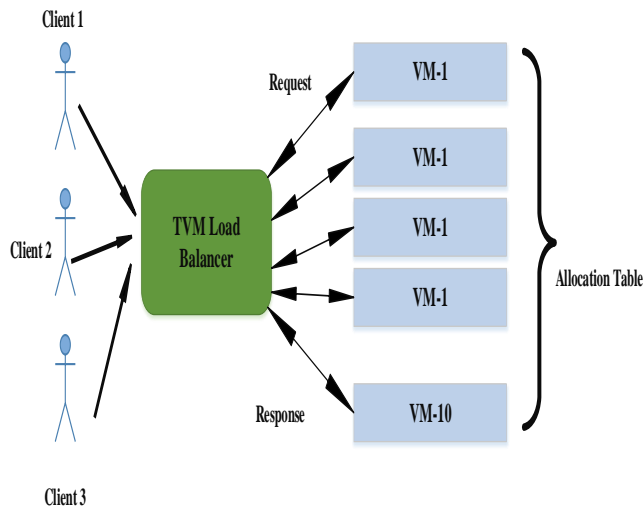


Fig. 4. Throttled algorithm load balancing [25].

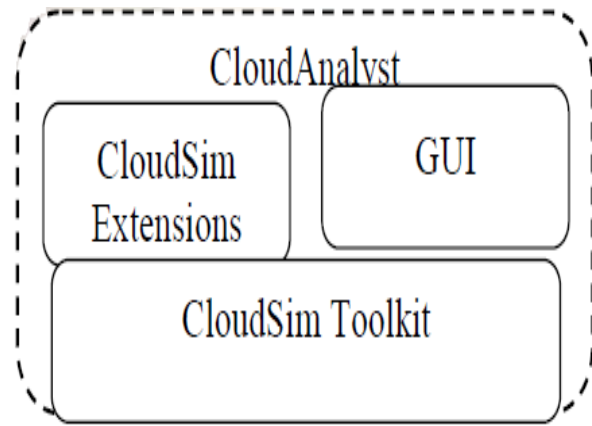


Fig. 5. CloudAnalyst architecture [12].

The CloudAnalyst is built on top of CloudSim tool kit [20], as shown in Fig. 5, by extending CloudSim functionality with the introduction of concepts that model Internet and Internet Application behaviors. It separates the simulation experimentations from programming tasks to enable users to quickly set up simulations and summarize results in useful formats. Different measures can be produced as an output of CloudAnalyst, such as, response time of the simulated application, usage patterns of the application, time taken by datacenters to service a user request, and the cost of operation.

IV. EXPERIMENTS

The purpose of our experiments is to show how different service broker policies, load balancing algorithms, and datacenter settings affect the performance of Internet large scale applications hosted on cloud datacenters. Social network services are one of the most popular Internet applications that vary with geographic location, sources of service requests, and time of day. They are large scaled applications that can benefit from cloud technology because they typically present non-uniform usage patterns [12]. One popular social network application today is Facebook, which has over 1.67 billion registered users around the world [22]. Table I shows Facebook usage and Facebook growth statistics by world geographic regions as of 30 June 2016. In our experiments, CloudAnalyst simulation tool was used to model and analyze the behavior of the Facebook application with the most recent worldwide users' statistical distribution.

A. Experiments Setup

We characterized six user bases representing the six primary areas of the world with the parameters depicted in Table II. For our simulation, we utilized a comparable speculative application with size equals 1/10 the size of Facebook. For simplicity, every client base was contained inside a solitary time zone and it was assumed that most clients utilized the application at night after working for approximately 2 hours. It was also assumed that 1% of the enrolled clients are web-based during peak hours at the same time and only a single tenth of that number of clients is online during off-peak hours. Moreover, every client makes another demand every 5 minutes when he or she is on the web. Table III shows the diverse CloudAnalyst parameter values used in our experiments.

TABLE I. FACEBOOK SUBSCRIBERS AND WORLD POPULATION STATISTICS AT JUNE 30, 2016 – UPDATE [22]

World Regions	Population (2016 Est.)	FACEBOOK users (as at 30 June 2016)
Africa	1,185,529,578	146,637,000
Asia	4,052,652,889	559,003,000
Europe	832,073,224	328,273,740
Latin America / Caribbean	626,054,392	326,975,340
Middle East	246,700,900	76,000,000
North America	359,492,293	223,081,200
Oceania / Australia	37,590,820	19,463,250
WORLD TOTAL	7,340,094,096	1,679,433,530

TABLE II. USER BASES USED IN THE EXPERIMENTS

User Base	Region	Time Zone	Peak Hours (Local time)	Peak Hours (GMT)	Simultaneous Online Users During Peak Hours	Simultaneous Online Users During Off-peak Hours
N. America	0	GMT - 6:00	7:00–9:00 pm	13:00–15:00	2230812	223081
S. America	1	GMT - 4:00	7:00–9:00 pm	15:00–17:00	3269753	326975
Europe	2	GMT + 1:00	7:00–9:00 pm	20:00–22:00	3282737	328274
Asia	3	GMT + 6:00	7:00–9:00 pm	01:00–03:00	5590030	559003
Africa	4	GMT + 2:00	7:00–9:00 pm	21:00–23:00	1466370	146637
Oceania/Australia	5	GMT + 10:00	7:00–9:00 pm	09:00–11:00	194633	19463

TABLE III. CLOUDANALYST PARAMETER VALUES USED IN THE EXPERIMENTS

Parameters	Assigned Value	
Simulation Duration	60 Min	
Virtual Machine (VM)	No. of VMs	Dependent on scenario
	Image size	10,000
	Memory	512 MB
	Bandwidth	1000 MB
Datacenter	Region	0
	Architecture	x86
	Operating system	Linux
	Virtual Machine Monitor (VMM)	Xen
	Memory per machine	2 GB
	Storage per machine	100 TB
	Available bandwidth per machine	1000000 MB
	Number of processors	4
	Processor Speed	10000 MB
VM Policy	Time shared	
User grouping factor in user bases	1000	
Request grouping factor in datacenters	100	
Executable instruction length	250	

B. Simulation Scenarios

We used the CloudAnalyst simulation tool to evaluate the performance of Facebook application described above in six different scenarios using: 1) two different service brokers policies; namely, closest datacenter (CDC) and optimum response time (ORT) and 2) three load balancing algorithms; namely; round robin, equally spread current execution, and throttled load balancer. Each scenario represents a different configuration for the datacenters. These scenarios are shown in Table IV.

TABLE IV. SIMULATED SCENARIOS

Scenario ID	Datacenter Settings
scenario1	One datacenter with 50 VMs, located at location 0
scenario2	Two datacenters with 25 VMs each, both located at location 0
scenario3	Three datacenters with 50, 75, and 100 VMs respectively, all located at location 0
scenario4	Three datacenters with 50, 75, 100 VMs, located at three different locations: 0, 1, and 2
scenario5	Six datacenters with 50 VMs, located at different locations: 0, 1, 2, 3, 4, and 5
scenario6	Six datacenters with 25, 25, 50, 50, 75, 100 VMs, located at locations: 0, 1, 2, 3, 4, and 5

In the first scenario, a single datacenter containing 50 VMs was employed to process all user requests from around the world. In scenario2, we added a second datacenter and reduced the number of VMs to 25 for each datacenter. In the third scenario, three datacenters with different numbers of VMs (50, 75, and 100) were used. In all these scenarios, the geographical location of the datacenter was location 0 (i.e., North America). The fourth scenario was the same as third one except that the datacenters were distributed over three locations, North America, South America, and Europe. In the 5th and 6th scenarios, we used six datacenters distributed over the six locations described in Table II. All datacenters in scenario5 have the same size (50 VM), while in scenario6, datacenters have sizes 25,25,50,50,75,100VM respectively.

C. Results and Discussion

During simulation of each scenario, readings during the 24 hours of the day for the response time of the application and the time spent for processing a user request by each datacenter are measured and the Min, Max, and overall average values are recorded as shown in Table V. A quick inspection to the results

revealed that, in general, the throttled load balancing algorithm outperforms the other two algorithms since its recorded overall average response time (ART) and overall average processing time (APT) are the shortest in all scenarios, this is also shown in Fig. 6 and 7. Table VI shows ART and APT of the throttled load balancing algorithm in all scenarios.

Table VI and Fig. 8 clearly show that the best values for ART and APT are obtained in scenario3 with the CDC service broker policy and in scenario6 with the ORT service broker policy. However, the cost of scenario3 is much less than scenario6 as can be noticed from Table IV. Finally, Fig. 9 shows a snapshot of the simulation results.

V. RELATED WORK

Various studies have been conducted on load balancing in the cloud computing environment. Randles et. al. [13] investigated the operation of three load balancing algorithms and discussed their disadvantages. They proposed arrangements for load adjusting. Khiyaita et. al. [14] reviewed cloud computing in depth and grouped load balancers in terms of their execution in a commonly circulated framework. They also discussed virtualization and examined the different interfaces in detail. Smith and Nayar [15] addressed the reasons for cloud appropriation and discussed how cloud computing helps different endeavors. Nandgaonkar and Raut [16] focused on the critical documentation required in distributed computing, administration from cloud suppliers, and operations in the cloud.

Padhy [11] assessed several outstanding current load adjusting algorithms that can be utilized. Jadeja and Modi [17] explored the building outline of distributed computing, its advantages, and a few drawbacks, for example, security, protection, and genuineness. Wickremasinghe et. al. [12] discussed the operation of a GUI-based apparatus called

CloudAnalyst that can be utilized for concentrating and executing gigantic scaled web applications. Finally, Mohialdeen [18] discussed various booking strategies, reviewed various scheduling algorithms in distributed computing, and discussed their application in cloud computing.

VI. CONCLUSIONS AND FUTURE WORK

This paper has studied the effect of service broker policies and load balancing algorithms on the performance of large scale applications in cloud computing environments. In order to accomplish this study, we have created different operation scenarios under different settings of datacenters and using two services broker policies; namely, closest datacenter and optimum response time and three load balancing algorithms; namely; round robin, equally spread current execution, and throttled load balancer. In our experiments, we modeled and analyzed the behavior of the popular Facebook application and evaluated its performance in each scenario using the CloudAnalyst simulation tools. The overall average response time of the application (ART) and the overall average time spent for processing a user request by a datacenter (APT) are measured. The results showed that the throttled load balancing algorithm outperforms the other two algorithms since its recorded readings (i.e., ART and APT) are the shortest.

This study would benefit CSP generate valuable insights on coordination between datacenters, service brokers policies, and load balancing algorithms when designing Cloud infrastructure services in geographically distributed areas to optimize the application performance and the cost to the owners. In addition, application designers may use this study to identify the optimal configurations for their applications. In the future, we plan to extent this work to include more parameters that may impact applications' performance and to examine other types of large scale applications.

TABLE V. EXPERIMENTAL RESULTS

			Response time of the application (ms)			Time spent in processing a request by a datacenter (ms)		
			Min	Max	Average	Min	Max	Average
Scenario 1	Closest Datacenter (CDC)	Round Robin	73.34	5398.37	1598	0.81	4770.27	1235.31
		Equally Spread	68	5396.20	1601.02	081	4771.94	1238.33
		Throttled	55.55	5279.04	970.86	0.81	4647.32	620.41
	Optimum Response Time (ORT)	Round Robin	71.25	5399.99	1597.85	1.47	4771.93	1235.05
		Equally Spread	72.13	5398.33	1601.43	1.47	4773.17	1238.65
		Throttled	55.46	5277.15	970.77	1.47	4660.90	620.27
Scenario 2	Closest Datacenter (CDC)	Round Robin	66.88	4450.79	1096.92	0.39	3851.89	740.94
		Equally Spread	59.77	4068.73	1096.53	0.26	3479.98	740.63
		Throttled	52.12	4295.62	730.02	0.82	3665.02	377.56
	Optimum Response Time (ORT)	Round Robin	58.66	4223.89	1277.05	0.41	3634.81	919.31
		Equally Spread	58.49	4487.76	1258.84	0.63	3909.48	901.27
		Throttled	53.27	4302.81	816.06	0.51	3686.20	463.70
Scenario 3	Closest Datacenter (CDC)	Round Robin	63.96	3038.29	785.29	0.96	2451.74	428.97
		Equally Spread	64.84	2150.56	716.48	0.95	1566.29	354.80
		Throttled	52.82	2821.96	577.52	0.78	2200.35	221.73
	Optimum Response Time (ORT)	Round Robin	58.21	3362.37	943.18	0.76	2762.19	587.08
		Equally Spread	64.64	3062.81	931.59	0.61	2502.19	575.49
		Throttled	52.15	3322.85	647.49	0.63	2687.67	292.63
Scenario 4	Closest Datacenter (CDC)	Round Robin	68.43	6507.21	1992.83	0.39	6111.12	1803.72
		Equally Spread	68.43	6605.26	1993.34	0.39	6209.48	1804.23
		Throttled	51.05	5364	1046.84	1.56	4973.69	868.30
	Optimum Response Time (ORT)	Round Robin	70.16	6561.78	1405.88	1.19	6159.04	1151.80
		Equally Spread	72.58	5460.62	1345.82	1.49	5044.94	1088.67

		Throttled	46.51	5370.93	778.58	1.81	4973.84	531.64
Scenario 5	Closest Datacenter (CDC)	Round Robin	47.71	2215.77	1209.78	1.48	2125.94	1141.38
		Equally Spread	47.71	2215.77	1209.81	1.48	2125.94	1141.40
		Throttled	47.71	2190.40	640.11	1.48	2087.26	576.26
	Optimum Response Time (ORT)	Round Robin	48.22	2129.41	842.68	0.95	1770.83	622.63
		Equally Spread	48.22	2214.30	835.95	0.67	1771.83	618.39
		Throttled	48.22	1976.28	520.21	0.95	1736.17	306.42
Scenario 6	Closest Datacenter (CDC)	Round Robin	48.49	2215.77	1194.19	2.26	2125.94	1125.31
		Equally Spread	48.49	2215.77	1194.24	2.26	2125.94	1125.35
		Throttled	48.49	2190.40	638.49	2.26	2087.26	574.66
	Optimum Response Time (ORT)	Round Robin	49.40	2215.75	825.95	1.25	1848.08	617.19
		Equally Spread	49.40	2324.31	818.39	1.25	1772.83	604.11
		Throttled	49.40	2070.27	517.48	1.42	1811.29	306.11

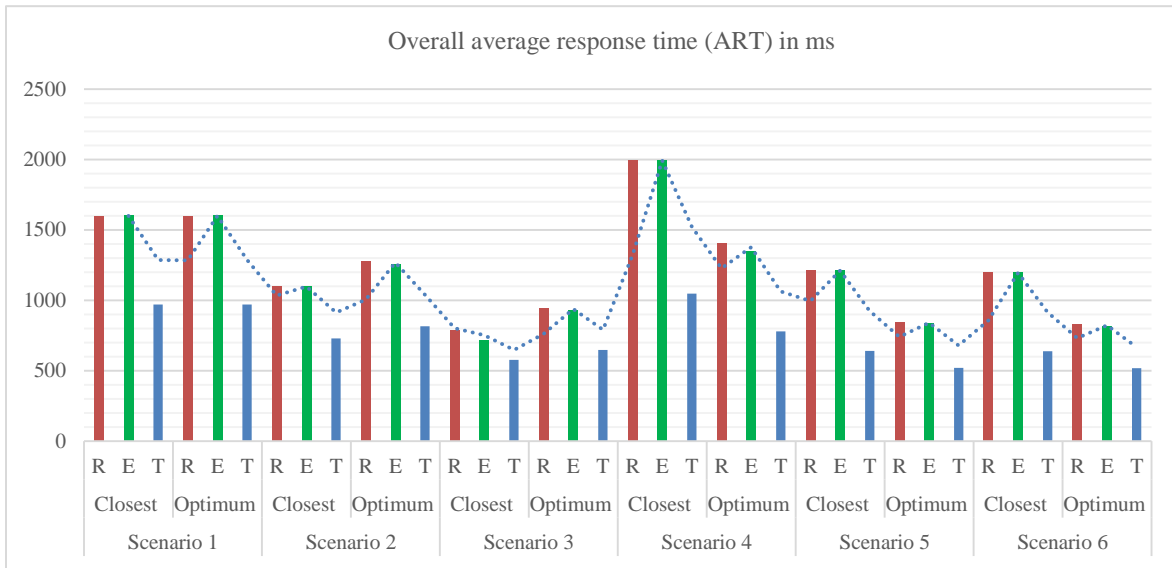


Fig. 6. Overall average response time (R: Round Robin algorithm, E: Equally Spread Current Execution algorithm and T: Throttled load balancing policy algorithm).

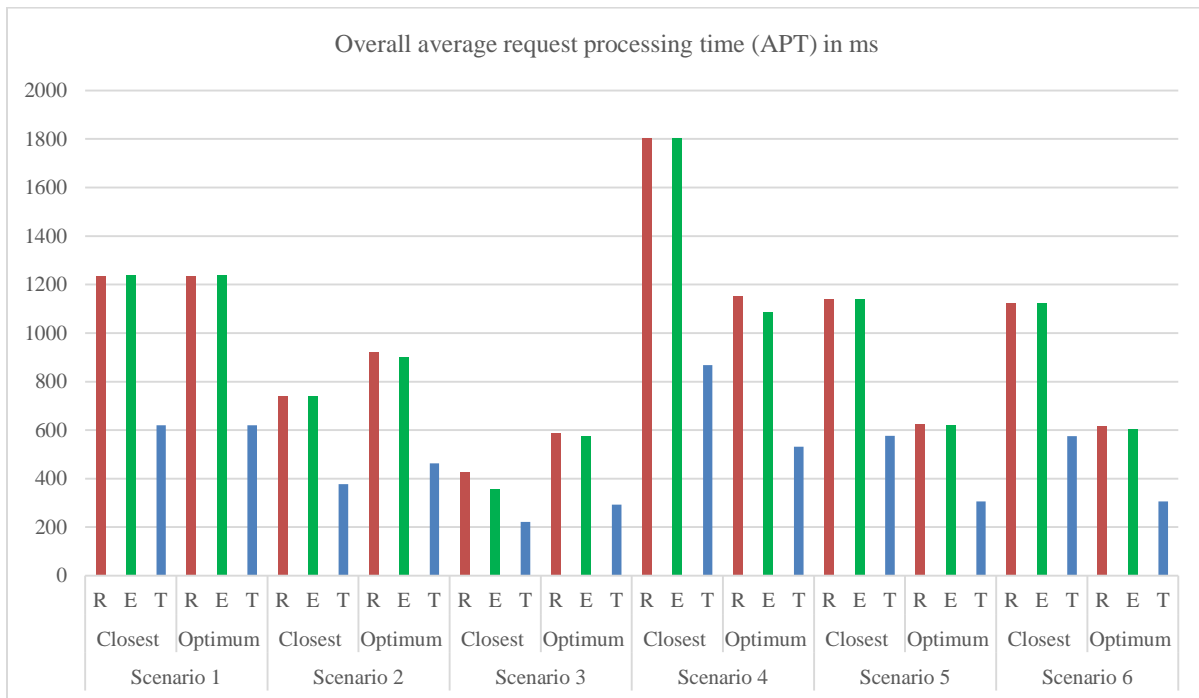


Fig. 7. Overall average request processing time (R: Round Robin algorithm, E: Equally Spread Current Execution algorithm, and T: Throttled load balancing policy algorithm).

TABLE VI. RESULTS FOR THE THROTTLED LOAD BALANCING ALGORITHM

Scenario ID	Scenario1		Scenario2		Scenario3		Scenario4		Scenario5		Scenario6	
	CDC	ORT	CDC	ORT	CDC	ORT	CDC	ORT	CDC	ORT	CDC	ORT
ART	970.86	970.77	730.02	816.06	577.52	647.49	1046.84	778.58	640.11	520.21	638.49	517.48
APT	620.41	620.27	377.56	463.7	221.73	292.63	868.3	531.64	576.26	306.42	574.66	306.11

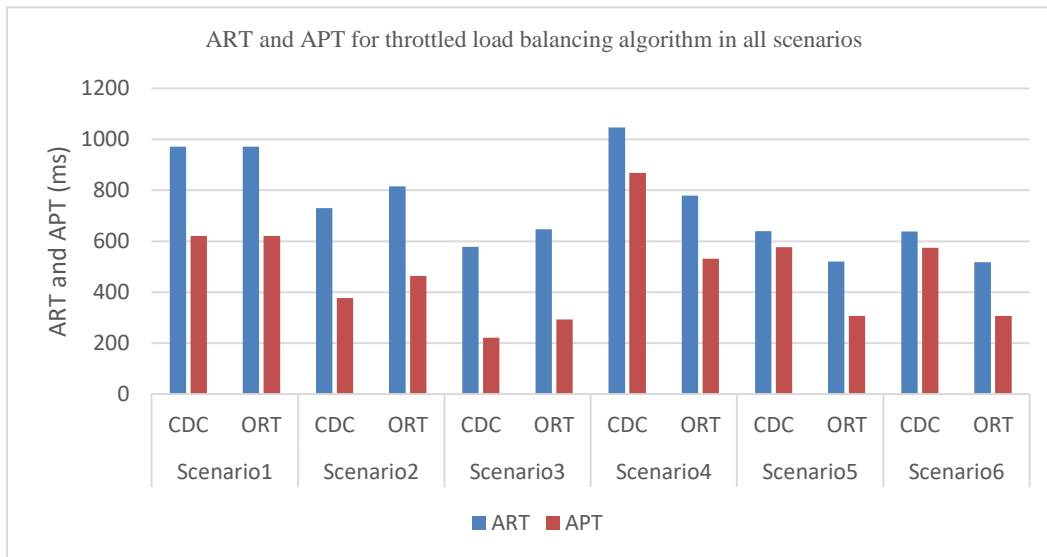


Fig. 8. ART and APT for throttled load balancing algorithm in all scenarios.

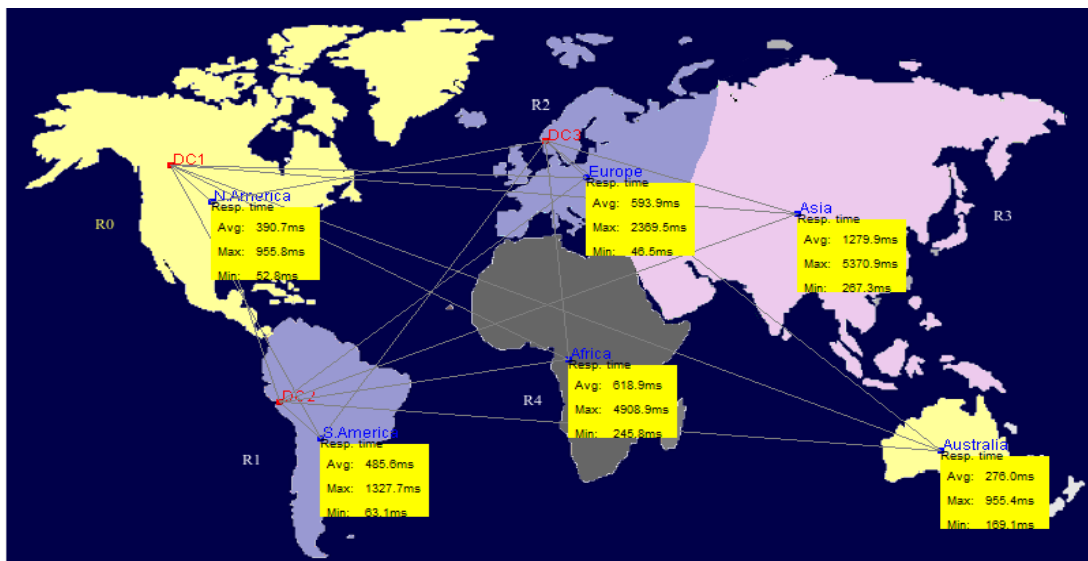


Fig. 9. Snapshot of simulation results.

ACKNOWLEDGMENT

This work is supported by the Research Center of College of Computer and Information Sciences (CRC) at King Saud University. The authors are grateful for this support.

REFERENCES

- [1] Q. Zhang, L. Cheng, and R. Boutaba, "Cloud computing: state-of-the-art and research challenges", Journal of internet services and applications, vol. 1, pp. 7-18, 2010.
- [2] M. D. Dikaiakos, D. Katsaros, P. Mehra, G. Pallis, and A. Vakali, "Cloud computing: Distributed internet computing for IT and scientific research," IEEE Internet Computing, vol. 13, 2009.
- [3] D. Blacharski and C. Landis, Cloud Computing Made Easy: Cary Landis, 2010.
- [4] S. Patel, R. Patel, H. Patel, and S. Vahora, "CloudAnalyst: "A Survey of Load Balancing Policies", International Journal of Computer Applications, vol. 117, 2015.
- [5] Z. Chaczko, V. Mahadevan, S. Aslanzadeh, and C. Mcdermid, "Availability and load balancing in cloud computing," in International Conference on Computer and Software Modeling, Singapore, 2011.
- [6] R. Lee and B. Jeng, "Load-balancing tactics in cloud," in Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), 2011 International Conference on, 2011, pp. 447-454.
- [7] A. Singh, P. Goyal, and S. Batra, "An optimized round robin scheduling algorithm for CPU scheduling," International Journal on Computer Science and Engineering, vol. 2, pp. 2383-2385, 2010.

- [8] D. Nusrat Pasha, A. Agarwal, and R. Rastogi, "Round Robin Approach for VM Load Balancing Algorithm in Cloud Computing Environment," *International Journal*, vol. 4, 2014.
- [9] M. M. D. Shah, M. A. A. Kariyani, and M. D. L. Agrawal, "Allocation of virtual machines in cloud computing using load balancing algorithm," *International Journal of Computer Science and Information Technology & Security (IJCSITS)*, vol. 3, pp. 2249-9555, 2013.
- [10] M. Nitika, M. Shaveta, and M. G. Raj, "Comparative analysis of load balancing algorithms in cloud computing," *International Journal of Advanced Research in Computer Engineering & Technology*, vol. 1, pp. 120-124, 2012.
- [11] R. P. Padhy, "Load balancing in cloud computing systems," *National Institute of Technology, Rourkela*, 2011.
- [12] B. Wickremasinghe, R. N. Calheiros, and R. Buyya, "Cloudanalyst: A cloudsimsim-based visual modeller for analysing cloud computing environments and applications," in *Advanced Information Networking and Applications (AINA)*, 2010 24th IEEE International Conference on, 2010, pp. 446-452.
- [13] M. Randles, D. Lamb, and A. Taleb-Bendiab, "A comparative study into distributed load balancing algorithms for cloud computing," in *Advanced information Networking and applications Workshops (WAINA)*, 2010 IEEE 24th International Conference on, 2010, pp. 551-556.
- [14] A. Khiyaita, H. El Bakkali, M. Zbakh, and D. El Kettani, "Load balancing cloud computing: State of art," in *Network Security and Systems (JNS2)*, 2012 National Days of, 2012, pp. 106-109.
- [15] J. E. Smith and R. Nayar, "Introduction to virtual Machines," *Virtual machines: versatile platforms for systems and processes*, pp. 9-10, 2004.
- [16] S. V. Nandgaonkar and A. Raut, "A comprehensive study on cloud computing," *International Journal of Computer Science and Mobile Computing*, vol. 3, pp. 733-738, 2014.
- [17] Y. Jadeja and K. Modi, "Cloud computing-concepts, architecture and challenges," in *Computing, Electronics and Electrical Technologies (ICCEET)*, 2012 International Conference on, 2012, pp. 877-880.
- [18] I. A. Mohialdeen, "Comparative study of scheduling algorithms in cloud computing environment," *Journal of Computer Science*, vol. 9, pp. 252-263, 2013.
- [19] <http://opensourceforu.com/2016/11/best-open-source-cloud-computing-simulators/>
- [20] R. N. Calheiros, R. Ranjan, C. A. De Rose, and R. Buyya, "Cloudsim: A novel framework for modeling and simulation of cloud computing infrastructures and services," *The 38th International Conference on Parallel Processing (ICPP)*, Vienna, Austria, 2009.
- [21] B. Wickremasinghe, "CloudAnalyst: A CloudSim-based Tool for Modelling and Analysis of Large Scale Cloud Computing Environments", MEDC Project Report, distributed computing project, CSSE dept., University of Melbourne. 2009.
- [22] <http://www.internetworldstats.com/>.
- [23] <http://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-145.pdf>
- [24] Mahajan, Komal, A. Makroo, and D. Dahiya. "Round robin with server affinity: a VM load balancing algorithm for cloud based infrastructure," *Journal of information processing systems* 9, no. 3 (2013): 379-394.
- [25] N. Pasha, A. Agrawal, R. Rastogi, "Round Robin Approach for VM Load Balancing Algorithm in Cloud Computing Environment", *IJARCSSE*, Volume 4, pages 34-39 Issue 5, May 2014.
- [26] <https://www.codeproject.com/Articles/430701/Singleton-Pattern-Load-Balancer-Demonstration>

Multiple Trips Pattern Mining

Riaz Ahmed Shaikh

Department of Computer Science
Shah Abdul Latif University
Khairpur, Sindh, Pakistan

Rafaqat Hussain Arain

Department of Computer Science
Shah Abdul Latif University
Khairpur, Sindh, Pakistan

Imran Memon

College of Computer Science
Zhejiang University, Hangzhou
Zhejiang, 310027, China

Kamelsh Kumar

Department of Computer Science
Sindh Madressatul Islam University
Karachi, Sindh, Pakistan

Hidayatullah Shaikh

Department of Computer Science
Shah Abdul Latif University
Khairpur, Sindh, Pakistan

Safdar Ali Shah

Department of Computer Science
Shah Abdul Latif University
Khairpur, Sindh, Pakistan

Abstract—In recent years, photograph sharing is one of the most mainstream web service, for example, Flickr, trip advisor and numerous other web services. The photograph sharing web services give capacities to include Geo coordinates, tags, and user ID to photographs to make photograph organizing easily. This study focuses on Geotagged photographs and discusses an approach to recognize user multiple trips pattern, i.e., common arrangements of visits in towns and span of stay and also elucidating labels that describe the multiple trips pattern. First, we segment collection of photos into multiple trips and categorize the photos manually based on photo trips into multiple trips, themes such as Landmark, Nature, Business, Neutral and Event. Our method mines multiple trips pattern for multiple trips theme categories. The experimental result of our technique beats other methods and accurate segmentation of photo collections into numerous trips with the 85% of accuracy. The multiple trips categorize about 91% correctly using tags, photo id, titles of digital photos, user id and visited cities as features. In last, we demonstrate the motivating examples showing an application with which one can find multiple trips pattern from our datasets and other different queries visit duration, destination and multiple trips' theme on trips.

Keywords—Multiple trips pattern mining; multiple trips classification; geo-tagging

I. INTRODUCTION

In recent years, there have been great innovations of digital camera and camera phone to sharing digital photos assigned tags, time stamps, geographical reference and visual information on web services such as Flickr, Facebook, Picasa and Panoramio and many other websites. Most popular internet application is a social networking service; millions of users share their information on these web services [1]-[5]. Users share their travel experience like photo and video of these social media web services. Photo sharing web services comprise billion of images accessible everywhere taken on earth. Increases volume of these images is defined various forms, including Geo tagged information, photographs, time and other variety of textual information. Increase volume of

Geo-reference, social media resource such as photos and videos document together with Geo tagged facilities. Textual metadata and temporal references; these enhanced the multimedia provided wealth data to solve the vision and media task. To discover graphic related information, knowledge about human societies and open new opportunities provide by multimedia. In computer vision research works most of people at single location make rich signified from image and also recognize image where an image was taken as well as image contents [6]-[12]. In this paper, we downloaded 10 million, data from Flickr using public Flickr API. This method automatically segments photo collection every user into multiple trips and also categorized multiple trips into the multiple trip's theme.

Following is our contribution in this paper:

- Our proposed algorithm empowers the novel knowledge detection, several trips pattern that depends on users' experiences, from the collection of Geo-tagged photos.
- We performed experiments for the performance evaluation of our proposed algorithm. The experimental result of our technique beats other methods and accurate segmentation of collection of photos into numerous trips with the 85% of accuracy. The multiple trips categorize about 91% correctly using tags, photo id, titles of digital photos, user id, and visited cities as features.

Table I used to outline the differences between our work and other significant studies. Here we show some essential issues about mining user trip and recommendation system, including: location Query Locations (LQ), timestamp (TS), GPS coordinates longitude/latitude (GPSC) and whether it considers the following ideas: Title assign to photo (TP), user travelling sequences (UTS), Trip classification (TC), categorization (CT), Visiting Time (VT) Recommendation system (RS), Geo tagged photo (GTP), user preference (UP), Distance (DI) sun database (SDB), Automatic detection (ATD) [3]-[18].

TABLE I. DIFFERENCES BETWEEN OUR WORK AND EXISTING WORKS

	LQ	TS	GPSC	TP	UTS	TC	CT	VT	RS	GTP	UP	DI	SDB	ATD
Zheng Liu et al. 2012	•	•	•	•			•		•	•		•		•
Takeshi et al. 2013					•									
Platt et al. 2002	•	•	•				•	•						
Toshihiko et al. 2013	•			•				•	•	•	•	•		
Jianxiong et al. 2010		•	•				•		•				•	•
Yan-Ying Chen et al. 2013	•			•	•		•		•	•	•	•		•
David et al. 2009	•		•	•					•	•				
Zheng et al. 2012		•	•		•	•	•	•	•	•				•
Richang et al. 2010		•	•				•			•				
Claudiu S et al. 2010			•	•			•		•	•				•
Yue Shi et al. 2013	•		•	•					•	•		•		•
Adrian et al. 2009			•	•			•			•	•			•
Zheng et al. 2011	•		•	•	•			•	•	•	•	•		•
This Paper	•	•	•	•	•	•	•	•	•	•	•	•	•	•

II. PROBLEM DEFINITION

A. Characteristics of Photo Trajectory

We utilize Geo tagging photographs in light of the fact that it is extremely rough than GPS trajectories. GPS trajectories acquire by GPS receiver through portable devices. You can see various users take photos during their movements because users take photos only when they see some attractive or interesting thing, in general users do not upload unattractive photos on web and various users visited different locations within the city. This is a great advantage of multiple trips mining because users take photos visiting different locations during their trips. The large data available on photo sharing web services, but it is hard to collect data from GPS trajectories to group of users.

B. Dataset

We downloaded photo collections of crawled users using Flickr public API and arranged data set such as time, user id, Geo coordinates, title and tags assigned to Geo tagged photos. Single photo assigned multiple tags at that time also possible same tag allocated to several photos.

We model the photo sets as $\mathbb{p} \triangleq \{p\}$ where p has characteristics of time captured $photot_p$, the ID of owner u_p location captured denoted by geo coordinate and assign title to photos. According to above description the photo collections each user represented by $\mathbb{p}_u \subseteq \mathbb{p}$, where all photos $p \subseteq \mathbb{p}$ are acquired by the similar user with the user ID u_p arranged in sequence that is \mathbb{p}_u can as temporal and spatial order.

We separated tags from Geo tag photos; use the variable ℓ to represent by the tag and \mathcal{L} to represent by all set of tags. Each tag assigned multiple photos and multiple photos have assigned similar tag. We representation \mathcal{L}_s to set of tags appear in any subset $\mathbb{p}_s \subseteq \mathbb{p}$ of photo. The subset of photo related with particular tag $\mathbb{p}_1 \subseteq \mathbb{p}$. Photo associate with tags in the subset of the particular tag represent by $\mathbb{p}_{,\ell}$. The users set related to photos in $\mathbb{p}_{,\ell}$ and \mathbb{u}_s as set of all users related to photos in \mathbb{p}_s .

C. The Data Model for Multiple Trips Pattern

Elementary elements for photos \mathbb{p} and tags \mathcal{L} , the trips formally define as:

$$\mathcal{PT} = \varepsilon, \mathcal{T}, \mathbb{L}, \mathcal{A} \quad (1)$$

Due to above equation the element defines as ε represented by locations sequence of visits user. Locations have hierarchy, such as scene, city, province and country. Our study focuses on the city level. Further, $\varepsilon = (\mathcal{C}_1, \dots, \mathcal{C}_n)$ represented cities sequences visited by users. $\mathcal{T} = (\mathcal{T}_1, \dots, \mathcal{T}_{n-1})$ represented visited sequence duration determined by two consequences cities. Then, $\mathcal{P} = \mathcal{P}_{\mathcal{C}_1, \dots, \mathcal{P}_{\mathcal{C}_n}}$ represented set of photo captured in visited cities and $\mathbb{L} = (\mathcal{L}_{\mathcal{C}_1}, \dots, \mathcal{L}_{\mathcal{C}_n})$ represented tags assign to photos. Finally, \mathcal{A} represented multiple trip's theme, users are chief object on the trips determined. Fig. 1 demonstrates the multiple trips of user which contain three trips 1) visiting famous landmarks in China, 2) everyday gathering with friends in Germany, 3) and last shopping in Paris. We aim to detect multiple trips pattern on each trip into multiple trips themes. As the frequently visited cities and duration of visiting and characteristic of tags order for multiple trip pattern. We describe the multiple trips, tag and characteristic trips pattern, tags denoted what to enjoy and what to see during to multiple trips pattern [19]. Content analysis using shape and image edge detection are useful techniques for object recognition [20], [21].

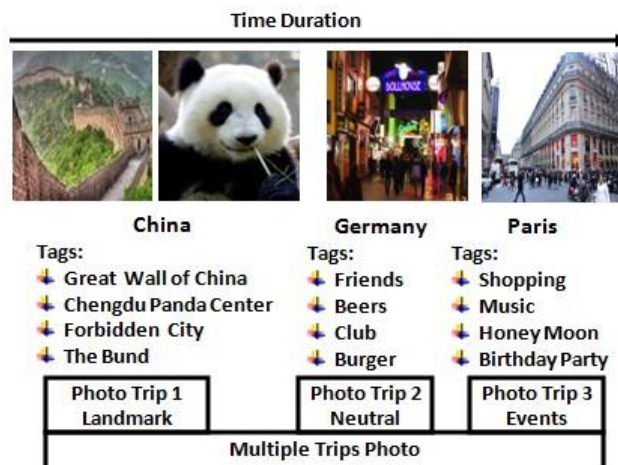


Fig. 1. User's multiple trips photo.

III. APPROACH

The idea to detect event from photo collections is to use gaps of capture time duration because the information about the location of photos did not available at that time. The peoples' movement reflects directly transition of photo capture location and also necessary to efficiency improve event detection. The Geo-tagged photos each tag assigns well define the trip. The user assigns similar labels on groups of successive photos [1] which ought to uphold photo collections algorithm of segmentation. We can use to capture the location and time gaps, tag on segmentation of photo collections. In [3] proposed algorithm separate photo into the event. It detects the time, which regarded the event changes. The algorithm first sort photo time captured if the change time gaps between two

photos, then sorted list as change event consider if it's much higher than local gap:

$$g_{N \geq k + \frac{1}{2d+1} \sum_{i=-d}^d g_{N+i}} \quad (2)$$

The appropriate threshold is K and "d" is a size of window. If N + i to photograph past the end of photo collection, expression is overlooked and the denominator 2d+1 is decreased for each disregarded term to keep the average standardized. The time_{gap} and distance_{gap} compute the distance and transition time of two successive photos. We calculate the gaps between two consecutive photos: p_k and p_{k+1} show the distance and time gaps following equation:

$$\begin{aligned} \text{distance}_{\text{gap}} &= \log(D\emptyset) \\ \text{time}_{\text{gap}} &= \log(t_{p_{k+1}} - t_{p_k}) \end{aligned} \quad (3)$$

Finally, the gap g_N Is considered as changing trip when much higher the local gap average calculates by

$$g_{N \geq k + \frac{1}{2d+1} \sum_{i=-d}^d g_{N+i}} \quad (4)$$

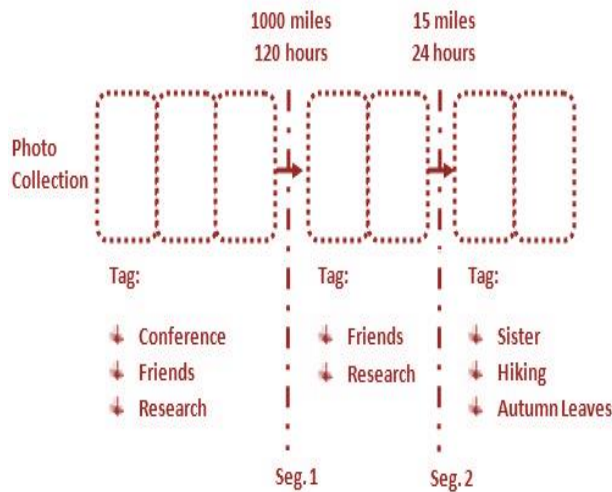


Fig. 2. Example of photo collection segmentation.

In Fig. 2 shown the example where we consider photo collections containing of photo capture when attending conference and user research in the lab with their friends. Our algorithm of segmentation collection as Seg.1; the large gaps caught time and area while the labels are comparable with taking after photograph. In Seg.2, segmentation tags allocated to photographs in the final group did not share labels to previous photos.

IV. PROPOSED MODEL: THE MULTIPLE TRIPS CLASSIFICATION

Proposed model is very adequate for the multiple trips classification, after segmentation of photograph collection and connected with city names and duration of visit to various trips we characterize numerous trips as indicated by different trip's theme. The numerous trips theme empowers clients to find various trip patterns as per their objectives on the trip, and distinct theme result's trip in dissimilar descriptions and patterns, despite the fact that we visit the same sites.

A. Multiple Trips Theme

We surveyed various websites of travel forums, agencies and blogs regarding travel experience. On the basis of our survey we categorize various trips in five key objectives.

Landmark: Visit of popular sightseeing, world heritage, for example, the Great Wall of China, Chengdu Panda center, etc. as shown in Fig. 1.

Nature: Visit of such places which are popular for rich nature, for exp. Tiananmen Square in Beijing and Beihai parks, to communicate with nature.

Event: Visit for attending events, like Ninfeng in Shanghai and a wedding ceremony.

Neutral: Everyday small trips to visit adjacent areas for example, gatherings with friends and relatives at hotel or bar, and to enjoy delicious food with family such as famous hotspots in Chengdu.

Business: To visit places for the business purpose.

Note that we collected multiple trips only those users who have visited more than three trips. We describe these five types of themes of trip and develop an algorithm to categorize various trips into theme of multiple trips. We ought to note that user's trips going to the similar areas may belong with different classes, since individuals go to the similar areas with distinctive aims

B. Classifier

For the classification of multiple trips Support Vector Machines (SVMs) has been used. To prepare support vector machine we assigned labels manually to 20000 randomly selected various trips photographs by their multiple themes of trips. We categorize various trips into five classes and prepare a SVM [15] for each of the various themes of trip. To carry out numerous trips photograph categorization on a set of test, we execute five classifiers on different trips and pick the classifications with most elevated score (most prominent optimistic separation from the support vector machines isolating hyperplane). We performed separation of the labeled various trips into partitioning so as to prepare and testing sets various tours of the same clients so as to evade the likelihood that related numerous trips of the similar individual show in training sets.

V. EVALUATION

10 million Geo-tagged photographs are downloaded from Flickr website. The several tags allocated to the photographs are 82.9 million whereas the quantity of novel tags between them is 4.7 million. Quantity of photo proprietors is 79.7K, between the 49.1K clients whose home town areas are accessible (on Flickr, clients can enlist their locations), around 18% of photographs are caught inside of clients' city, around 60% of photographs are caught inside of clients' country, 22% of photographs are caught outside the county of client. As the sizes of towns are not stable, we can't straightforwardly state that all 17% of photographs were caught around homes of the client. Although, the measurements demonstrate that the users are capturing photographs in different areas, which incorporates both places adjacent and abroad. In this way, we

hope to have the capacity to recognize a various trips pattern from the collection of photographs.

A. Segmentation Accuracy

In this section, we examined an accuracy of segmentation of the collections of photographs and the effect of the tags and distance/time gaps.

B. Evaluation Metrics

In this section, we compare segmentation algorithm with existing algorithms (Platt et al 2002) and three differences of our proposed algorithm which are given below: 1) suppose only the distance gaps, 2) one tag only, 3) supposes captured time and distance gap both (beta = 0 parameter) in the sequence to contribute from every element. During experiment, set parameter such as (d, k, alpha, beta) of these methods are the best as performance in a different training set of photographs. The fact is, we arbitrarily select eighteen clients who seem like active clients of Flickr and physically divided their photograph collection into several trips. Least number of testing user's about 2240 and highest number of test user's is about 5022 with the aggregate number of test user's is 115208.

We use compare metrics in existing work of occasion clustering for collections of photo [18], [19], following equations is for precision and recall of detecting boundaries for event.

Precision = (no.of corectlydet ectedboundaries) / (Total No.of det ectedboundaries) (5)

Recall = (no.of corectlydet ectedboundaries) / (Total No.of det ectedboundaries) (6)

The trade-off among recall and precision is that lower recall end to result in high precision and vice versa.

The high precision joins some real truth diverse trips into several trips and high call tends to over segmentation for the collections of photo. To maintain balance between the recall and precision is vital. F measurement (F1) is the harmonic mean of recall and precision, (3) used for the assessment of balance between precision and recall. The range of values is from 0 to 1 for F1, whereas a higher value is good.

F1 = (2(precision)(recall)) / ((precision) + (recall)) (7)

C. Results and Discussions

Table II shows the result of photo collections segmentation. To over the segment photo collections means higher recall and lower precision values Platt's algorithm results show that segmented photo collection into a smaller piece of the trip. On the other side, tend to combine various trips into multiple trips the algorithm count distance gaps only and higher precision lower recall resulting shows us. The based on a distance gap algorithm and plats algorithm; an algorithm considers both distance gap and time gap. The algorithm outperforms the other algorithm by considering the title and assign tags to

photos resulting shows F1 best. It turns out complement other algorithm, the performance of tags counting algorithm were poor. By itself, the many tags on the photographs are insufficient for segmentation. Being merge with different frameworks the tags offer additional advantage to them, since straightforwardly reflect what is going on in the photographs.

TABLE II. SEGMENTED PHOTO COLLECTION

Table with 4 columns: Methods, Precision, Recall, F1. Rows include Tag, User id, Time, Distance and Time, Distance, and Our proposed method.

The spatial feature based matrix of confusion is shown in Table III.

TABLE III. SPATIAL FEATURE BASED CONFUSION MATRIX (%)

Confusion matrix table with 6 columns: (La), (N), (E), (NE), (B) and 6 rows: Landmark, Nature, Event, Natural, Business.

IV. CONCLUSION AND FUTURE WORK

The described idea of multiple trips and proposed multiple trips pattern mining algorithms that detected novel multiple trips from the collection of Geo-tagged photos on web. First photo collections can be dividing into multiple trips by segmentation and categorize them multiple trip's theme. After this identify several trips patterns as arrangements of most often traveled areas and their visit spans. Mine labels, tolerating their geological scope to incorporate a delineation of a various trip pattern with the objective that clients interpret them as what to expect and see if grasp the patterns. Various trips arrangement, results demonstrate that texture components of digital photos, i.e., titles and tags of photos are the very competent components and spatial components of users' trips, such as visited towns, can supplement the textural component. Utilizing these elements can categorize around 91% of numerous trips accurately. In future work different ANN classifiers may be used.

REFERENCES

[1] Crandall, D., J., Backstrom, L., Huttenlocher, D., Kleinberg, J., "Mapping the world's photos", International Conference on World Wide Web, pp. 761-770, 2009. [2] Zheng, Y., T., Zha, Z., J., Chua, T., S. "Research and applications on Geo-Referenced multimedia a Survey" Multimedia Tools and Applications, Vol. 51(1), pp. 77-98, 2010. [3] John C. Platt, M., C., Brent A., F., "Photo TOC: Automatic Clustering for Browsing Personal Photographs", Information, Communications and Signal Processing, pp. 1-5, 2003.

- [4] Zheng, Y., T., Zha, Z., J., Chua, T., S., "Mining Travel Patterns from Geo-Tagged Photos", *ACM Transactions on Intelligent Systems and Technology*, Vol. 3(3), pp. 1-18, 2012.
- [5] Xiao, J., Hays, J., Ehinger, K., A., Oliva, A.,; Torralba, A., "SUN database: Large-scale Scene Recognition from Abbey to Zoo", *Computer Vision and Pattern Recognition (CVPR)*, pp. 3485-3492, 2010.
- [6] Yan-Ying Chen, A., J., C., Winston, H., H., "Travel Recommendation by Mining People Attributes and Travel Group Types from Community-Contributed Photos", *IEEE Transactions on Multimedia*, Vol. 15(6), pp. 1283-1295, 2013.
- [7] Liu, Z., Yan, H., Han, H., "Mining Large-Scale Social Images with Rich Metadata and Its Application", *Journal of Software*, Vol. 7(4), pp. 749-756, 2012.
- [8] Yamasaki, T., Gallagher, A., Chen, T., "Personalized Intra- and Inter-City Travel Recommendation Using Large-Scale Geotags", *GeoMM '13 Proceedings of the 2nd ACM international Workshop on Geotagging and its applications in multimedia*, pp. 25-30, 2013.
- [9] Shi, Y., Serdyukov, P., Hanjalic, A., Larson, M., "Nontrivial landmark recommendation using geotagged photos", *ACM Transactions on Intelligent Systems and Technology*, Vol. 4(3), 2013.
- [10] Claudiu, S., Firan, M., G., Nejd, W., Paiu, R., "Bringing Order to your Photos: Event-Driven Classification of Flickr Images Based on Social Knowledge", *CIKM'10*, pp. 189-198, 2010.
- [11] Popescu, A., Grefenstette, G., Moëllic, P., "Mining Tourist Information from User-Supplied Collections", *Proceedings of the 18th Conference on Information and Knowledge Management*, pp. 1713-1716, 2009.
- [12] Hong, R., Li, G., Nie, L., Tang, J., Chua, T., "Exploring Large Scale Data for Multimedia QA an Initial Study", *CIVR '10 Proceedings of the ACM International Conference on Image and Video Retrieval*, pp. 74-81, 2010.
- [13] Kurashima, T., Iwata, T., Irie, G., Fujimura, K., "Travel Route Recommendation Using Geo-tagged Photos", *Knowledge and Information Systems*, Vol. 37(1), pp. 37-60, 2012.
- [14] Giannotti, F., Nanni, M., Pedreschi, D., "Efficient Mining of Temporally Annotated Sequences", *SAC '06 Proceedings of the 2006 ACM symposium on Applied computing*, pp. 593-597, 2006.
- [15] Tsochantaridis, I., Hoffman, T., Joachims, T., Altun, Y., "Support Vector Machine Learning for Interdependent and Structured Output Spaces", *Proceedings of the 21st International Conference on Machine Learning, Banff*, pp. 104, 2004.
- [16] Hao, Q., Cai, R., Wang, X., Yang, J., Pang, Y., Zhang, L., "Generating Location Overviews with Images and Tags by Mining User-Generated Travelogues", *MM '09 Proceedings of the 17th ACM international conference on Multimedia*, pp. 801-804, 2009.
- [17] Rattenbury, T., Naaman, M. "Methods for Extracting Place Semantics from Flickr Tags", *ACM Transactions on the Web*, Vol. 3(1), 2009.
- [18] Cooper, M., Foote, J., Girgensohn, A., Wilcox, L., "Temporal Event Clustering for Digital Photo Collections", *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, Vol. 1(3), pp. 269-288, 2005.
- [19] Alexander, C., Loui, S., M., Savakis, A., "Automated Event Clustering and Quality Screening of Consumer Pictures for Digital Albuming", *IEEE Transactions on Multimedia*, Vol. 5(3), pp. 390-402, 2003
- [20] Riaz Ahmed Shaikh, Jian-Ping Li, Asif Khan, Kamlesh Kumar, "Content Analysis Using Shape and Spatial Layout with Markov Random Field", *Indian Journal of Science and Technology*, Vol 9(7), pp. 1-6, February 2016.
- [21] Kamlesh Kumar, Jian-Ping Li, Saeed Ahmed Khan, "Image Edge Detection Scheme Using Wavelet Transform", *International Conference on Wavelets Active Media and Information Processing*, vol. 11, pp. 261-256, December. 2014.

E-shape Multiband Patch Antenna for 4G, C-band and S-band Applications

Mehr-e-Munir¹, Khalid Mahmood², Saad Hassan Kiani³

Department of Electrical Engineering, Iqra National University, Peshawar Pakistan^{1,3}
Department of Electrical Engineering, University of Technology, Pakistan²

Abstract—In this study, a new E shape mounted on minowaki island patch antenna on FR4 substrate is presented for communication systems applications. With insertion of shortening pin between patch and ground plane, the proposed structure resonated on 6 frequencies; hence producing Hex-band response with good realized gain and directivity radiation values and patterns. Co axial cable is used as means of excitation to excite proposed structure with minimum impedance mismatch losses. The proposed design is miniaturized up to 60.66% and can be used for GSM, GPRS, 4G, WLAN and other S-band and C-band applications.

Keywords—Minowaki island patch; miniaturization; E shape; gain; directivity

I. INTRODUCTION

The development in wireless communication, microwave technology is increasing day by day with the passage of but these technologies requires smaller size of antennas which can be used in a numbers of applications, such as in 4G, S-band, C-band, mobile applications and other applications. As many telecommunication systems and radar communication system used several frequencies i.e. dual band antenna is more beneficial single band [1].

A new proposed antenna patch structure is suggested in this study for wireless communication application, the design results the multiband frequencies and reduced size antenna [2]. The size reduction is obtained by various ways such as using high permittivity substrate which gives good results of miniaturization but the cost of high permittivity substrate is expensive and unsuitable for low cost consumers' application [3]. Meta materials and magneto dielectrics are used for size reduction purpose but with these materials and dielectrics the miniaturization is archived but disadvantage is that this miniaturization gives the lower gain also the cost of these materials are expensive and complex to manufacture [4]. The Minowaki island shape or fractal patch results in miniaturization but this lower size of antenna gives lower gain [5].

Defected ground structure such as H, L, U shaped is used for size reduction of antenna which gives good gain and return loss but these structures gives very narrow bandwidth [6]. Defected patch structure such as H, U, Pi, E shaped is used for multiband operation of antenna which gives good return loss but these structure gives very lower gain [7]. Shorting pin method is also used for miniaturization purpose also as

direction of current is changed with respect to location of the pin [8].

The use of Artificial Magnetic conductors (AMC), Split ring Resonators (SRR) has also been found effective in shifting the fundamental resonating frequency to lower levels but these techniques cause the gain to diminishing levels after few iterations. High cost and complex geometry also keeps researcher on edge at using Meta materials [9].

This paper presents a novel study of multiband enabled reduced sized antenna by introducing Defected ground structure in ground plane and E shape edged on resonating fractal patch. FR4 is used as a substrate having permittivity 4.2 which is easily available in the market, better efficiency, high bandwidth, low water absorption [10], [11]. Size reduction of patch antenna is obtained by mounting slots on ground plane while edging the patch with fractal and E slots has produce Hex Band response. All of the resonating frequencies have shown satisfactory performance parameters results with good radiation patterns and minimum mismatch losses. The proposed designed is carried out CST 2014 and can be used for S- and C-band applications.

II. ANTENNA CONFIGURATION

A. Width of Patch

Following formula is used for deriving patch width.

$$W = \frac{c}{2f_0 \sqrt{\frac{\epsilon_r + 1}{2}}} \quad (1)$$

B. Length of Patch

Following formula is used for deriving patch length.

$$L = L(e_{ff}) - 2\Delta L \quad (2)$$

Where,

$$L(e_{ff}) = \frac{c}{2f_0 \sqrt{\epsilon(re_{ff})}} \quad (3)$$

And

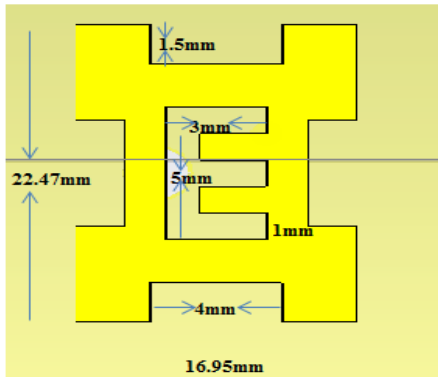
$$\epsilon(re_{ff}) = \frac{\epsilon_r + 1}{2} + \frac{\epsilon_r - 1}{4} \left(1 + \frac{12h}{W}\right)^{-1} \quad (4)$$

A conventional antenna of 4.1GHz is designed after calculating width and patch of antenna. Regarding patch and ground structural dimensions with all slots information is covered in Table I.

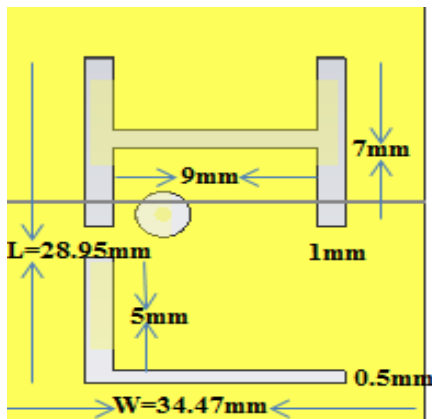
The resonating patch length and width is taken 16.95 and 22.47mm, respectively. Fractal slots are introduced having length and width 5.1 and 2.3mm. In center of Resonating patch, and E shape is mounted with 3mm length and 1mm of width. FR4 height is taken 2mm and height of ground and patch is taken to be 0.8 and 0.245mm, respectively.

For miniaturization, ground plane is introduced with H and L Slots. The length and width of slots is taken as 7 and 9mm and L slot length is taken as 5mm and width 9mm.

Studies have shown that E shape slotting on various patch locations exhibit different response as the current is circulated around its corners [12]. The detailed geometry of proposed structure is shown in Fig. 1(a) and 1(b), respectively.



(a)



(b)

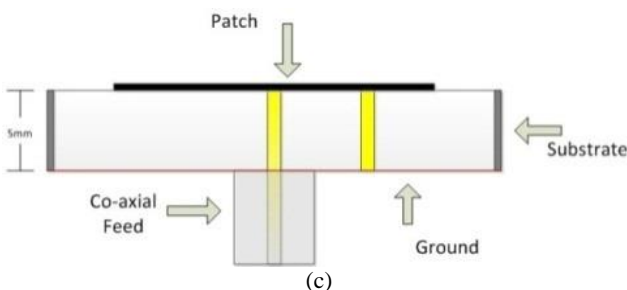


Fig. 1. (a) Front view, (b) ground view, (c) bottom view.

Fig. 1(c) shows the bottom view of our structure as evident from figure it is clear that coaxial probe feed is used as means of transmission for exciting patch antenna. The antenna is well

impedance matched with Voltage Standing Wave Ratio of 1.20 ensuring input power being delivered efficiently.

III. RESULTS AND DISCUSSION

In order to evaluate performance of our proposed design, different parameter results like realized gain, voltage standing wave ratio, directivity, reflection coefficient, efficiency, bandwidth were evaluated.

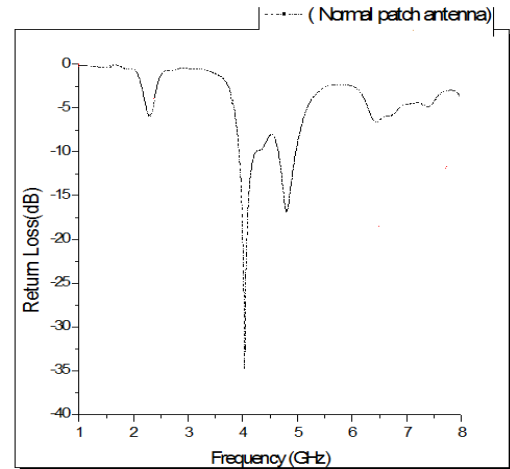


Fig. 2. Return loss of conventional antenna.

In Fig. 2, the antenna is designed for 4.1GHz and the miniaturized antenna is operating on 2.66GHz which shows in Fig. 3, due to the combination of different techniques i.e. defected patch structure, defected ground structure, shorting pin method, fractal patch structure. The frequencies are shifted downward due to the implementation of these techniques.

As our proposed design fundamental resonating frequency shifted downward to 2.66GHz while having the patch dimensions of 4.1, our proposed design showed size reduction up to 60.67% since conventional design for 2.66GHz would require dimensions of 950mm² and our design has dimensions of 375.86mm² only.

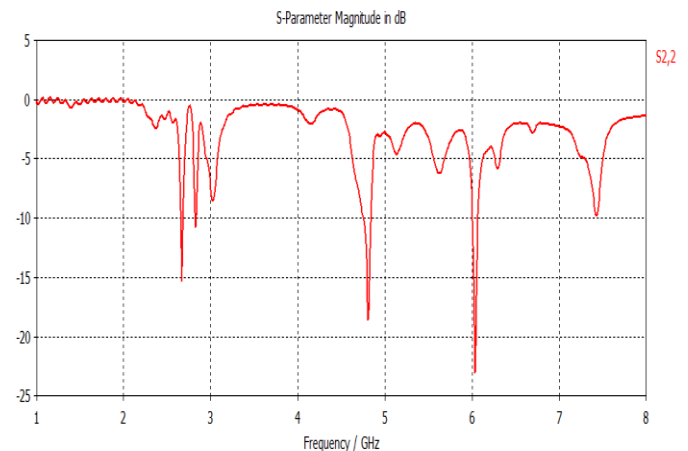


Fig. 3. Return loss of miniaturized antenna.

Reflection coefficient or Return loss graph is shown in Fig. 3. Return loss impedance bandwidth of -10dB clearly

shows six different frequencies exhibiting Hex Band Response.

All of the resonating frequency performance parameters are covered in Table I and mismatch losses are shown in Table II.

TABLE I. SIMULATION RESULTS

FREQUENCY (GHZ)	DIRECTIVITY	GAIN	BANDWIDTH	RETURN LOSS
2.66	4.25	3.43	70	-15.2
2.87	3.68	3.75	40	-20.7
3.01	4.15	3.69	30	-8.5
4.80	5.41	4.71	110	-18.7
6.03	5.41	4.71	190	-23.7
7.42	5.46	4.28	81	-10.1

TABLE II. MISMATCH LOSSES

FREQUENCY (GHZ)	2.66	2.87	3.01	4.80	6.03	7.42
VSWR	1.02	1.25	1.11	1.09	1.05	1.21

Fig. 4 mentions and shows all the radiation frequencies 1D radiation patterns. From Radiating Frequencies, it is clear that our proposed structure is resonating at six different directions which is a good response for an antenna since this behavior leads to usage of structure for different application purposes.

As in 2.66GHz the main lobe magnitude is 3.4dB, main

Lobe direction is 80.0 deg and angular width is 86.2 deg while side lobe level is -6.8dB.

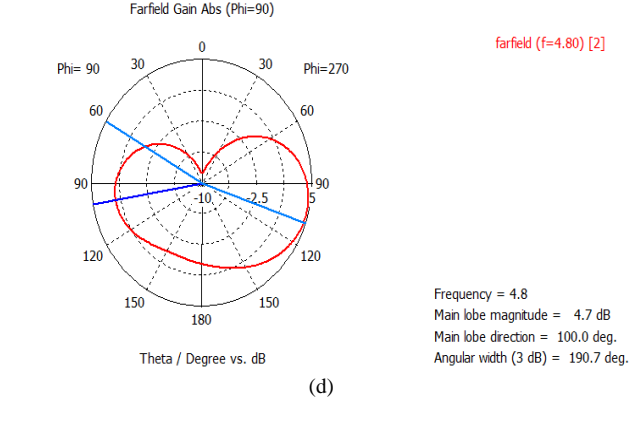
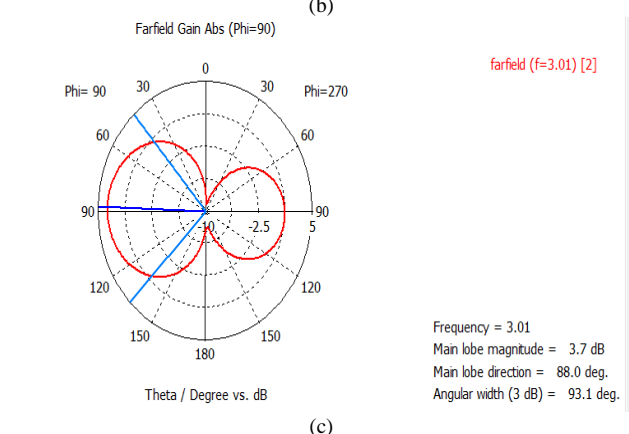
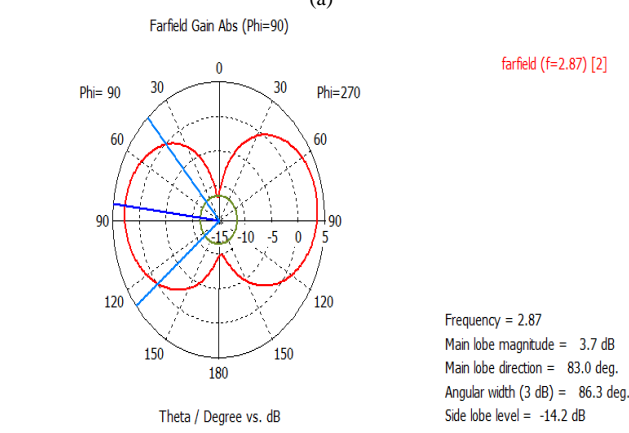
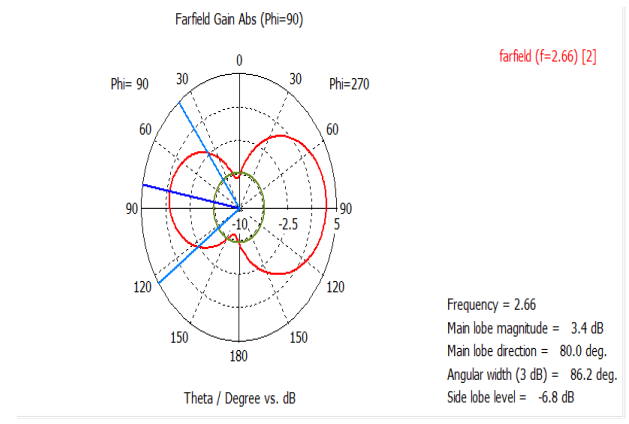
In 2.87GHz the main lobe magnitude is 3.7dB, main lobe direction is 83.0 deg and angular width is 86.3 deg while side lobe level is -14.2dB.

In 3.01GHz the main lobe magnitude is 3.7dB, main lobe direction is 88.0 deg and angular width is 93.1 deg.

In 6.03GHz the main lobe magnitude is 2.7dB, main lobe direction is 00.0 deg and angular width is 103.4 deg while side lobe level is -4.7dB

In 4.80GHz the main lobe magnitude is 4.7dB, main lobe direction is 100.0 deg and angular width is 190.7.3 deg.

In 7.42GHz the main lobe magnitude is 2.8dB, main lobe direction is 180.0 deg and angular width is 42.7 deg while side lobe level is -11.6dB.



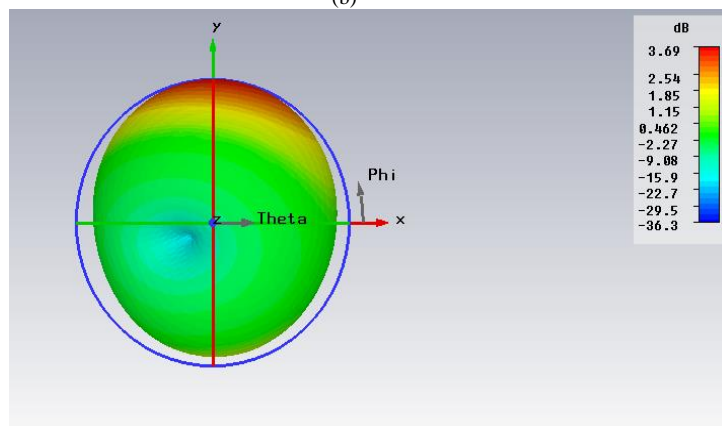
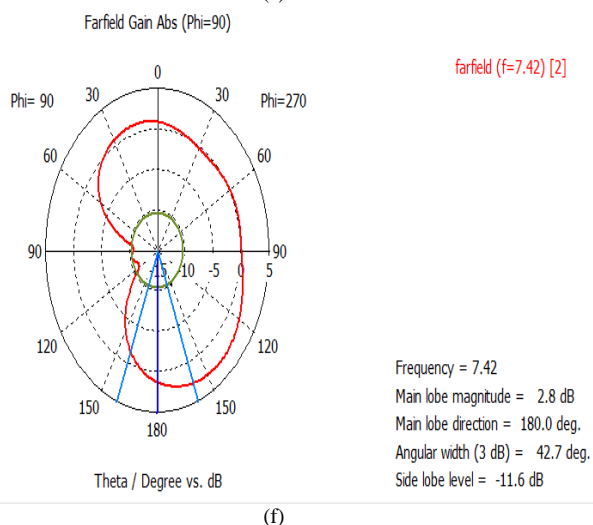
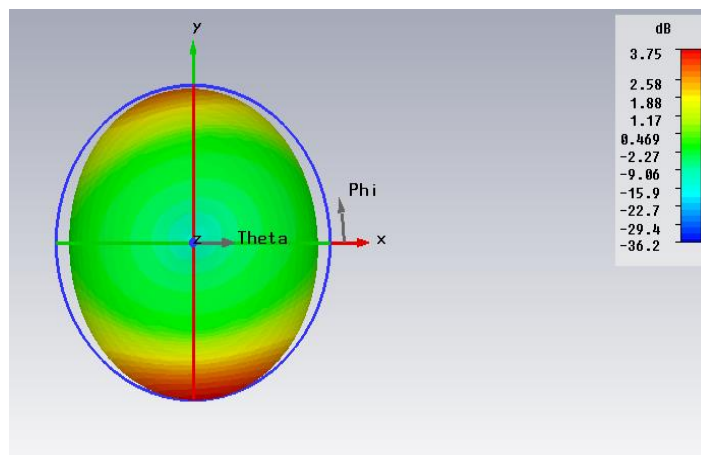
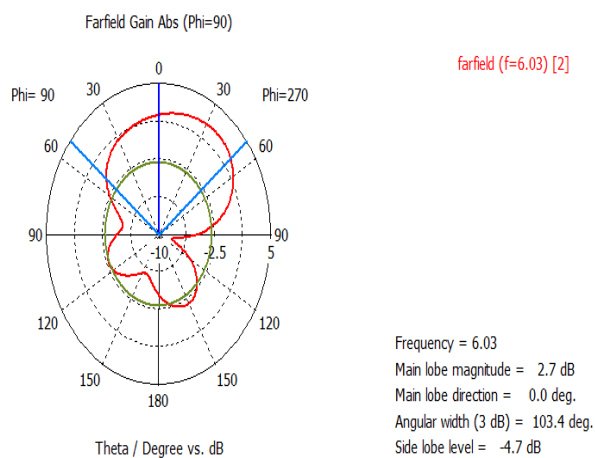
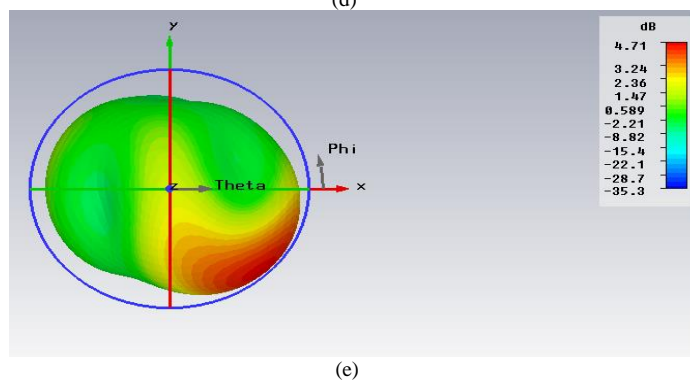
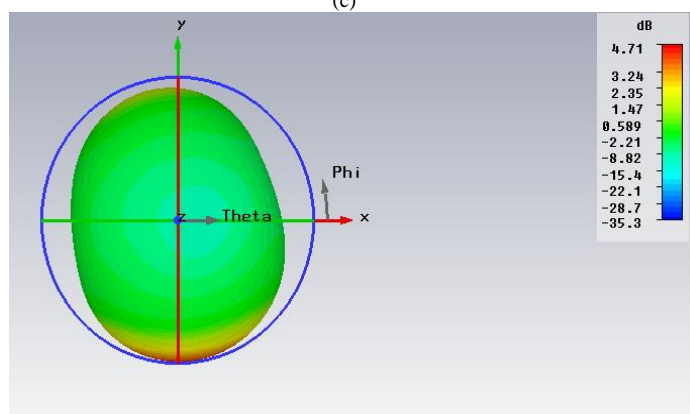
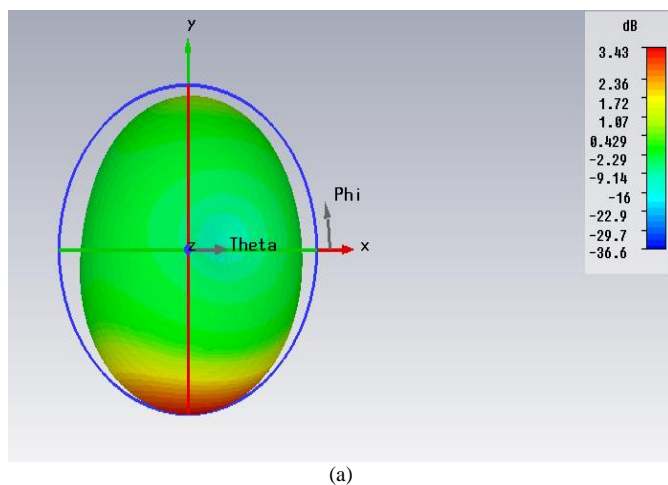


Fig. 4. 1D Radiation Pattern of the proposed miniaturized antenna for (a) at f=2.66 GHz (b) at f=2.87GHz(c) at f=3.01GHz (d) at f=4.80GHz (e) at f=6.03GHz and (f) at f=7.42GHz.

Fig. 5 shows 3D radiation pattern of miniaturized proposed antenna. All the images are extracted from Computer Simulation technology 2014.



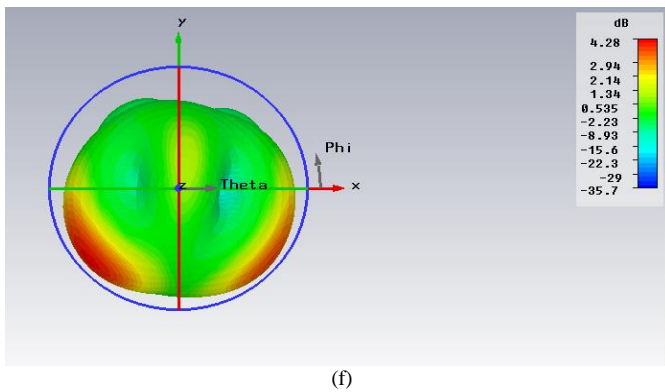


Fig. 5. (a) 3D Gain of 2.66GHz, (b) 3D Gain of 2.87GHz, (c) 3D Gain of 3.01GHz, (d) 3D Gain of 4.80GHz, (e) 3D Gain of 6.03GHz, (f) 3D Gain of 7.42GHz,

IV. CONCLUSION

In this paper, a novel fractal E slot mounted patch is proposed with H and L slots in ground. With Slots in ground miniaturization response is observed and due to fractal slotting, the antenna exhibit Hex Band Response with reduction up to 60%. Antenna performance parameters showed excellent results as gain varied form 3.79dB to 4.7dB with bandwidth up to 180MHz and minimum mismatch losses. The proposed antenna can be used for different S and C band application Systems.

REFERENCES

[1] S. Hassan, K. Mahmood and M. Munir, "U Patch Antenna using Variable Substrates for Wireless Communication Systems", *International Journal of Advanced Computer Science and Applications*, vol. 7, no. 12, 2016

[2] M. MUNIR, S. S. QURESHI, S. H. KIANI, K. MAHMOOD, J. KHAN, "Performance Analysis between Single and Dual Substrate Patches for

Wireless Communication and Applications", *Sindh University Research Journal*, vol.49, no.1, 2017

[3] Kiani, S. H., Mahmood, K., Shafeeq, S., Munir, M., & Khan, K. M. (2016). A Novel Design of Miniaturized Patch Antenna Using Different Substrates for S-Band and C-Band Applications. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 7(7).

[4] Hwang, S., Lee, B., Kim, D. H., & Park, J. Y. (2018). Design of S-Band Phased Array Antenna with High Isolation Using Broadside Coupled Split Ring Resonator. *Journal of Electromagnetic Engineering And Science*, 18(2), 108-116.

[5] S. Hassan, K. Mahmood, M. Munir and A. James, "A Novel Design of Patch Antenna using U-Slot and Defected Ground Structure", *International Journal of Advanced Computer Science and Applications*, vol. 8, no. 3, 2017.

[6] Mehre-e Munir; Ahsan Altaf; Muhammad Hasnain," Miniaturization of microstrip fractal H-Shape patch antenna using stack configuration for wireless applications", in *IEEE 2nd International Conference on Recent Trends in Information Systems (ReTIS)*, Kolkata, 2015

[7] Mehre-e-Munir; Ahsan Altaf; Syed Imran Hussain Shah, "Miniaturization of multiband patch antenna using stack configuration", in *17th IEEE International Multi Topic Conference*, Karachi, 2014

[8] Mehr-e-Munir; Khalid Mahmood," Miniaturized microstrip patch antenna using stack configuration for S-band, C-band & mobile applications", in *International Conference on Emerging Technologies (ICET)*, Peshawar, 2015

[9] Mehr-e-Munir; Umar Farooq "Multiband microstrip patch antenna using DGS for L-Band, S-Band, C-Band & mobile applications", in *13th International Conference on Modern Problems of Radio Engineering, Telecommunications and Computer Science (TCSET)*, Ukraine, 2016

[10] Li, Yujian, and Kwai-Man Luk. "60-GHz substrate integrated waveguide fed cavity-backed aperture-coupled microstrip patch antenna arrays." *IEEE Transactions on Antennas and Propagation* 63, no. 3 (2015): 1075-1085.

[11] D. Sievenpiper, H. P. Hsu, J. Schaffner, G. Tansonan, R. Garcia, and S. Ontiveros, "Low profile, four sector diversity antenna on high impedance ground plane," *Electron. Lett.* vol. 36, pp. 1343-1345, 2000.

[12] Zhang, X. Yang, "Study of a slit cut on a microstrip antenna and its applications," *Microwave and Optical Technology Letters*, vol.18, no.4, pp.297- 300, 1998.

An Optimized Inset Feed Circular Cross Strip Antenna Design for C-Band Satellite Links

Faisal Ahmed Dahri¹

Postgraduate student in Department of Telecommunication Engineering, Mehran University of Engineering & Technology, Jamshoro, Pakistan

Riaz A. Soomro², Sajjad Ali Memon³

Assistant Professor in Department of Telecommunication Engineering, Mehran University of Engineering & Technology, Jamshoro, Pakistan

Zeeshan Memon⁴

Lecture in Department of Telecommunication Engineering, Dawood University of Engineering & Technology, Karachi, Pakistan

Majid Hussain Memon⁵

Assistant Professor in Department of Electronics Engineering, Quaid-e-Awam University of Engineering, Science and Technology, Nawabshah, Pakistan

Abstract—This proposed antenna model and progressing the investigation of an inset fed wideband circular slotted patch antenna is suitable for 5.2 GHz satellite C-band applications. A circularly shaped slot has been chosen to be etched on diminutive square patch (4.4cm*5.64cm) of the inset feed antenna. The object of this work is to develop an efficient and inexpensive transducer system to facilitate its compatibility with monolithic microwave integrated circuits; expenses are minimized for its fabrication and trail low profile for C-band satellite links. This paper focuses on the circular profile of the microstrip patch antenna intended for the proficient gain to enhance the performance of the satellite communication. The return loss of -21.79dB with the directivity 8.22dB and gain of 8.17dB have been estimated. The efficiency of 97% with VSWR of 1.22 compensates each other with better simulation results.

Keywords—Circular slot; high gain; C-band; satellite communication; efficiency

I. INTRODUCTION

Rapid progress in modern satellite communication has urged meticulous research concerning high directive single and multiple feed planar antennas. The microstrip patch antenna is selected based on characteristics such as their miniature size, low weight, lower profile, ease to operate simple configuration and fabrication. Many satellite antennas are intended to offer versatility and scope pursuing antenna gain and direction; to face hurdles and hindrances of seasonal variations and enhance the performance of systems in unidirectional and multidirectional scenarios [1]. With the swift progress of worldwide technologies, the worlds of miniature size microwave structures along with low-cost designs are globally appreciated. Therefore, heaps of research work stimulated the expansion of ultra-wideband antennas for diverse areas [2]-[4]; planar circle antennas are utilized for single-band and ultra-wideband procedure which is fed by coplanar lines to prevent from persistent ground current effects. It comprises the slotted circle as the radiating part, the slotted circle to slim the antenna dimension and ameliorate the overall impedance transmission capacity and better impedance matching [5]. Microstrip patch antenna comprises minute

directing patch based on the ground plane divided by a dielectric substrate. The radiating patch has been designed with dependable materials, e.g. gold or copper and their ability to adopt any conceivable geometrical patterns [6]. It has been observed on patch antenna for wireless communication, coplanar patch antenna has been used due to its high radiation efficiency and wide bandwidth that suits with satellite transponders [7]; a wideband smaller scale strip patch antenna is intended for satellite correspondence. It accomplished multi-band usefulness throughout the cluster system together with creating circular polarization and greater radiation efficiency for C-band [8]. The unique patch shapes such as E-shaped antenna by [9], [10] found its aptness with wider bandwidth spectrum, multi-band behavior and Wi-Fi applications to minimize the antenna size with glass cum fiber substrate and antenna width appeared to be reduced in the magnitude from 19.2 mm to 15.8 mm.

The transmission capacity of the antenna has been expanded from 8.5% to 17.5% with a coverage frequency range of (5.1-6.0) GHz separately. In satellite communications, a microwave frequency of C-band is defined as the standard range from (4.0-8.0) GHz, approximately the entire C-band communication satellites utilize the band of frequencies from (3.7 - 4.2) GHz for their downlinks whereas uplink takes up the spectrum of frequencies from (5.9-6.4) GHz. The research theme, demands the reduction of antenna size effectively without affecting its radiation pattern and directivity. Therefore the proposed work covers the development of reasonable antenna design that can capture the satellite traffic of C-band by sustaining the efficiency and directivity of the system (Fig. 1).

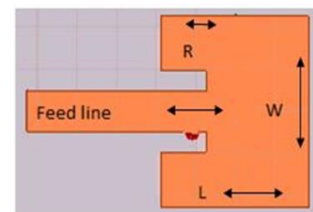


Fig. 1. Basic configuration profile of inset feed patch antenna.

II. RELATED STUDIES

To acclimatize with the increasing demands of wireless communication and to compensate the extended application prerequisites, it turns into a very strong issue to configure and design low profile, portable, wider bandwidth planar antenna at minimum cost. Many design goals have been achieved by the researchers in few decades. Circular radiation polarization antennas have been discovered with wider applications in satellite communication because of their inattentiveness to the polarization rotation of ionosphere. Helical antennas with axial mode are known for a long time offer circular polarization over a wide transmission capacity and eliminate the need for polarizer [11], [12]. The authors have investigated the performance of microstrip patch antenna for C-band satellite communication. The concept represented in this paper was based on resonant frequency to achieve multiband operations. The tradeoff was highlighted that if the width of patch increases the radiation efficiency, bandwidth and power radiation decreases [13]. In [14], F-shaped antenna is proposed for satellite communication. Simulation results were not quite suitable to accommodate the c-band traffic. The circular slotted antenna was recommended for wireless communication. The achieved gain was 4.3dB over the 3 GHz to 8 GHz frequency band. It is demonstrated that the lower the voltage standing ratio for better the transmission capacity with more power deliverance [15]. The effects of ground and patch length were optimized at different lengths and to get optimum value for better impedance matching. This work has not optimized the gain and return loss improvements. In [16], authors proposed a triple equalitarian triangular slot antenna for Ku-band satellite communication. Although results were shown good but gain and bandwidth efficiency is low. Divesh Mittal et al. [17] has discussed the performance of microstrip antenna for c-band communication. The authors failed to provide high gain and directivity. The satellite antennas ought to be directive with better gain results. The step structured slots added on the patch was proposed to provide broadband and smaller size antenna [18]. The circular slot microstrip patch antenna reported in many papers [19] but the circular slot with plus sign was not introduced in the presented literature. The multiband antenna was proposed for different wireless applications. The simulation results were better to accommodate WiMax traffic [20]. Microstrip patch antennas have many applications as E-shaped antenna proposed for intelligent transportation traffic data transfer and the Yagi shaped microstrip was designed for telemedicine applications [21], [22]. This paper [23] presents the arc-shaped strips along dual inverted L-shape partial ground plane antenna for X-band and WLAN applications. Sharma, et al. [24] demonstrated a tri-band nature of antenna used for C, X, and Ku Satellite communication and recommends that cutting the corner of the antenna with defective ground offer a solution for patch antenna drawbacks. The results satisfy the bandwidth need.

In this paper, we have proposed an antenna for satellite communication with improved gain, directivity and the effects of tuning the slot positions for getting optimum position for better return loss and VSWR.

III. ANTENNA DESIGN

There are numerous methods to design and analyze a microstrip patch antenna such as transmission line pattern hollow designs. The selected antenna utilizes transmittal line design where the width W and length of patch L are analyzed as broadening the transmittal line to resonate through an electric field varying sinusoidal as well as its length L [25]. There are three factors for the configuration of a rectangular small-scale patch antenna, Firstly the resonance band of the antenna besides the high polarization material with steady Rogers 5880 (tm) substrate; the substrate thickness (h) controls the fringing effects of the structure. Fig. 2 demonstrates a comprehensible planar rectangular Small scale strip patch antenna that operates at a frequency of 5.2 GHz with measurements of length and width are given in Table I. Rogers 5880 (tm) substrate bears the relative permittivity of 2.2 with loss tangent 0.001; selection has been made due to its ease, simple accessibility, the simplicity of fabrication. The radiating patch is aligned with an inset feed line crossing the center of the patch along the line of symmetry. The winding of the patch is adjusted by cutting one spherical slot and a pair of rectangular slots are combined together to minimize the impedance mismatching losses for satellite C-band applications, one of which is imprinted into the non-radiating center of the patch and two are united together. Pair of Slots is symmetrical in nature due to the fact that it reduces the cross polarization. The slits initiate boosting the current streak and reduce frequency operation.

A. Design Calculation of Inset Feed Patch Antenna

The inset fed antenna has been designed to resonate at 5.2 GHz. The Patch antenna length and width has been calculated by (1) and (2); whereas λ_0 & μ_0 are the wavelength and permeability of free space respectively with f_r defines the resonance frequency of the antenna. Using following formulae (1)-(8) essential parameters for design has been calculated.

$$W = \frac{\lambda_0}{2} \left(\frac{2}{\epsilon_r + 1} \right)^{\frac{1}{2}} \quad (1)$$

$$L = \frac{1}{2f_r \sqrt{\epsilon_{eff}} \sqrt{\mu_0 \epsilon_0}} - 2\Delta L \quad (2)$$

The extension length ΔL h is the height of patch and ϵ_{reff} effective permittivity are given as [12].

$$\Delta L = 0.412h \left[\left(\frac{\epsilon_{reff} + 0.3}{\epsilon_{reff} - 0.258} \right) \left(\frac{W/h + 0.264}{W/h + 0.813} \right) \right] \quad (3)$$

Where,

$$\epsilon_{eff} = \frac{\epsilon_r + 1}{2} + \frac{\epsilon_r - 1}{2} \left[\left(1 + \frac{12W}{h} \right) \right]^{-\frac{1}{2}} \quad (4)$$

According to the dimension of the width (W) of the patch, the dimension for the inset length of the patch is to be calculated. Conductance G_1 is given by [26].

If $W \gg \lambda_0$

$$G_1 = \frac{1}{120} \left(\frac{W}{\lambda_0} \right) \quad (5)$$

If $W \ll \lambda_0$

$$G_1 = \frac{1}{90} \left(\frac{W}{\lambda_0} \right)^2 \quad (6)$$

Now, take into account that the characteristic impedance of the microstrip line feeder is R_{in} . Therefore, equating the

given equations [26] to get similarities between the input impedance R_{in} of the patch and feeder (i.e. inset length, Y_0).

$$R_{in}(y = 0) = \frac{1}{2(G_1 \pm G_2)} \quad (7)$$

And then

$$R_{in}(y = y_0) = \frac{1}{2(G_1 \pm G_2)} \cos^2\left(\frac{\pi}{L} y_0\right) \quad (8)$$

B. Simulation of Inset Feed Patch Antenna with a circular slot

It has been simulated small-scale strip inset fed patch antenna by means of defined parameters of an inset fed patch antenna and employing those parameters for the simulation work. The physical measurements of proposed antenna are specified in Tables I and II separately. Fig. 2 demonstrates the configuration of a microstrip patch antenna structure in the HFSS simulation software with activation of electric fields.

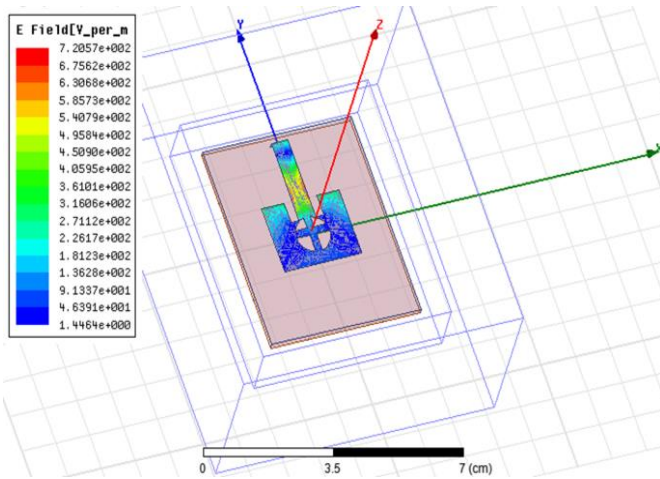


Fig. 2. Circular cross slot patch antenna with the emission of E-field distribution.

TABLE I. SPECIFICATION LIST OF INSET FED PATCH ANTENNA

Parameters	Units
Operating frequency	5.2 GHz
Ground plane dimension (L*W)	4.4cm*5.64cm
Patch length	1.85 cm
Patch Width Value	2.28 cm
Dielectric Constant Value	2.2
Dielectrical Material	RogersRT/duroid 5880(tm)
Substrate Thickness height	0.15784 cm
Input impedance	50Ω
Effective permittivity	1.23
The wavelength of free space	0.05769 cm

TABLE II. SPECIFICATIONS FOR LENGTH OF INSET FED PATCH ANTENNA

Parameters	units
Inset Distance	0.566cm
Inset Gap	0.243cm
Feed width	0.485cm
Feed length	2.257cm
Slot position	0,0,0.4

IV. RESULTS AND DISCUSSION

This segment furnishes the simulation results. The return Loss S_{11} of the circular dimension antenna is shown in Fig. 3. The transmission capacity is accomplished from 4.9 GHz to 5.3 GHz of 8.7%. The improved transmission capacity has resulted in the effects of the operating frequency. The efficiency of the antenna is 97%. The Directivity accomplished is 8.22 dB and the gain of the antenna is 8.04 dB at 5.2 GHz. The distinctive simulation results incorporate 2D and 3D radiations as appeared in Fig. 4 to 7.

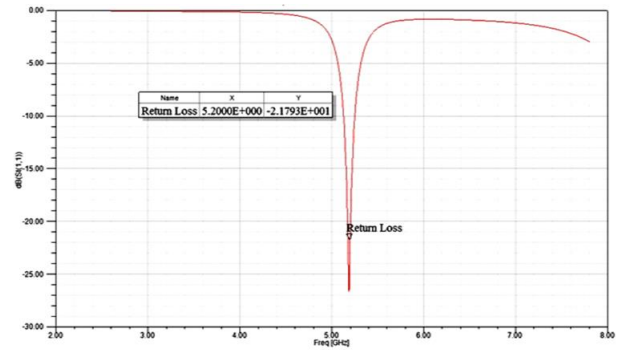


Fig. 3. Simulated Return Power Loss of circular slotted patch antenna with 5.2 GHz resonance frequency.

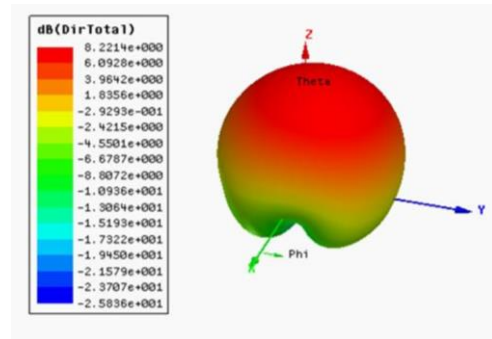


Fig. 4. 3D radiation pattern of the circular slotted inset fed patch antenna with directive value at the resonance frequency of 5.2GHz.

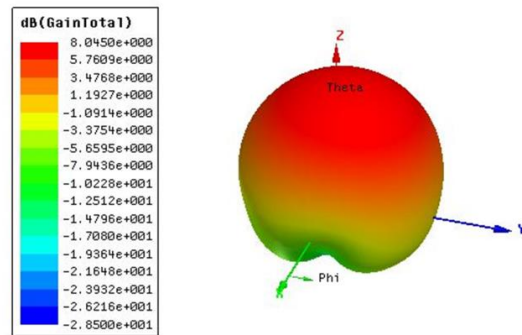


Fig. 5. The 3D radiation pattern of the circular slotted inset fed patch antenna with Gain values at the resonance frequency of 5.2GHz.

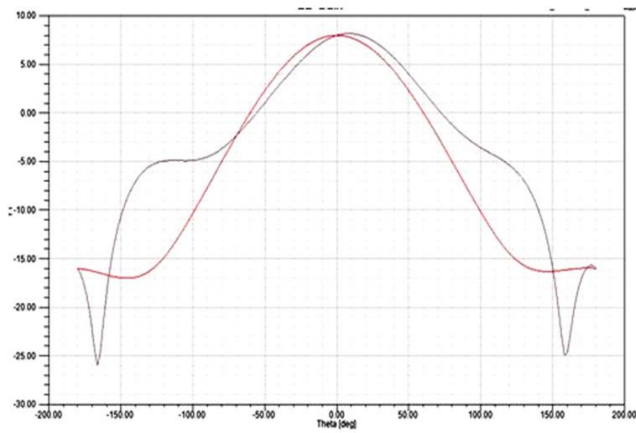


Fig. 6. Simulated 2D radiation pattern and with a Co-polarization pattern of the circular slotted Inset fed patch antenna.

TABLE III. TUNED CIRCULAR SLOT SHAPED ANTENNA

Radius	0.4	0.4	0.5
Slot Position	0,0,0	0,-0.5,0	0,-0.2,0

The impact of the different measurements such as slots width and position of the circular slotted antenna with simulation results has examined and discussed. The expansion in width of the slot would tune the return loss S_{11} execution to higher operating frequency. Correspondingly, this impact is quite obvious as the antenna is resonating at a higher frequency given that current flows through the center while the lower frequency exists owing to the movement of current through the circular region; another important circular slotted patch antenna design parameter is the slot location, which can alter the results for different return losses. The effects are shown in Fig. 8 and 9. The radius of the slot escalates from (0.4-0.5) cm indicated in Table III suggest slot positions, it shows results at the same frequency.

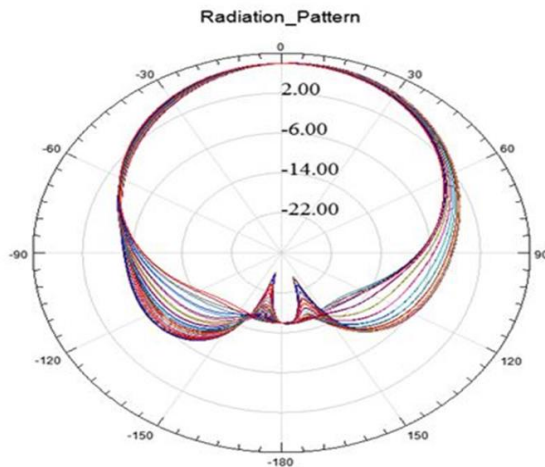


Fig. 7. Polarization Simulated radiation pattern at 5.2GHz for circular slotted inset fed antenna.

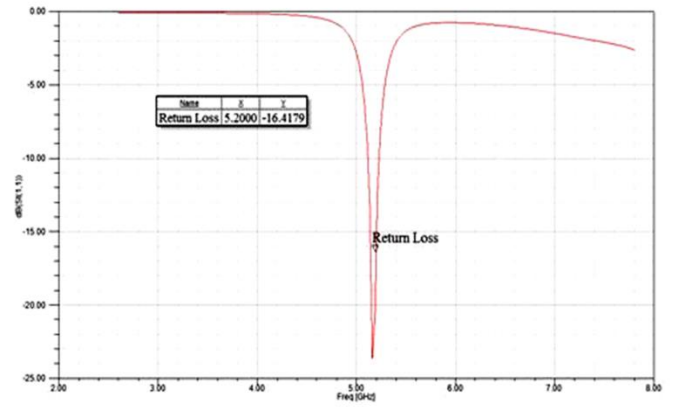


Fig. 8. Simulated return power loss of circular slotted patch antenna with the resonant frequency of 5.2 GHz.

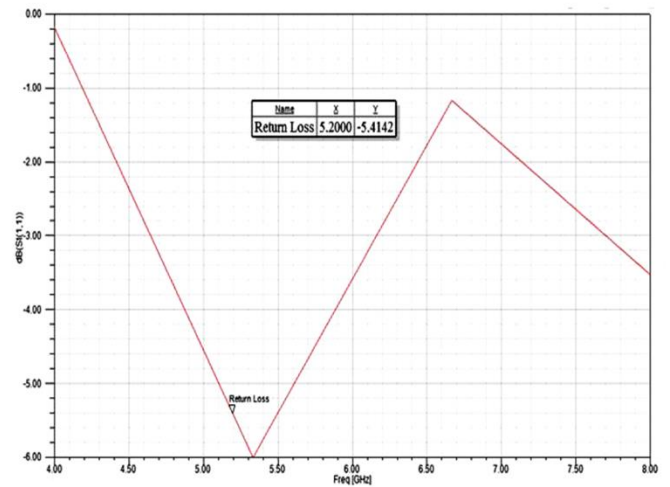


Fig. 9. Simulated return power loss of circular slotted patch antenna with 5.2 GHz.

TABLE IV. COMPARATIVE SUMMARY OF SATELLITE COMMUNICATION ANTENNAS

Ref	Frequency [GHz]	Return loss [dB]	Gain [dB]	VSWR
[13]	5.2	-15	4.4	1.5
[14]	4.1	-13.11	-	1.56
[15]	4.4	-32	1.8	1.2
[18]	5	-19.74	4	1.25
[20]	6.2	-15	-	<2
This work	5.2	-21.79	8.04	1.22

V. CONCLUSION

In the proposed work, a substantial structure of circular slot rectangular microstrip patch antenna has been designed successfully that manifests the efficiency of 97% with operating frequency of 5.2 GHz. This exclusive design and profile of the antenna mostly serve the role of improved results for the satellite C-band application as shown in Table IV. The simulated results are good which validate proposed antenna for C-band traffic. Adequate outcomes have been found including directivity of 8.22dB with the high gain of 8.04 dB; VSWR of 1.22 compensates with a return loss of -21.79dB to facilitate the efficient transmission of

electromagnetic energy from earth station to transponder antenna and then headed for receiving station which confirms the suitability of circular cross-sectional patch antenna for C-band satellite links. In future, According to obtained simulated results of proposed design will be fabricated and simulated in other simulators for the performance test. The intended method may be extended and modified for other type satellite bands antennas.

REFERENCES

- [1] Jorge Sosa-Pedroza, Fabiola Martinez-zuniga and Mauro Enciso-Aguilar (2010). "Planar Antennas for Satellite Communications, Satellite Communications", Nazzareno Diodato (Ed.), InTech, Instituto Politecnico Nacional. Escuela Superior de Ingenieria Mecanica y Electrica Mexico.
- [2] Lin, C.-C. and H.-R. C.Lin, "A 3-12 GHz UWB planar triangular monopole antenna with ridged ground-plane", Progress In Electromagnetics Research, Vol.83, 2008, pp.307-321
- [3] Li, X., L. Yang, S.-X. Gong, and Y.-J. Yang, "Ultra-wideband monopole antenna with four-band-notched characteristics", Progress in Electromagnetics Research Letters, Vol. 6, 2009, pp.27-34,
- [4] Khan, S. N., J.Hu, J. Xiong, and S.He, Circular fractal monopole antenna for low VSWR UWB applications, Progress In Electromagnetics Research Letters, Vol. 1,2008, pp. 19-25.
- [5] M.K. Mandal, Z. N. Chen, "Compact Dual-Band and Ultra Wideband loop antennas", IEEE Transactions on Antennas and Propagation, vol. 59, 2011, pp. 2774-2779.
- [6] Williamsburg Virginia, Microstrip antennas, IEEE international symposium on Antennas and Propagation, 1972 pp 177-180.
- [7] Arvind Singhjadan, Jalaj Sharma, Ajay Prajapat & Avanish Bhadauria, "Coplanar Rectangular patch antenna for X Band Applications using inset fed technique", International Journal of Electronics and Communication Engineering and Technology (IJECET) November 2013, Volume; 4 Issue; 7, page: 95-102.
- [8] M.R.Ahsan, M.T. Islam, M. Habib Ullah, W.N.L. Mahadi, & T.A.Latef , "Compact Double-P Slotted Inset-Fed Microstrip Patch Antenna on High Dielectric Substrate" , Scientific World Journal Volume 2014 (2014), Article ID 909854, 6 pages.
- [9] A. A. Razzaqi, M. Mustaqim and B. A. Khawaja, "Wideband E-shaped antenna design for WLAN applications", Emerging Technologies (ICET), 2013 IEEE 9th International Conference on, Islamabad, 2013, pp. 1-6.
doi: 10.1109/ICET.2013.6743550.
- [10] Vivek Kumar Agarwal, Anand Kumar Shaw, Mrinmoy Kr. Das, Jayati Mukherjee & Kaushik Mandal, " A Novel Compact Dual-Frequency Microstrip Antenna", Science Direct Procedia Technology 4(2012) 427 - 430.
- [11] Kraus, J.D., Helical beam antennas," Electronics, Vol. 20, 109-111, 1947.
- [12] Nakano, H., Y. Samasa, and J. Yamauchi, "Axial mode helical antennas," IEEE Transactions on Antennas and Propagation, Vol. 34, No. 9, 489-509, 1986.
- [13] Afzal, W., Y. Mehmood, M. Z. Baig, and M. Z. Babar. "Optimization of Microstrip Patch Antenna for C-band and WLAN Applications." 2nd International Conference on Engineering & Emerging Technologies (ICEET),2015.
- [14] Muhammad Afsar Uddin, Fatama Akter, Md. Jahangir Alam. "High Bandwidth F-Shaped Microstrip Patch Antenna for C-band Communications" International Journal of Advanced Research in Computer and Communication Engineering, Vol. 6, No.3, pp 143-146, 2017.
- [15] Gupta, Prerna, and S. N. Vijay. "Design and parametric study of rectangular microstrip patch antenna for C-Band Satellite Communication." In Recent Advances and Innovations in Engineering (ICRAIE), 2016 International Conference on, pp. 1-5. IEEE, 2016.
- [16] Samsuzzaman, M., Mohammad Tariqul Islam, and J. S. Mandeep. "Design of a compact new shaped microstrip patch antenna for satellite application." Advances in Natural and Applied Sciences 6, no. 6 (2012): 898-903.
- [17] Divesh Mittal, Aman Nag, Ekambir Sidhu "Design and Performance analysis of Microstrip Patch antenna for C band applications", International Journal of Engineering Trends and Technology (IJETT), V48(5),242-246 June 2017.
- [18] Gour, Puran, Ravishankar Mishra, A. S. Zadgaonkar, and Pradeep Kumar Sahu. "Circular Polarized Broad Band Microstrip Patch Antenna for C band." International Journal of Computer Technology and Applications (IJCTA), Vol. 9, No. 41, pp 191-197, 2017.
- [19] Thaker, Nidhi M., and Vivek Ramamoorthy. "A Review on Circular Microstrip Patch Antenna with Slots for C Band Applications." International Journal of Scientific & Engineering Research 5, no. 12 (2014).
- [20] Badr, Samah, and Ehab KI Hamad. "Design of multiband microstrip patch antenna for WiMax, C-band and X-band applications." Aswan Engineering Journal (AswEJ), 2018.
- [21] Asif Ali, Nasrullah Pirzada, Muhammad Moazzam Jawaid and Sajjad Ali Memon, "Design and Simulation of a Rectangular E-Shaped Microstrip Patch Antenna for RFID based Intelligent Transportation" International Journal of Advanced Computer Science and Applications(IJACSA), 9(4), 2018.
- [22] Dahri, F. A, Soomro R. A, Memon Z. A "Design of Wearable Microstrip Yagi Array Antenna aimed for Telemedicine Applications", Academic Journal of Management Science (AJMS), Vol. 5, No.2, pp 45- 52, 2017.
- [23] Zhi, R., Han, M., Bai, J., Wu, W., & Liu, G. "Miniature Multiband Antenna for WLAN and X-Band Satellite Communication Applications". Progress In Electromagnetics Research, 75, 13-18, 2018.
- [24] Sharma, I. B., Lohar, F. L., Maddila, R. K., Deshpande, A., & Sharma, M. M. "Tri-Band Microstrip Patch Antenna for C, X, and Ku Band Applications". In Optical and Wireless Technologies (pp. 567-574). 2018, Springer, Singapore.
- [25] H. Pues and A Van de Capelle, Accurate transmission-line model for the rectangular microstrip antenna Proc. IEEE, vol. 131, pt. H, no. 6, pp. 334-340.
- [26] C.A. Balanis, Antenna Theory, 2nd Ed., John Wiley and sons, inc., New York.

Hybrid Ensemble Framework for Heart Disease Detection and Prediction

Elham Nikookar

Department of Computer Engineering, Faculty of
Engineering
Shahid Chamran University of Ahvaz
Ahvaz, Iran

Ebrahim Naderi

Computer Department,
University of Applied Science and Technology,
Ahvaz, Iran

Abstract—Data mining techniques have been widely used in clinical decision support systems for detection and prediction of various diseases. As heart disease is the leading cause of death for both men and women, detection and prediction of the heart disease is one of the most important issues in medical domain and many researchers developed intelligent medical decision support systems to improve the ability of the CAD systems in diagnosing heart disease. However, there are almost no studies investigating capabilities of hybrid ensemble methods in building a detection and prediction model for heart disease. In this work, we investigate the use of hybrid ensemble model in which a more reliable ensemble than basic ensemble models is proposed and leads to better performance than other heart disease prediction models. To evaluate the performance of proposed model, a dataset containing 278 samples from SPECT heart disease database is used that after applying the model on the data, 96% of classification accuracy, 80% of sensitivity and 93% of specificity are obtained that indicates acceptable performance of the proposed hybrid ensemble model in comparison with basic ensemble model as well as other state of the art models.

Keywords—Data mining; hybrid ensemble; base classifier; classification accuracy; sensitivity; specificity

I. INTRODUCTION

The World Health Organization has estimated that 12 million deaths occur worldwide, every year due to the heart diseases [1]. Although, in the last few decades many computational tools have been designed to improve the abilities of physicians for making decisions about condition of disease in their patients [2], low performance of current heart disease detection models is remained a matter of concern and potential of data mining algorithms which are motivated by the need of an expert system, have not be highlighted in any research yet.

Artificial intelligence techniques as a subfield of data mining have been increasingly used in solving problems in medical domains such as in oncology, urology, liver pathology, cardiology, gynecology, thyroid disorders and perinatology [2]. The primary concern of artificial intelligence in medicine is construction of an intelligent system that can assist a medical doctor in performing expert diagnosis as well as predicting probability of a disease in a patient more accurately. Besides, artificial intelligence algorithms have great potential for exploring the hidden patterns in the datasets of the various disease related subjects by adjusting the data mining model for utilizing such patterns for clinical diagnosis

[1] and this potential has led to building expert systems that can be used in CAD systems for prediction and detection of diseases in patients. One of the concepts that have been emerged in recent years is the idea of combining classifiers as a new direction for the improvement of the performance of individual classifiers [3]. These classifiers could be based on a variety of classification methodologies and could achieve different rate of correctly classified samples. Such classifiers which are called ensemble classifiers have potential to lead to an increase in generalization performance by combining several base or weak classifiers and train them on the same task [4]. However, although in recent years, better models of ensemble classifiers such as hybrid ensemble classifier which have been proved to achieve better performance than basic ensemble algorithms has been introduced [5], there are almost no studies investigating application of hybrid ensemble models and their feasibilities in heart disease domain. Thus, in this study, we evaluate the performance of a hybrid ensemble model which uses five popular classification methods including Naïve Bayes, k-NN, Random Tree, SVM and Bayes Net as base classifiers and takes benefits of aggregating all these classifiers by forwarding their results to a novel fuser classifier which is chosen in this study between Adaboost, LogitBoost, MLP and Random Forest for the diagnosis of the heart disease disorders. To evaluate the performance of the proposed model, a comparative study is realized by using a dataset containing 267 samples which is available in public UCI Repertory website [6]. We finally show that the proposed method is capable of being used as a more powerful tool to assist the medical doctor in detection and prediction of the heart disease than the basic ensemble models as well as other state of the art models.

This paper is organized as follows. Section II presents the dataset that is used to train, test and evaluate the proposed model. In Section III a number of previous studies in heart disease detection and prediction domain is discussed which culminates with an identification of the knowledge gap and inconsistencies in the literature. Section IV explicitly explains the proposed model and Section V provides the performance evaluation measures used in this study. In Section VI single base classifier model which is investigated to be compared with proposed model is introduced and in Section VII experimental results are provided. Section VIII presents a general discussion of the study. Section IX concludes the study and Section X provides the recommendations for future studies.

II. DATASET

SPECT heart disease dataset is used in this paper which is available on university of California, Irvine (UCI) machine learning dataset repository [6]. The dataset is provided for investigating diagnose of cardiac Single Proton Emission Computed Tomography (SPECT) images using machine learning algorithms. SPECT, or less commonly, SPET, is a nuclear medicine tomographic imaging technique using gamma rays. It is very similar to conventional nuclear medicine planar imaging using a gamma camera (that is, scintigraphy). However, it is able to provide true 3D information. This information is typically presented as cross-sectional slices through the patient, but can be freely reformatted or manipulated as required.

SPECT heart disease dataset was obtained from Medical College of Ohio, OH, U.S.A. The database of 267 SPECT image sets (patients) was processed to extract features that summarize the original SPECT images. As a result, 22 continuous feature patterns were created for each patient. All continuous attributes have integer values from the 0 to 100 but were further processes to obtain 22 binary feature patterns. Each of the patients is classified into one of two categories: normal and abnormal. SPECT dataset was firstly utilized by Kurgan et al. [7] where they used CLIP3 algorithm which used to generate classification rules from Features. The performance evaluation of their proposed model was evaluated by classification accuracy and the maximum value that they achieved was 90.4% in their study.

III. BACKGROUND

Classification algorithms are generally very useful for medicinal issues, especially when applied for the heart disease detection and prediction purposes [8]-[16]. Many machine learning algorithms are applied in the medical domain in the course of recent decades. A large portion of these applications are specific and include machine learning procedures like using data mining for identification and detection of disease in patients [7] and application of neural network rules for the prediction of breast cancer [17]. For example, in [18] an intelligent model is proposed for the detection of heart disease based on wavelet packet neural networks (WPNN) and they reported 94% of correct classification rate for abnormal and normal subjects. In [11] a system is proposed for diagnosis and prediction of heart disease based on Genetic Neural Network Using Risk Factors. In [9] the use of least-square support vector machines (LS-SVM) classifier for improving the performance of the proposed model of [13] is investigated. However, according to what previous studies reported, they did not investigate the use of hybrid ensemble methods to predict the occurrence of heart disease based on SPECT images of patients. Lack of research studies on this topic makes it unclear whether the hybrid ensemble models are capable of providing a model that utilizes the power of ensemble model by merging initial features of patients and predicted class labels by base classifiers. Therefore, the present study is focused on the idea of hybrid ensemble models and investigates the effectiveness of such models on the performance of a heart disease detection and prediction system.

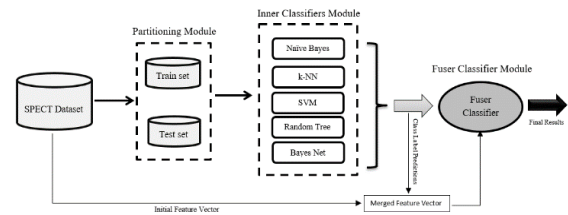


Fig. 1. The general flowchart of proposed hybrid ensemble model.

IV. METHOD

The aim of this paper is to propose a hybrid ensemble model for heart disease detection and prediction which focuses on predicting labels of each SPECT image based on feature vector of the images and the labels that base classifiers assign to each image. To facilitate understanding of the proposed framework, in this section we describe the details of layout of the proposed model. A schematic illustration of proposed hybrid ensemble model can be seen in Fig. 1. It consists of three modules, including *partitioning module*, *inner classifiers module* and *fuser module*. The initial dataset is first given to *partitioning module* to produce train and test subsets and prepare them for the next module. In *inner classifiers module* different classification algorithms are applied on the train and test datasets to produce input data for *fuser module* in which results of base classifiers next to initial feature vector of samples are considered simultaneously for building and adjusting components of the final classifier. In the rest of this section, a brief description of each component is given.

A. Partitioning Module

This module divides the initial dataset into test and train subsets by assigning 80 samples to train set and 187 samples to test set and provides mutually exclusive datasets which share no instance with each other and provides initial data for base classifiers in the next module.

B. Inner Classifiers Module

This module is constructed using five classification algorithms as base or weak classifiers including Naïve Bayes, k-NN, Random Tree, SVM and Bayes Net. All these base classifiers are applied on the train data using 10-fold cross validation as model validation technique to be adjusted for the best possible prediction about healthy or unhealthy situation of a patient. The reason of considering odd number of classifiers in *inner classifiers module* is based on the pigeonhole principle [19], which states that for natural numbers k and m , if $n=km+1$ objects are distributed among m sets, at least one of the sets will contain at least $k+1$ objects. For arbitrary n and m , it generalizes to $k+1=\lfloor(n-1)/m\rfloor+1$, where $\lfloor \cdot \rfloor$ is the floor function. It means that in the two-class problem (healthy 0, unhealthy 1) in which each classifier has to give its vote for the class of a sample, there is a need to have an odd number of classifiers to avoid equal 0 and 1 predictions for a sample. This odd number is considered five in this study. The increasing number of classifiers may obviously result in finding a more powerful model for the data but it has the risk of overfitting the model on the specific data which is used in this study. After applying all runs of 10-fold cross validation, test dataset is given to *inner classifiers module* to assign each test sample five labels by five base classifiers. To provide

input data for fuser classifier in next module, these labels are added to feature vector of test samples which leads to generating a new feature vector for each test sample including 22 binary features from initial feature vector and 5 features from *inner classifiers module* of proposed model.

C. Fuser Module

After training and testing the five classifiers in *inner classifiers module*, a new feature vector is built with 27 features including 5 predicted class labels by five base classifiers plus 22 initial features of samples. Then, the new dataset is used to find an optimal fuser classifier for the model. The candidates for fuser classifier are Adaboost, LogitBoost, MLP and Random Forest. As the fuser classifier needs to be trained to fit the data in the best form, the test dataset produced in *partitioning module* is divided itself into test and train subsets using a stratified training-test partition (80-20) and 10-fold cross validation is used as model selection technique in *fuser module* to adjust the fuser classifier and complete the hybrid ensemble model. The final result of the model is then produced for all test samples.

V. PERFORMANCE EVALUATION MEASURES

Performance evaluation is mandatory in all automated disease recognition systems and is conducted in this study to evaluate the ability of base classifiers as well as proposed hybrid ensemble model for predicting possibility of heart disease in patients based on SPECT images. Although precision and recall are more common in general data mining tasks, in medical domain, researchers prefer to assess how much sensitive and specific their proposed model is and the standard evaluation measures are sensitivity and specificity. Actually, in clinical context, a more sensitive model is preferable as the cost of overlooking a positive sample is very high and a more specific model is preferable as the cost of registering a sample as positive for the samples that are not the target of testing is very high [20].

$$Sensitivity = tp / (tp + fn)$$

$$Specificity = tn / (tn + fp)$$

The classification accuracy is also considered as evaluation measure in this study as it facilitates comparison of the results of present study with other state of the art models. The classification accuracy, CA, depends on the number of samples correctly classified (true positives plus true negatives) and is evaluated by the formula:

$$CA = (t/n).100$$

where t is the number of sample cases correctly classified, and n is the total number of sample cases.

VI. SINGLE BASE CLASSIFIER MODEL

To compare the results of our proposed model with the situation in which only a single base classifier is used, such as only SVM is investigated, this study separately applied all the five classifiers on the dataset. General flowchart of applying single base classifier is illustrated in Fig. 2. Same dataset diving procedure like *partitioning module*, i.e. 30% for train data and 70% for test data, as well as 10-fold cross validation have been used for providing data for each single base classifier training and testing.

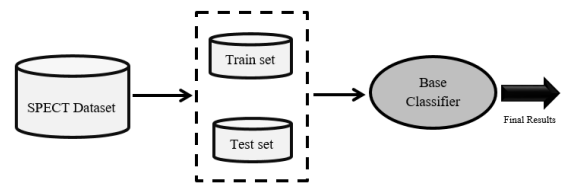


Fig. 2. The general flowchart of applying a single base classifier on the dataset.

VII. EXPERIMENTAL RESULTS

Weka, which is a collection of machine learning algorithms for data mining tasks [21] is used to train, test and evaluate the proposed model as it has two important characteristics; it is a free software system and it uses ARFF files that can be easily used and modified without data format problems. The results of applying base classifiers as well as results of experiments conducted to choose the best fuser classifier will be discussed in part 1 of this section. In part 2, results of applying different classifiers as fuser classifier will be investigated and a comparison between single base classifier model introduced in Section VI and the proposed hybrid ensemble model will be discussed in part 3. Comparing results of applying basic ensemble with results of the proposed model is conducted in part 4. In addition, a comparison between the proposed hybrid ensemble model and other heart disease detection and prediction systems will be discussed in part 5.

1) Results of Applying Single base Classifiers

As we mentioned earlier, five well-known classification algorithms including Naïve Bayes, k-NN, SVM, Random Tree and Bayes Net were used in *inner classifiers module* for constructing the core of proposed hybrid ensemble model. The experimental results of applying each base classifier on the train data is given in Table I.

As shown in Table I, the results indicate that the best base classifiers are k-NN and Random Tree considering sum value of three evaluation measures as decision criterion, However, we do not ignore predicted labels by any of base classifiers and in the next step, vote of each base classifier which is a predicted class label is kept to be used in *fuser module* for the purpose of constructing new feature vector for the hybrid ensemble model.

TABLE I. RESULTS OF APPLYING BASE CLASSIFIERS ON THE TRAIN DATA.

Method	Sensitivity	Specificity	CA
Naïve Bayes	55%	75%	69%
k-NN	63%	75%	76%
SVM	63%	48%	75%
Random Tree	68%	62%	78%
Bayes Net	55%	75%	69%

TABLE II. RESULTS OF APPLYING DIFFERENT FUSER CLASSIFIERS ON THE DATA.

Method	Sensitivity	Specificity	CA
Adaboost	71%	93%	85%
LogitBoost	85%	77%	93%
MLP	80%	93%	96%
Random Forest	77%	77%	90%

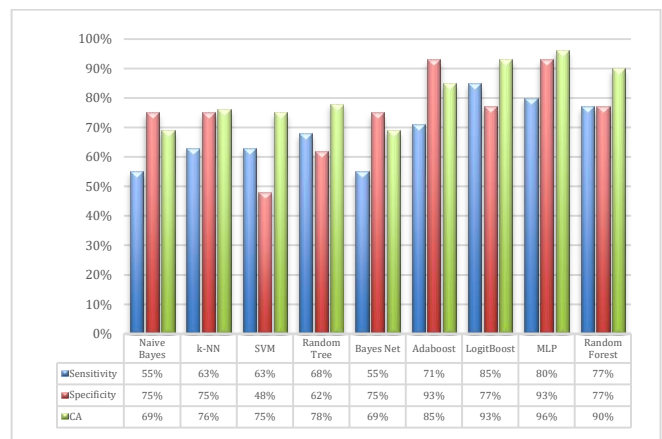
2) *Choosing best Fuser Classifier*

For selection of best fuser classifier many choices were available among diverse collection of classification algorithms. Among all choices, Adaboost, LogitBoost, MLP and Random Forest were chosen as they proved to produce acceptable results in most of the machine learning models. The experimental results of applying different fuser classifiers can be observed in Table II. It is needed to point that for each fuser classifier different configurations has been tested and the best result of each classifier is inserted in Table II.

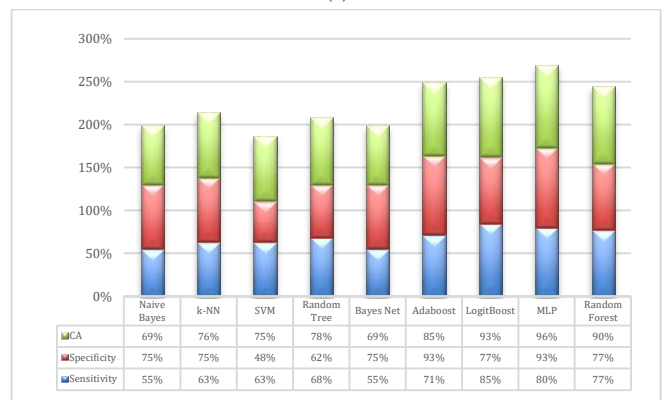
Based on the results, MLP is the best performing candidate to be chosen as fuser classifier. The parameters of the MLP that we applied on the data were set as follows: The backpropagation learning algorithm has been used in the feedforward single hidden layer of the neural network. The algorithm used for training the proposed MLP is the Levenberg–Marquardt (LM) algorithm [22]. A tangent sigmoid transfer function has been used for both the hidden layer and the output layer of the model. Besides, we used 10 neurons in the hidden layer, the initial weights were chosen randomly and in regression node, logistic regression was used.

3) *Compare Hybrid Ensemble Model with base Classifiers*

In Fig. 3, the results of applying different base classifiers as well as results after applying fuser classifier methods in *fuser module* can be observed. It is clear that the hybrid ensemble model enhances results of base classifiers on the data and there is a considerable difference between results of best performing base classifiers which are k-NN and Random Tree and the results of best fuser classifier, i.e. MLP. Therefore, the proposed hybrid ensemble model has its strength from both powerful base classifiers in *inner classifiers module* and fuser classifier which incorporates initial features of samples with predictions of base classifiers for samples. In fact, the idea of applying fuser classifiers for building up an effective ensemble classifier, in line with the idea of adding predictions of base classifiers to initial feature vector of samples, leads to final results of the proposed model. From the experimental results that are given in Fig. 3 in two charts, we conclude that the proposed hybrid ensemble model outperforms all base classifiers in terms of sensitivity, specificity and classification accuracy.



(a)



(b)

Fig. 3. Comparison between results of fuser classifiers with single base classifiers. Chart (a) shows each evaluation measure with separate bar and chart (b) shows all of the measures with one bar.

4) *Compare Hybrid Ensemble Model with basic Ensemble Model*

A basic ensemble in this study means an ensemble similar to proposed hybrid ensemble model with the difference that in basic ensemble model, feature vector of the samples which are fed to final module of the model only includes class labels predicted by base classifiers in *inner classifiers module* and does not include initial features of samples. The results of applying basic ensembles on the test data is shown in Table III. The results indicate that the idea of merging features with predicted class labels led to construction a model with better performance.

TABLE III. COMPARISON BETWEEN RESULTS OF APPLYING HYBRID ENSEMBLE (HE) AND BASIC ENSEMBLE (BE) ON THE TEST DATA

Method	Sensitivity		Specificity		CA	
	HE	BE	HE	BE	HE	BE
Adaboost	71%	53%	93%	80%	85%	69%
LogitBoost	85%	62%	77%	70%	93%	60%
MLP	80%	68%	93%	74%	96%	68%
Random Forest	77%	61%	77%	70%	90%	72%

TABLE IV. COMPARISON BETWEEN PROPOSED MODEL AND OTHER HEART DISEASE DETECTION AND PREDICTION MODELS.

Author	Method	CA
[23]	Naïve Bayes	81%
[10]	Fuzzy-AIRS-Knn based System	87%
[13]	Neural Network Ensemble	89%
[7]	CLIP3	90%
[18]	wavelet packet neural networks (WPNN)	94%
Proposed Model	Hybrid Ensemble Model	96%

5) Comparison with other Heart Disease Detection Methods

Although the experiment has achieved acceptable results by building a hybrid ensemble model, another important challenge is to compare current study with other previous methods. Related studies reporting same evaluation measures to the present study has been searched. The majority of the previous studies applied their models on private datasets and reported the results in different forms as there is no standard for this process. With all this among similar studies, as shown in Table IV, the proposed approach has provided better performance than the other techniques regarding to the classification accuracy which is the general performance measure that is used in all related studies.

VIII. DISCUSSION

The ability of an artificial intelligence model in predicting the possibility of heart disease is imperative for decreasing the mortality rate of heart disease. The ability in this study is expressed in terms of evaluation measures including sensitivity, specificity and classification accuracy that in our best configuration, the experimental results respectively show the values of 80%, 93% and 96% for these evaluation measures. This study highlights two important aspects. First, the effectiveness of using an ensemble classifier instead of base classifiers may be obvious. Second and the more important, the effectiveness of considering a combination of initial features of samples and class labels of samples predicted by base classifiers as the feature vector of fuser classifier instead of only considering predicted class labels by base classifiers which is common in basic ensemble classifiers. In the other words, the second aspect considers effectiveness of using hybrid ensemble classifier instead on basic ensemble classifier. For the first aspect, the use of hybrid ensemble classifier for heart disease detection and prediction has reached to 80%, 93% and 96% for sensitivity, specificity and classification accuracy which is 12%, 18% and 18% more than results of best base classifier (assuming highest values of measures for k-NN and Random Tree base classifiers). For the second aspect, based on Table III, it can be seen that better performance is achieved by applying a hybrid ensemble classifier instead of a basic ensemble classifier and the results show 12%, 19% and 18% improvement is sensitivity, specificity and classification accuracy respectively. These results show that the idea of proposed hybrid ensemble model has improved the ability and effectiveness of heart disease detection and prediction artificial intelligence models.

IX. CONCLUSIONS

The proposed heart disease detection and prediction model enables the physician to predict and diagnose the heart disease by investigating and analyzing Single Proton Emission Computed Tomography (SPECT) images of patients. The artificial intelligence models that use SPECT images have been underscored in the previous studies. However, there is a limited number of works that underscore use of a hybrid ensemble classifier in a heart disease detection and prediction artificial intelligence model. Therefore, this study introduces a new approach that merges initial features of samples and base classifier predictions to produces a new feature vector for fuser classifier. It culminates with the formulation of a new model, which is considered as a novel of the present study.

In order to build a reliable model, this study investigated different fuser classifiers and considered comparison between basic ensemble and hybrid ensemble as well as comparison between hybrid ensemble and base classifiers. The results obtained from different configurations of the model indicate that the proposed model is a more reliable system that can support clinical decision makers by providing more reliable information. The proposed model is an effective artificial intelligence model for predicting heart disease, especially in terms of sensitivity and specificity that are clinically important evaluation measures. This improvement would increase the performance of the heart disease CAD systems in the clinical environments. As a conclusion, this study confirms that merging initial features of samples with predicted class labels of samples by different classification algorithms would be advantageous for the clinical decision makers.

X. FUTURE WORKS

Our study raises a number of opportunities for future researches on heart disease prediction models. As mentioned in Section I, this study uses five classifiers in inner classifiers module. This limitation is due to a tradeoff between model simplicity and maximum possible values of evaluation measures. Although this study outlines the model simplicity, however, it is a challenge to add more classifiers to reach better performance. Future researches may also tackle the proposed model by applying more fuser classifiers. In addition, another opportunity for future researches would be extending the proposed model for other types of diseases.

REFERENCES

- [1] Soni, J., et al., Predictive data mining for medical diagnosis: An overview of heart disease prediction. *International Journal of Computer Applications*, 2011. 17(8): p. 43-48.
- [2] Sarwar, A., V. Sharma, and R. Gupta, Hybrid ensemble learning technique for screening of cervical cancer using Papanicolaou smear image analysis. *Personalized Medicine Universe*, 2015. 4: p. 54-62.
- [3] Mohapatra, S., D. Patra, and S. Satpathy, An ensemble classifier system for early diagnosis of acute lymphoblastic leukemia in blood microscopic images. *Neural Computing and Applications*, 2014. 24(7-8): p. 1887-1904.
- [4] A.R., E. Naderi, and A. Osareh. Parallel weak learners, a novel ensemble method. in *Computational Intelligence and Computing Research (ICCIC)*, 2010 IEEE International Conference on. 2010. IEEE.
- [5] Woźniak, M., M. Graña, and E. Corchado, A survey of multiple classifier systems as hybrid systems. *Information Fusion*, 2014. 16: p. 3-17.

- [6] Dua, D. and E.K. Taniskidou, UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. 2017.
- [7] Kurgan, L. and K.J. Cios. Ensemble of classifiers to improve accuracy of the CLIP4 machine-learning algorithm. in *Sensor Fusion: Architectures, Algorithms, and Applications VI*. 2002. International Society for Optics and Photonics.
- [8] Thanigaivel, R. and K.R. Kumar, Review on heart disease prediction system using data mining techniques. *Asian Journal of Computer Science and Technology (AJCST)* Vol, 2015. 3: p. 68-74.
- [9] Çomak, E., A. Arslan, and İ. Türkoğlu, A decision support system based on support vector machines for diagnosis of the heart valve diseases. *Computers in Biology and Medicine*, 2007. 37(1): p. 21-27.
- [10] Polat, K., et al. A new classification method to diagnosis heart disease: supervised artificial immune system (AIRS). in *proceedings of the turkish symposium on artificial intelligence and neural networks (TAINN)*. 2005.
- [11] Amin, S.U., K. Agarwal, and R. Beg. Genetic neural network based data mining in prediction of heart disease using risk factors. in *Information and Communication Technologies (ICT), 2013 IEEE Conference on*. 2013. IEEE.
- [12] Das, R. and A. Sengur, Evaluation of ensemble methods for diagnosing of valvular heart disease. *Expert Systems with Applications*, 2010. 37(7): p. 5110-5115.
- [13] Das, R., I. Turkoglu, and A. Sengur, Effective diagnosis of heart disease through neural networks ensembles. *Expert systems with applications*, 2009. 36(4): p. 7675-7680.
- [14] Dai, W., et al., Prediction of hospitalization due to heart diseases by supervised learning methods. *International journal of medical informatics*, 2015. 84(3): p. 189-197.
- [15] Nahar, J., et al., Computational intelligence for heart disease diagnosis: A medical knowledge driven approach. *Expert Systems with Applications*, 2013. 40(1): p. 96-104.
- [16] Kaur, B. and W. Singh, Review on heart disease prediction system using data mining techniques. *International journal on recent and innovation trends in computing and communication*, 2014. 2(10): p. 3003-3008.
- [17] Kurgan, L.A., et al., Knowledge discovery approach to automated cardiac SPECT diagnosis. *Artificial intelligence in medicine*, 2001. 23(2): p. 149-169.
- [18] Turkoglu, I., A. Arslan, and E. Ilkay, An intelligent system for diagnosis of the heart valve diseases with wavelet packet neural networks. *Computers in Biology and Medicine*, 2003. 33(4): p. 319-331.
- [19] Trybulec, W.A., Pigeon hole principle. *Formalized Mathematics*, 1990. 1(3): p. 575-579.
- [20] Altman, D.G. and J.M. Bland, Diagnostic tests. 1: Sensitivity and specificity. *BMJ: British Medical Journal*, 1994. 308(6943): p. 1552.
- [21] Hall, M., et al., The WEKA data mining software: an update. *ACM SIGKDD explorations newsletter*, 2009. 11(1): p. 10-18.
- [22] Das, R., A comparison of multiple classification methods for diagnosis of Parkinson disease. *Expert Systems with Applications*, 2010. 37(2): p. 1568-1572.
- [23] Cheung, N., Machine learning techniques for medical analysis. School of Information Technology and Electrical Engineering. 2001, B. Sc. Thesis, University of Queensland.

M/M/1/n+Flush/n Model to Enhance the QoS for Cluster Heads in MANETs

Aleem Ali¹, Neeta Singh²

School of ICT, Gautam Buddha University,
Greater Noida (U.P.), India

Poonam Verma³

School of Information and Communication Technology
Gautam Buddha University, Greater Noida (U.P.), India

Abstract—Clustering in MANET is important to achieve scalability in presence of large networks and high mobility in order to maintain the Quality of Services (QoS) of the network. Improving the QoS is the most important and crucial issue in the area of MANET. Keeping this in mind, the research paper presents an M/M/1/n+Flush/n queueing model to perform better parametric results for cluster heads in MANETs. In an effort to make the M/M/1/n+Flush/n queueing model, the paper establishes the expressions for utilization (U_i) of the Cluster Head (CH), mean queue length (L_q), mean busy period (E_Q), mean waiting time (\bar{Q}) and average response time (\bar{R}) of the CH. The analytical results are further verified using MATLAB simulations which reveal better outcomes.

Keywords—Mobile Ad hoc Network (MANET); Cluster Head (CH); queueing approach; Quality of Services (QoS), flushing technique

I. INTRODUCTION

Due to rapid adaptation in Communication technology, regular desktop computing is changing itself at a very fast pace. This has lead to a situation where a huge number of diverse communication technologies transmits and exchange information over various network platforms [1]. In such an environment, the devices keep on adapting and reconfiguring themselves individually and jointly (forming cluster) to support the requirements of mobile users [2], [3]. Within the next generation of wireless network, there'll be a necessity for the quick deployment of mobile nodes on a spontaneous basis. MANETs accomplishes such extemporaneous communication among all nodes within the network without the occurrence of federal administration; however, all nodes can be treated as routers. This results in MANET's two of the most crucial qualities; adaptable and quick to deploy [4], [5]. Typically, a MANET is come to existence by spreading mobile nodes (MNs) in desolate areas or in adversity dominated areas where already existing networks have been exterminate or generally, are not possible. Therefore, one of the most reliable analytical approaches to predict and evaluate the system performance is to develop a stochastic model. These models also make available essential guidelines for the designing, implementation, and optimization of MANETs technologies [6].

Moreover, unlike simulation techniques, queueing models necessitate comparatively lesser information about the network. Additionally, in view of the fact that they are very fast to run, they offer an easy means to carry out "what-if"

analyses. This helps in identifying tradeoffs among the various performance measures in the network and find attractive solutions rather than just predicting performance for a given scenario.

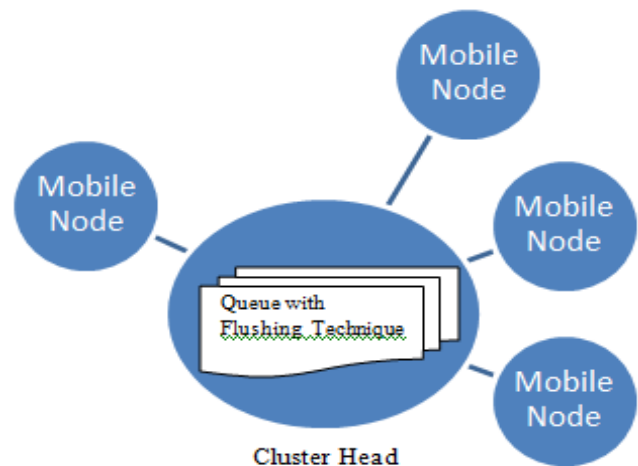


Fig. 1. Cluster head in MANET.

Queueing approach determines the performance measures in a stochastic scenario that requires some background in probability theory [7-10]. Our main objective in this paper is to determine the utilization (U_i), mean waiting time (\bar{Q}), mean queue length (L_q), distribution of the number of packets in the queue and average response time of the cluster heads when applying flushing technique in MANET (Fig. 1).

The research paper is prepared as follows: Related work and background is introduced briefly in Section II. Section III presents the proposed model's description in detail. Evaluation of the proposed mathematical model to calculate the performance measures such as utilization of the CH, mean waiting time, mean number of packets in the queue, throughput and average response time is discussed in Section IV. Section V presents the results of the proposed model. Finally, Section VI provides the conclusion of this contribution.

II. RELATED WORK AND BACKGROUND

In this section, we give a short overview of few of the recent advances in the area of queueing theory. We list the research papers that are concerned with basic queueing models including finite queue capacity with a single server as well as multiple servers.

Typical outcomes of the queuing models have widely been studied and further extended in the area of the ad hoc network. Single server queuing models [11] and multiple server queuing models [12] consisting different service disciplines have been incorporated in recent years. These models adopted various disciplines like First-In-First-Out (FIFO) [13], Last-In-First-Out (LIFO) [14], random service [15], priority service and service in batches [16]. Likewise, various arrival processes have been taken into consideration, for instance Poisson input [17], superposition of various input processes [18]. Many generalizations for queuing model have already been projected like, queues in which many customers arrive simultaneously and getting serviced in batches, priority queues and finite queues, in which the queue size is considered limited [19], [20].

In the queuing model, usually minimizing the time that packets have to wait in the queue and maximizing the utilization of the nodes are sometimes contradictory goals. It is noticed from the review of the literature that it is very difficult to perfectly find out the performance of the queuing model with the help of the classical mathematical techniques [21]-[24]. Simulation modelling has widely been a part of queue modelling and various researches in the literature preferred the use of simulation models [25]-[30]. This is the reason which has motivated us to develop a queuing model using MATLAB simulation.

In MANET, various mathematical techniques have been used for studying various performance measures. In [32], the authors addressed the end-to-end delay analysis in a single-hop wireless network. The author in [31] extended the research work mentioned in [32], to address the end-to-end delay analysis in multi-hop wireless network [47]-[50] under unsaturated traffic condition in view of the hidden and exposed terminal problem. Every node in the wireless network was modelled according to M/G/1 queue and further helped to determine the service time distribution function. With the help of the service time distribution function for a single hop, the probability distribution function (PDF) of a single hop delay and its first and second moment of service were determined.

In [33], explicit delay distribution of M/M/m/k, M/M/m/K/n queueing model with FCFS service discipline along with the mixed loss-delay system is analyzed. The model includes the balking and finite population size models as the special case. Performance evaluation of queueing system was shown in [34] and mixing time and loss priorities in a single server queue was further presented in [35].

III. PROPOSED QUEUEING MODEL M/M/1/N+FLUSH/N

In this research paper, a single node i.e., cluster head with finite queue capacity 'n' and a finite population of size 'n' is formulated. Packets arrive according to a Poisson distribution with mean rate λ and the service duration follows an exponential distribution having service rate μ . All the packets wait in the queue of the CH till they are serviced completely so as to depart from that particular CH within MANET. A cluster within MANET can be depicted as a queueing model (M/M/1/n+Flush/n) having fixed transmission range [36]-

[38]. Packets in MANET can be delivered to the CH through many intermediary nodes known as cluster members.

The M/M/1/n+Flush/n queueing model for CH comprises of two main sections, namely, queue with flushing technique and queue without flushing technique. This helps in providing a clear depiction of the network behavior under the impact of flushing technique.

Mathematical Model:

Effective arrival rate (λ_{eff}) is given by,

$$\lambda_{eff} = \begin{cases} (n - K)\lambda & 0 \leq K \leq n \\ 0 & K > n \end{cases} \quad (1)$$

with the service rate;

$$\mu_n = \mu \text{ for } K > 1 \quad (2)$$

where λ_{eff} Effective arrival rate at a CH
 μ_n Actual service rate at a CH
n Total packets in the queue of the CH
 λ Average arrival rate of packets at the CH

Assuming Poisson distribution, if P (t) is considered in the network during the time interval of length t then, Poisson distribution is represented by the following equation [15]:

$$P(S = s|\lambda) = \frac{e^{-\lambda} \cdot \lambda^n}{n!} \quad s = 0, 1, 2, \dots, \infty \quad (3)$$

Where, P is the probability that n packets exactly reach at a CH within the very short time interval ' $\Delta t (\Delta t \rightarrow 0)$ '.

In this model, flushing is applied to the queue of CH which is of finite length. Whenever the queue is full and can no longer hold the packet in it, the flushing of packet occurs. At the same time, if the CH within the MANET experiences failure, all the packets are flushed out and are queued in the flushing queue. This results in no packet loss at CH's end, thus improving the QoS of MANET. The flushing occurrences depend on the number of packets in the queue relative to the threshold value.

A. Queue with Flushing Technique

The CH's queue with flushing threshold subsystem has a signal input port labelled thresh that represents the threshold for flushing the queue. Inside the subsystem, an FCFS Queue block stores packets, while the control space usage subsystem compares the queue length to the threshold. If the queue length exceeds the threshold, the Enabled gate block (EGB) permits enough packets to depart from the queue until the queue length no longer exceeds the threshold. In this model, flushing succeeds as long as the EGB's OUT port is not blocked.

B. Queue without Flushing Technique

The CH's queue without flushing threshold subsystem requires only a signal input port and it is not labelled to any thresh used for representing the flushing of the packets in the queue. Therefore, the queue can hold the packets equal to the maximum population size (n), which may result in packet drop.

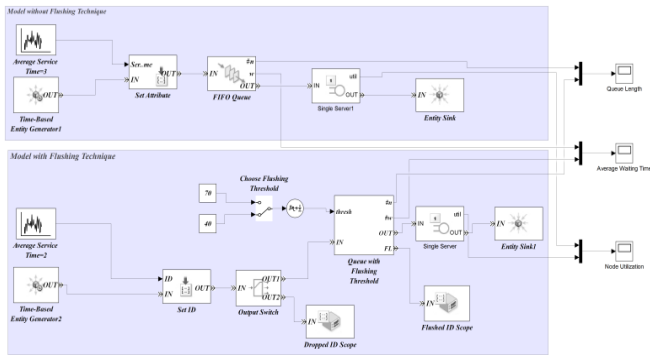


Fig. 2. M/M/1/n+Flush/n queueing model: Simulink.

IV. PERFORMANCE MEASURES FOR M/M/1/N+FLUSH/N

In this section, various performance measures [39], [40] are evaluated based on M/M/1/n+Flush/n queueing network model.

A. Model Input: Arrival rate, Average Service Time

We formulate a queueing network model for MANET with their underlying multi-hop packet forwarding. The network consists of a cluster head, with finite buffer size. In this model, packets arrive at the CH in the network follow the Poisson distribution with rate λ packets/second on an average, where packet generation at the node is independent (not depend on the behavior of earlier packets). The packet is forwarded with service rate μ packet/sec by the CH follows markovian distribution [41]-[44]. The steady-state probabilities are obtained iteratively.

For the proposed model,

Traffic intensity is given by

$$\rho = \lambda/\mu \quad (4)$$

Where, ρ represents the typical proportion of time for which the CH is occupied. Also, it is required that $\rho < 1$ for the stable queue.

Therefore, for the distribution we have the following steady-state probabilities [45]-[46];

$$n\lambda P_0 = \mu P_1 \quad (5)$$

$$(n - k\lambda + \mu)P_k = (n - K + 1)\lambda P_{k-1} - 1 + \mu P_k + 1 \quad (6)$$

$$1 \leq K \leq n - 1 \quad (7)$$

$$\mu P_K = \lambda P_{K-1} \quad (8)$$

$$P_K = P_0 \frac{\lambda^0 \lambda^1 \dots \lambda^{K-1}}{\mu^1 \mu^2 \dots \mu^K} \quad (9)$$

Let us denote Poisson distribution by $P(K, \lambda)$ and its cumulative distribution function by $Q(K, \lambda)$.

Therefore,

$$P(K, \lambda) = \left[\frac{e^{-\lambda}}{K!} \right] \lambda^K \quad 0 \leq K \leq \infty \quad (10)$$

$$Q(K, \lambda) = \sum_{i=0}^K P(i, \lambda) \quad 0 \leq K \leq \infty \quad (11)$$

By doing elementary calculations, we can obtain that

$$\frac{P(n-K, R)}{Q(n, R)} = \frac{\frac{\lambda^n}{(n-K)!} (\lambda/\mu)^K}{\sum_{i=0}^n \frac{\lambda^i}{i!} (\lambda/\mu)^{n-i}} \quad (12)$$

Where,

$$R = \frac{\mu}{\lambda} = (\rho)^{-1} \quad (13)$$

Now, we can determine that

$$\frac{P(n-K, R)}{Q(n, R)} = P_K \quad (14)$$

Thus for the distribution, we have

$$P_K = P_0 \frac{\lambda^n}{(n-K)!} \rho^K \quad (15)$$

which can be written as:

$$P_K = (n-K+1) \rho P_K \quad (16)$$

$$\sum_{n=0}^k p_n = 1 \quad (17)$$

The probability that the CH is empty (i.e., there is no packet in the queue) is given by,

$$P_0 = \frac{1}{\sum_{n=0}^k \frac{\lambda^n}{(n-K)!} \rho^K} \quad (18)$$

B. Model Output: Performance Criteria

In this paper, we take the utilization (U_i) of the CH, throughput (λ_i), mean queue length (L_q), mean waiting time (\bar{Q}) and average response time \bar{R} as the output measures for evaluating the performance of the proposed model.

1) Utilization (U_i) of the CH

Node utilization is the amount of packets which can a node hold within a specific time duration. It is one of the main measures to determine the performance of the node. The utilization of the CH can be obtained as:

$$U_t = 1 - P_0 = 1 - B(n, R) = 1 - \frac{1}{\sum_{n=0}^k \frac{\lambda^n}{(n-K)!} \rho^K} \quad (19)$$

With the help of cumulative distribution function, we can write this expression as:

$$U_t = \frac{Q(n-1, R)}{Q(n, R)} \quad (20)$$

Hence, throughput for the CH in the MANET can be given by,

$$\lambda_t = \mu \cdot U_t \quad (21)$$

$$\lambda_t = \mu \left(1 - \frac{1}{\sum_{n=0}^k \frac{\lambda^n}{(n-K)!} \rho^K} \right) \quad (22)$$

2) Mean queue Length (L_q)

It is the total number of packets waiting in the queue and can be determined as follows:

$$L_q = \sum_{i=0}^K (K-1) P_K = n - (1 + R) \cdot U_t \quad (23)$$

Here, mean busy period (E_Q) of the CH can be derived as

$$(E_Q) = \frac{1-P_0}{n\lambda.P_0} \quad (23)$$

Where, P_0 is the probability of empty queue, n is the maximum buffer size and λ is the mean arrival rate of the packets at the queue of the *CH* in the MANET.

3) Mean Waiting Time (\bar{Q})

It is the time that a packet has to wait in the queue to get its chance for getting service. It can be determined as the follows:

$$U_t = \frac{\frac{1}{\lambda}}{\frac{1}{\lambda} + \bar{Q} + \frac{1}{\mu}} \quad (24)$$

Hence, by Little’s law, we have

$$\bar{Q} = \text{Mean queue length } (L_q) / \text{Actual arrival rate}$$

$$\bar{Q} = \text{Mean queue length } (L_q) / \lambda . U_t . R$$

and it can be derived by simple calculations as

$$\bar{Q} = \frac{1}{\mu} (\frac{n}{U_t} - R) \quad (25)$$

4) Average Response Time

It is the total amount of time a packet has to wait to get serviced (including waiting time and service time). The average response time for all the packets in the network results in

$$\bar{R} = \lambda (\bar{Q} + \frac{1}{\mu}) \quad (26)$$

$$\bar{R} = (n - U_t . R (1 + \frac{1}{R})) + U_t \quad (27)$$

$$\bar{R} = (n - U_t . R) \quad (28)$$

It is worth noting that utilization of the *CH* plays a very crucial role in determining the rest of performance measures of the network.

V. RESULTS AND DISCUSSION

The simulation is performed using MATLAB to analyse the performance of the proposed model. Simulation parameters values are shown in Table I. Analytical results are validated by simulations which show better outcomes.

TABLE I. SIMULATION PARAMETERS

Parameter	Value
Simulation Time	300 sec
Arrival Rate Distribution	Poisson
Service Rate Distribution	Exponential
Queue Capacity	100 packets
Traffic Generator	CBR
Threshold Upper	70
Threshold Lower	40
Packet Size	512 byte
Number of Nodes	1 (Cluster Head)

The performance analysis is based on three performance metrics regarding QoS requirements, i.e., the node utilization, mean number of packets in the queue (buffer capacity) and mean waiting time of the packet at the node. The buffer capacity of FIFO queue in the *CH* is considered of limited length. Lower and upper flushing threshold values are set to 40 and 70, respectively.

Mean number of packets in the *CH* is compared between traditional and proposed *M/M/1/n+Flush/n* queueing model in Fig. 2. Simulation result shows the average number of packets in the queue remains same for the very short period of time, but after that, for the traditional technique, it increases linearly. For our proposed model it keeps on fluctuating near about a fixed buffer capacity, leaving no scope for queue overload. It is observed that queue length of the models does not exceed the average length of the queue at any time instant.

Simulation result in Fig. 3 shows the comparison between the average waiting time for the proposed *M/M/1/n+Flush/n* model and conventional *M/M/1/n/n*. Initially, till 100 sec, the value of the average waiting time remains same for both the techniques. But in the conventional model, after 100 sec it increases exponentially and keeps on fluctuating depending on the size of the queue. It can be clearly seen that the waiting time for each packet in the queue for the proposed *M/M/1/n+Flush/n* model is comparatively lesser than the conventional model. The peak value of average waiting time for the conventional model is more than 17 msec. On the contrary, for the proposed *M/M/1/n+Flush/n* model, it does not exceed 4 msec. It means that a packet for our proposed model has to wait for almost one fourth lesser time experienced by conventional model. This shows that the proposed model outperforms the conventional model.

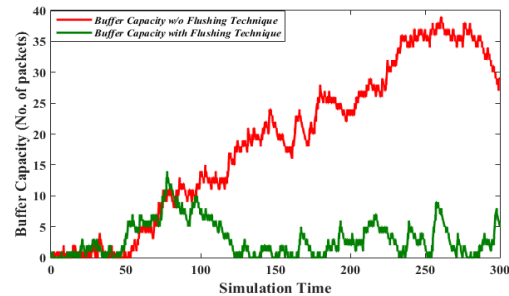


Fig. 3. Number of packets for proposed *M/M/1/n+Flush/n*.

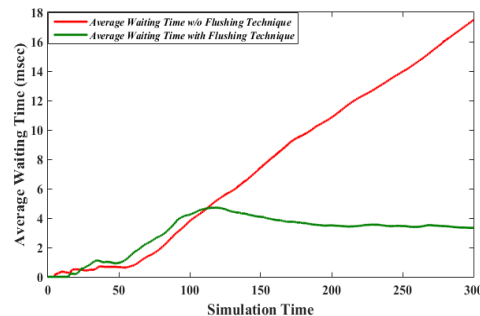


Fig. 4. Mean waiting time: Proposed *M/M/1/n+Flush/n* vs conventional *M/M/1/n/n*.

In Fig. 4, simulation result presents the comparison between CH utilization for M/M/1/n+Flush/n model and the conventional M/M/1/n/n. Utilization mainly depends upon the queue content, arrival and service pattern of the packets, and thus, the more is queue content, the larger is the value of utilization and vice-versa. It can be analyzed from Fig. 4 that the proposed model outperforms the conventional model in terms of CH utilization.

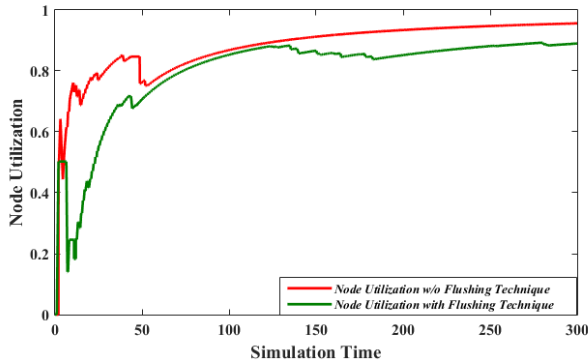


Fig. 5. Node utilization: Proposed M/M/1/n+Flush/n vs conventional M/M/1/n/n.

Table II presents the relationship between average waiting time and utilization with respect to the arrival rate of packets. It can be clearly seen that with increasing values of arrival rate, the average waiting time and utilization is also bound to increase. This shows a positive relationship between the performance parameters. Likewise, enhancement in average response time is caused due to the increment of waiting time. This is further validated in Fig. 3 which is reflecting outperformance of M/M/1/n+Flush/n model over conventional one.

TABLE II. AVERAGE WAITING TIME AND UTILIZATION

Arrival Rate (Packet/sec)	AWT in msec	Utilization of Node
1	1.080	24.9600
2	14.9980	99.9900
3	54.9993	299.9967
4	114.9997	599.9983
5	114.9997	599.9983

VI. CONCLUSION

This research paper projected an M/M/1/n+Flush/n queuing model for improving the performance measures of cluster heads in MANET, in order to enhance its QoS. The authors have developed the equations of the steady state probabilities based on which some crucial performance measures like CH utilization, queue length, average waiting time, node and average response time are obtained. Analytical results are verified with the help of simulations using MATLAB Simulink. Simulation result shows that our proposed model M/M/1/n+Flush/n outperforms the conventional model (M/M/1/n/n). In Fig. 3, 4 and 5, we have

found that results of queue capacity, average waiting time and node utilization are far better than the conventional model.

REFERENCES

- [1] S. Mishra, I. Zhang, S. C. Mishra, Guide to Wireless Ad Hoc Networks, Springer-Verlag, London, 2009.
- [2] Jonathan Loo, J. L. Mauri and J. Hamilton, Mobile Ad hoc Networks: Current Status and Future Trends, Taylor and Francis, CRC Press, 2011.
- [3] P. M. Ruiz, G. L. Aceves, J. J. Ad-Hoc, Mobile and Wireless Networks, Proceeding 8th International Conference, ADHOC-NOW 2009, Springer, Murcia, Spain, September 22-25, 2009.
- [4] S. Chakrabati and A. Mishra, "QoS Issues in Ad Hoc Wireless Networks," IEEE Communication Magazine, vol. 39, Issue 2, pp. 142-148, 2001.
- [5] M. Karimi and D. Pan, "Challenges for Quality of Service (QoS) in Mobile Ad-Hoc Networks (MANETs)," in proc. 2009 IEEE 10th Annual Wireless and Microwave Technology, 2009, Clearwater, FL, USA, pp.1-5.
- [6] K. Erciyes, O. Dagdeviren, O., Cokuslu, O. Yilmaz and H. Gumus, Modeling and Simulation of Mobile Ad hoc Network, Taylor and Francis group, 2010, CRC Press.
- [7] Ivo Adan and Jacques Resing, Queueing Theory, Eindhoven University of Technology Netherlands, February 14, 2001.
- [8] A. Lee and I. Ra, "A Queueing Network Model Based on Ad hoc Routing Networks for Multimedia Communications," Applied Mathematics and Information Sciences, vol. 6, no. 1, USA, Jan. 2012.
- [9] A. H. Zakaria, Md. Yazid, A.S.M. Noor, R.O.M Saleh, "Performance Analysis of MANETs using Queueing Theory", Proceedings of the First International Conference on Advanced Data and Information Engineering (DaEng-2013), Springer, pp.-555-562, 2014.
- [10] Andreas Willig, A Short Introduction to Queueing Theory, Technical University Berlin, Telecommunication Networks Group, Berlin, July 21, 1999.
- [11] D.J. Daley, The Correlation Structure of the Output Process of Some Single Server Queueing Systems, The Annals of Mathematical Statistics 39 (3) (1968) 1007-1019.
- [12] J. Keilson, A. Kooharian, On the General Time Dependent Queue with a Single Server, The Annals of Mathematical Statistics 33 (2) (1962) 767-791.
- [13] S. Karlin, J. Mc Gregor, Many server queueing processes with Poisson input and exponential service times, Pacific Journal of Mathematics 8 (1) (1958) 87-118.
- [14] V. Limic, A LIFO queue in heavy traffic, The Annals of Applied Probability 11 (2) (2001) 301-331.
- [15] L. Flatto, The waiting time distribution for the random order service M/M/1 queue, The Annals of Applied Probability 7 (2) (1997) 382-409.
- [16] B McK. Johnson, The Ergodicity of Series Queues with General Priorities, The Annals of Mathematical Statistics 36 (6) (1965) 1664-1676.
- [17] S. Karlin, J. Mc Gregor, Many server queueing processes with Poisson input and exponential service times, Pacific Journal of Mathematics 8 (1) (1958) 87-118.
- [18] D. Oakes, Random Overlapping Intervals-A Generalization of Erlang's Loss Formula, The Annals of Probability 4 (6) (1976) 940-946.
- [19] P.D.Finch, On the Transient Behavior of a Queueing System with Bulk Service and Finite Capacity, The Annals of Mathematical Statistics 33 (3) (1962) 973-985.
- [20] B McK. Johnson, The Ergodicity of Series Queues with General Priorities, The Annals of Mathematical Statistics 36 (6) (1965) 1664-1676.
- [21] D.G. Kendall, Some Problems in theory of queues, Journal of the Royal Statistical Society (B) 13 (2) (1951) 151-157 and 184-185.
- [22] D.G. Kendall, Stochastic Processes Occurring in the Theory of Queues and their Analysis by the Method of the Imbedded Markov Chain, Annals of Mathematical Statistics 24 (3) (1953) 338-354.
- [23] R. Schassberger, On the Waiting Time in the Queueing System GI/G/1, The Annals of Mathematical Statistics 41 (1) (1970) 182-187.

- [24] K.Udagawa, G. Nakamura, On a certain queuing system, Kodai Mathematical Seminar Reports 8 (3) (1956) 117-124.
- [25] K. Schabowicz, B. Hola, Mathematical-neural model for assessing productivity of earthmoving machinery, Journal of Civil Engineering and Management 13 (1) (2007) 47-54.
- [26] K. Schabowicz, B. Hola, Application of artificial neural networks in predicting earthmoving machinery effectiveness ratios, Archives of Civil and Mechanical Engineering 8 (4) (2008) 73-84. Sivakami Sundari et al. / Simulation of M/M/1 Queuing ... 293
- [27] J. Xu, S. Wong, H. Yang, C. Tong, Modeling Level of Urban Taxi Services Using Neural Network, Journal of Transportation Engineering 125 (3) (1999) 216-223.
- [28] Yousefi zadeh, Homayoun, A neural-based technique for estimating self-similar traffic average queuing delay, IEEE Communications Letters 6 (10) (2002) 419-421.
- [29] Yousefizadeh, Homayoun, E.A. Jonckheere, Dynamic neural based buffer management for queuing systems with self similar characteristics, IEEE Transactions on Neural Networks 16 (5) (2005) 1163-1173.
- [30] DH.M. Zhang, G. Stephen, Ritchie, Freeway ramp metering using artificial neural networks, Transportation Research Part C: Emerging Technologies 5 (5) (1997) 273-286.
- [31] E. Ghadimi, A. Khonsari, A. Diyanat, M. Farmani, and N. Yazdani, "An analytical model of delay in multi-hop wireless ad hoc networks," Wireless Networks, vol. 17, no. 7, pp. 1679-1697.
- [32] F. Alizadeh-Shabdiz, and S. Subramaniam, "Analytical Models for Single-Hop and Multi-Hop Ad Hoc Networks," Mobile Networks and Applications, vol. 11, no. 1, pp. 75-90, 2006.
- [33] H. Takagi, Y. Ozaki. Explicit delay distribution in FCFS M/M/m/K and M/M/m/K/n queues and mixed loss-delay system, in proc. Asia Pacific Symposium on Queuing Theory and Its application to Telecommunication Networks, pp-1-11, 2006.
- [34] Oliver C. Ibe, M/G/1 Vacation Queuing Systems with Server Timeout, American Journal of Operations Research, 2015, 5, 77-8.
- [35] A. Gravey, G. Hebuterne, Mixing time and loss priorities in a single server queue, TELE TRAFFIC AND DATA TRAFFIC, Elsevier, pp. 147-152, 1991.
- [36] I. Adan and J. Resing, Queuing Systems, Eindhoven University, The Netherlands March 26, 2015.
- [37] A. O. Allen, Probability, Statistics, and Queuing Theory with Computer Science Applications, 2nd ed., Orlando, USA, Academic Press, 1978.
- [38] Moshe Zukerman, "Introduction to Queuing Theory and Stochastic Teletraffic Models," EE Department, City University of Hong Kong, 6 June 2016.
- [39] Aleem Ali and Neeta Singh, "Qos Analysis in MANETs Using Queuing Theoretic Approaches: A Review," International Journal of Latest Trends in Engg. and Technology (IJLTET), vol. 7, issue 1, pp. 120-124, May 2016.
- [40] G. Bolch, S. Greiner, H. D. Meer and K. S. Trivedi, Queuing Networks and Markov Chains: Modeling and Performance Evaluation with Computer Science Applications, New York, USA, John Wiley & Sons, 1998.
- [41] B. Filipowicz, J. Kwiecien, Queuing systems and networks. Models and applications, Bulletin of the Polish Academy of Technical Sciences, 56(4), Poland, 2008, 379-390.
- [42] Y. Yeh, Handbook of Healthcare Delivery Systems, 1st ed., Taylor & Francis Group, CRC Press, 2011.
- [43] D. Gross and C. M. Harris, Fundamental of Queuing Theory, 3rd ed., John Wiley and Sons Inc., NY, USA, 1998.
- [44] U. Narayan Bhat, Simple Markovian Queuing Systems: An Introduction to Queuing Theory, Springer, India, 2008.
- [45] A. Ali, N. Singh. "Queuing approach based MANET performance analysis", IEEE Confluence 2017, 7th International Conference on Cloud Computing, Data Science & Engineering, Noida, India, pp. 422-424, Jan 2017.
- [46] A. Ali, N. Singh. "An analytical model for performance analysis of mobile ad hoc network using queuing approach", IEEE ICRITO 2017, 6th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions), Noida, India, pp. 396-399, Sep 2017.
- [47] D., C., Stiliadis, D., Ramanathan, P., "Proportional Differentiated Services: Delay Differentiation and Packet Scheduling", ACM SIGCOMM Computer Communication Review, vol. 29, no. 4, pp. 109-120, 1999.
- [48] J. Wei, Q. Li, C. Z. Xu, "VirtualLength: A New Packet Scheduling Algorithm for Proportional Delay Differentiation", China, pp. 331-336, 2003.
- [49] A. Striegel, G. Manimaran, "Packet scheduling with delay and loss differentiation", Computer Communications, vol 25, no. 1, pp. 21-31, January 2002.
- [50] T. Nandagopal, N. Venkitaraman, R. Sivakumar, and V. Bharghavan. "Delay differentiation and adaptation in core stateless networks" In Proceedings of IEEE INFOCOM, pp. 421-430, April 2000.

Binary PSOGSA for Load Balancing Task Scheduling in Cloud Environment

Thanaa S. Alnusairi

College of Computer Sciences and Information
Aljouf University, Skaka, Aljouf

Ashraf A. Shahin^{1,2}, Yassine Daadaa¹

¹College of Computer and information Sciences

¹Al Imam Mohammad Ibn Saud Islamic University (IMSIU), Riyadh, Saudi Arabia

²Department of Computer and Information Sciences, Institute of Statistical Studies & Research,
Cairo University, Cairo, Egypt

Abstract—In cloud environments, load balancing task scheduling is an important issue that directly affects resource utilization. Unquestionably, load balancing scheduling is a serious aspect that must be considered in the cloud research field due to the significant impact on both the back end and front end. Whenever an effective load balance has been achieved in the cloud then good resource utilization will also be achieved. An effective load balance means distributing the submitted workload over cloud VMs in a balanced way, leading to high resource utilization and high user satisfaction. In this paper, we propose a load balancing algorithm, Binary Load Balancing – Hybrid Particle Swarm Optimization and Gravitational Search Algorithm (Bin-LB-PSOGSA), which is a bio-inspired load balancing scheduling algorithm that efficiently enables the scheduling process to improve load balance level on VMs. The proposed algorithm finds the best Task-to-Virtual machine mapping that is influenced by the length of submitted workload and VM processing speed. Results show that the proposed Bin-LB-PSOGSA achieves better VM load average than the pure Bin-LB-PSO and other benchmark algorithms in terms of load balance level.

Keywords—Gravitational search algorithm; load balancing; particle swarm optimization; task scheduling; task-to-virtual machine mapping; virtual machine load

I. INTRODUCTION

In the last few years, cloud computing has emerged as a new computing paradigm that primarily aims to provide reliable, customized, and Quality of Service guaranteed dynamic computing environments for end users. Simply, cloud computing is the technology that provides a shared pool of computing resources in the base of on-demand services. In other words, cloud computing is the delivery of computing services such as hosts, storage, databases, networking, software, and more over the Internet. In fact, there are three basic models of services, namely Infrastructure as a Service (IaaS), Platform as a Service (PaaS), and Software as a Service (SaaS). First, in the service model IaaS, a cloud provider delivers datacenters, hosts and virtual machines, storage, networks, and operating systems to cloud users on a pay-as-you-go basis. Second, the service model PaaS delivers

services that supply an on-demand environment for cloud users such as developing, testing, delivering, and managing software applications. It is mainly used by application and software developers. Third and finally, the service model SaaS delivers software applications over the Internet to cloud users on-demand and typically on a subscription basis. It is essential to cloud providers to tend the management operations in both task-level and resource-level services. The task-level scheduling allocates a task to a virtual machine (which we address in this study), while the resource-level scheduling allocates a virtual machine to a host. The other important issue is to keep cloud resources balanced. Therefore, they also tend to schedule the incoming application requests to virtual machines in order to complete submitted tasks at the expected time in a balanced way. Numerous objectives have been addressed in the literature, such as minimizing makespan, maximizing load balancing, minimizing flowtime, and minimizing monetary cost.

Considering that task scheduling is NP-complete, many heuristics have already entered the scene, and some have emerged. For instance, Greedy heuristic, Genetic heuristic, Swarm Intelligence-based heuristics such as Ant colony inspired algorithms, Bee Colony inspired algorithms, Fish-inspired algorithms, the Gravitational Search algorithm, and Particle Swarm algorithms. Swarm Intelligence (SI)-based algorithms are population-based and stochastic search algorithms, as these are evolutionary algorithms. In this work, Swarm Intelligence based algorithms are used due to their amazing results achieved in different problems. The SI concept refers to the collective behavior that emerges from the swarms of social insects. Swarms can solve complex problems that exceed the capabilities of their insects without central supervision. What is important, in SI-based scheduling algorithms, is that social insects collectively solve complex problems that are beyond their individual capabilities in an intelligent and decentralized way. As a result, these collective, intelligent, and decentralized behaviors of insects have become a model for solving the problem of task-level scheduling [1], [2]. Due to the impressive performance of SI-based algorithms, researchers have been attracted to this

strategy over the last years. Therefore, we have been inspired by the hybrid of a gravitational search algorithm and particle swarm optimization to propose a bio-inspired task scheduling algorithm that can solve the problem of load balancing task scheduling.

By definition, cloud environments are continuously changed since cloud resources are usually dynamically reallocated per demand. This behavior must be captured by the proposed load balancing task scheduling algorithm to manage the process of allocating virtual machines to tasks. The task-level load balancing can explicitly improve makespan, throughput, and scheduling time. On the other hand, resource-level load balancing has also succeeded in increasing the performance in terms of response time, VM migration number and time, and resource utilization [3]. In this paper, we propose a load balancing task scheduling algorithm that has been inspired by the hybrid of particle swarm optimization and gravitational search algorithm to find a near-optimal Task-to-Virtual machine mapping to achieve best-balanced resource allocation as well as minimize makespan. The proposed algorithm is a dynamic load balancing algorithm that is more suitable for cloud environments due to its dynamic nature. Although the dynamic load balancing algorithm can consider all changes during runtime, it achieves better results than static algorithms [4].

The rest of the paper is organized as follows. Section II discusses related works. Sections III, IV, and V present a brief introduction to the standard GSA, Standard PSO, and Standard PSOGSA algorithms, respectively. Section VI explains the proposed binary load balancing PSOGSA. Section VII discusses the complexity of the proposed algorithm. Section VIII describes the neighborhood topology of the proposed algorithm. The objective function is discussed in Section IX. Section X explains in detail the technical processes of the proposed algorithm. The experimental results are presented in Section XI. Finally, Section XII concludes the work and suggests some directions for future work.

II. RELATED WORK

This literature review presents many works under the umbrella of Load Balancing to resolve problems related to different performance parameters, such as throughput, CPU utilization, overhead, fault tolerance, migration time, number of migrations, response time, and makespan. Some of these parameters involve the efficiency of task-level schedulers, and some involve resource-level schedulers.

ParticleZ in [4] solves the problem of task scheduling in grids through four phases: job submission, queuing, node communication, and job exchanging. In particular, the role of PSO appears in the phases of communication and exchange. In those phases, particles (nodes) search to find the best position (node with minimum load among neighbor positions). Further, to guarantee load balancing all the time, each node (particle) exchanges its loads with its neighbors in a parallel way. In load search space, the lower the load is, the higher the velocity of the particle. On the other hand, to establish a fair load distribution between each neighbor's set, the exchanged loads have to be under a predefined threshold (second lightest load of neighbor's set - lightest load).

Aslanzedh et al. [5] have been inspired by the endocrine system to improve the load balancing technique in cloud management. In fact, they have combined the endocrine system and PSO (Endocrine-PSO), aiming to schedule tasks as well as minimize makespan in load-balanced resources. The Endocrine-PSO algorithm has employed the functions of hormone regulation (load balance side) and PSO (task scheduling side) to perform the scheduling process efficiently. By endocrinology science, there is the push-pull procedure, which describes the hormone regulation process in the human body (push means stimulating a hormone from a gland, while pull means inhabitation the hormone to another gland). Technically, Endocrine-PSO provides two particles: one for the push operation (St), and the other for pull (Dt). The St and Dt carry values of stimulating hormones and inhibiting hormones respectively. Results show that the Endocrine-PSO can find the best mapping, either for choosing the best task-to-VM schedule or migrating the tasks from overloaded under-loaded VM.

Different criteria have been evaluated in [6] for the PSO-based scheduling algorithm that has developed to increase the efficiency of the load balancing scheduling process in clouds. This algorithm (LBMPSTO) aims to minimize makespan, transmission time, and transmission cost, and to maximize reliability and load balancing. LBMPSTO guarantees the reliability of clouds by rescheduling tasks that have failed to be scheduled as well as guarantees load balancing between tasks and available VMs.

Jena [7] has proposed a Nested and Multi-objective PSO Framework for task scheduling in a cloud environment. Furthermore, other criteria have been addressed in this paper. The MOPSTO (Multi objective PSO) algorithm has been proposed to minimize energy and makespan. The author has hybridized PSO with an evolutionary algorithm to create the proposed multi-objective algorithm. Additionally, another concept has been considered to make solutions of MOPSTO valuable spread solutions that are selected based on Pareto dominance.

Dasgupta et al. [8] have modified the load balancing genetic algorithm with a new objective function that guarantees the user's QoS preferences as well as minimizes response time. The authors have contributed the weights of the objective function to satisfy users' preferences. This algorithm outperforms its rivals such as SCH, RR, and FCFS. However, a limitation in this load balancing algorithm has been observed: it considers that all jobs have the same priority, which is not the case of real-world jobs.

Xin Lu and Zilong Gu [9] have proposed an ACO-inspired load-adaptive cloud resource scheduling algorithm to maximize CPU utilization. It has solved two issues, the detection of hotspot node (the overloaded VM) and adaptive resource scheduling. The proposed model would monitor the CPU usage, memory, and bandwidth of all VMs within a cluster, and if a hotspot VM is detected, the scheduling process starts. The resource scheduling process is performed to find the idle node that contributes to lighten the load over the hotspot node. The authors added an expansion factor to the global update to enable faster convergence of the ant to the

path that has the desired resources by expanding its pheromone intensity. Results show that the proposed algorithm easily detects overloaded VMs and finds the nearest idle node.

More related research has addressed the load balancing issue in grid environments. Ludwig et al. [4] have introduced the AntZ approach. They have enhanced the previous works [10] and [11] to adopt the problem of load balancing efficiently. They take advantage of decay rate in [11] and mutation rate in [10] and then combine them into AntZ. Specifically, the mutation rate addresses the problem of load balancing due to its effect on guiding ants to the best path.

Dhinesh Babu and Venkata Krishna [12] have proposed a honey bee inspired load balancing algorithm (HBB-LB). The proposed HBB-LB strategy schedules tasks by taking into consideration the VMs' load balance aiming to minimize makespan, response time, and number of migrations. The proposed algorithm divides VMs based on its computing capacity into overloaded VMs, under-loaded VMs, and balanced VMs. The balancing process is performed by removing tasks from the overloaded VMs and submitting them to the under-loaded VMs with respect to task priority. The removed tasks act as honeybees, and the under-loaded VMs represent a profitable nectar source. Results show that the HBB-LB algorithm works robustly without heavy overhead, and also works efficiently in heterogeneous cloud systems.

Lili Xu and Kun Wang in [13] have proposed a green cloud task scheduling algorithm (GCTA) based on an improved binary PSO variant. In the proposed algorithm, they tried to enhance the binary PSO solution by avoiding matrix operations and using pipelined numbers for virtual machines. The authors have compared the performance of the proposed GCTA algorithm with a sequential scheduling algorithm and found that the proposed GCTA strategy achieves better performance.

III. STANDARD GSA

The GSA treats masses as search agents. The Newtonian laws of gravity and motion define how all masses move in the direction of other masses and the speed at which they do so. The greater the mass, the slower the movement and the greater the attraction to the other masses is. Since, in the GSA, a greater mass means a better solution, the GSA is seen as an excellent way to guarantee convergence with the optimum. Every mass has a position; it also has inertial, active gravitational and passive gravitational masses. Theoretical physics defines these properties in the following way [14]:

Active gravitational mass: measures how strong an object's gravitational field is. Objects with small active gravitational mass have weaker gravitational fields than objects with greater active gravitational mass.

- Passive gravitational mass: measures how strong is the interaction of an object with the gravitational field. Objects with small passive gravitational mass are subject to a weaker force than objects in the same gravitational field with larger passive gravitational mass.

- Inertial mass: measures the strength of the resistance offered by an object to changes in its motion state as a result of the application of force. Objects with large inertial mass will undergo a slower state change as a result of the application of force than objects with small inertial mass.

If we assume the existence of s masses, then the position vector of the k^{th} mass object at time(t) $X_k(t)$ will be as set out in (1):

$$X_k(t) = \{x_1^k, x_2^k, \dots, x_i^k, \dots, x_n^k\} \quad (1)$$

As well as positional property, each mass also possesses velocity and acceleration, and these may be represented using a vector.

The acceleration vector $Acc_k(t)$ of mass object k at time t is a vector of n elements as follows:

$$Acc_k(t) = \{acc_1^k, acc_2^k, \dots, acc_i^k, \dots, acc_n^k\} \quad (2)$$

while the velocity vector also has n elements as follows:

$$V_k(t) = \{v_1^k, v_2^k, \dots, v_i^k, \dots, v_n^k\} \quad (3)$$

Additionally, the vector of global best positions at time t is:

$$Xgbes_i(t) = \{xgbest_1^k, xgbest_2^k, \dots, xgbest_i^k, \dots, xgbest_n^k\} \quad (4)$$

Equation (5) sets out the force exerted on object i by the object j :

$$F_{ij}^d(t) = G(t) + \frac{M_i(t) \times M_j(t)}{R_{ij}(t) + \epsilon} (x_{id} - x_{jd}) \quad (5)$$

In this equation, M_j is the value of the mass related to mass object j , M_i is the value of the mass related to mass object i , ϵ is a small constant, and $R_{ij}(t)$ is the straight-line distance in Euclidean space between mass object i and mass object j . Equation (6) gives us the value of $G(t)$ as a function of initial value G_0 at iteration t :

$$G(t) = G_0 \times \exp\left(\frac{-\alpha \times t}{t_{\max}}\right) \quad (6)$$

In this equation, G_0 is the initial gravitational constant and α is a user-defined descending constant, t is the current iteration, and t_{\max} is the maximum number of possible iterations. $F_d^i(t)$ is the total force exerted in the d^{th} direction on mass object i and is a sum (randomly weighted) of the d components of other mass objects' forces:

$$F_d^i(t) = \sum_{j=0}^n rand_i \times F_d^{ij}(t) \quad (7)$$

In this equation, $rand_i$ is a uniform random variable in the interval $[0, 1]$.

The object i accelerates in the d^{th} direction at time t at the rate of $Acc_d^i(t)$, calculable according to (8):

$$Acc_d^i(t) = \frac{F_{i,d}(t)}{M_{ii}(t)} \quad (8)$$

where M_{ii} is the mass object i inertial mass. Equations (9) and (10) calculate, respectively, this object's next velocity and position at time $t+1$:

$$v_d^i(t+1) = rand_i \times v_d^i(t) + acc_d^i(t) \quad (9)$$

$$x_d^i(t+1) = x_d^i(t) + v_d^i(t+1) \quad (10)$$

where $rand_i$ is, once again, a uniform random constant in the interval $[0, 1]$. Its purpose is to give the search a

randomized characteristic. Current velocity and current position are, respectively, expressed as $v_d^i(t)$ and $x_d^i(t)$.

IV. STANDARD PSO

The usual use of the population-based algorithm PSO (Particle swarm optimization) is the efficient solution of problems of optimization, and PSO is one of various techniques of swarm intelligence used to solve problems of optimization.

In this class of techniques, “particles” (search agents) fly in the optimization problem’s search space. This activity is a representation of the process of searching – it is, in effect, a journey that searches for the best position that can be taken by a particle. Each search agent, or particle, is a candidate for the role of optimal optimization problem solution, and each changes velocity and position to look for an improved position in the search space. These changes of velocity and position follow the rules deduced originally from behavioral models representing the flocking of birds as proposed by Kennedy and Eberhart in [15]. In each case, calculation of next velocity and position (respectively $v_{i,d}(t+1)$ and $x_{i,d}(t+1)$) is specified (again respectively) by (11) and (12):

$$v_d^i(t+1) = (w \times v_d^i(t)) + (c_1 \times rand_i \times (pbest_d^i(t) - x_d^i(t))) + (c_2 \times rand_i \times (gbest_d^i(t) - x_d^i(t))) \quad (11)$$

$$x_d^i(t+1) = x_d^i(t) + v_d^i(t+1) \quad (12)$$

Once again, $rand_i$ is a uniform random constant in the interval [0,1] and is used to randomize the search. The $pbest_d^i$ represents the current mass object’s personal best position on the d^{th} direction, while $gbest_d^i$ represents the i^{th} mass object’s global best position on the d^{th} direction at iteration t . Current velocity and position are represented respectively by $v_{i,d}(t)$ and $x_{i,d}(t)$.

In fact, Kennedy and Eberhart have proposed another variant of PSO. In [16], they have proposed the binary version of PSO that is proposed to solve discrete problems. The significant difference in binary PSO is the way in which positions can be updated. Updating of positions is specified by finding the value of the sigmoid function for each mass’ velocities as in the following (13):

$$Sig(v_{i,d}^{t+1}) = \frac{1}{1 + \exp(-v_{i,d}^{t+1})} \quad (13)$$

Values that returned from the sigmoid function are normalized, as defined in (14):

$$x_{ij}^k = \begin{cases} 1, & \&rand_i < v_{i,d}^{t+1} \\ 0, & \text{Otherwise} \end{cases} \quad (14)$$

where $rand_i$ is a uniform random constant in the interval [0,1]. Here, the sigmoid function is used to transfer a real-valued velocity $v_{i,d}$ to a probability value in the range of [0, 1] [23].

V. STANDARD HYBRID PSO GSA

The Hybrid PSO GSA metaheuristic is a low-level bio-inspired heterogeneous hybrid algorithm. Seyedali Mirjalili and Siti Zaiton Mohd Hashim have proposed the Hybrid PSO GSA in [17] as a novel algorithm. In fact, they have

hybridized the standard PSO, and Standard GSA mentioned in the last two sections, to balance the exploration and exploitation abilities of GSA and PSO. The core idea of the Hybrid PSO GSA is to combine the exploration of GSA and the exploitation of PSO.

In other words, the strong points of both PSO and GSA were taken into consideration to improve the weakness of GSA exploitation ability as well as PSO exploration. As tested in [17], the Hybrid PSO GSA has very good exploration and exploitation abilities, which are due to its ability to avoid becoming stuck in local optima and tending to converge to the best solution quickly.

The combination of PSO exploitation and GSA exploration is translated into a new velocity equation (1). That Hybrid PSO GSA velocity integrates the velocity of both GSA and PSO to boost the balance between global search capability of GSA and local search capability of PSO. The Hybrid PSO GSA velocity equation considers the acceleration of the mass object rather than $pbest$ as in PSO velocity, which indicates that the Hybrid PSO GSA relies on the global search of PSO with the local search of GSA. The velocity of mass object i on the d^{th} dimension at next iteration ($t+1$) is $v_{i,d}(t+1)$ and its position $x_d^i(t+1)$ is calculated according to (15) and (16), respectively:

$$v_d^i(t+1) = (w \times v_d^i(t)) + (c_1 \times rand_i \times acc_d^i(t)) + (c_2 \times rand_i \times (gbest_d^i(t) - x_d^i(t))) \quad (15)$$

$$x_d^i(t+1) = x_d^i(t) + v_d^i(t+1) \quad (16)$$

where $rand_i$ is a uniform random constant in the interval [0,1]. This random number is used to give a randomized characteristic to the search, $acc_d^i(t)$ is the acceleration of the current mass object on the d^{th} direction, and $gbest_d^i$ is the global best position of i^{th} mass object on d^{th} direction in iteration t . The $v_{i,d}(t)$ and $x_{i,d}(t)$ are its current velocity and position, respectively.

A good balance between exploration and exploitation can be achieved by controlling terms of the velocity equation based on its factors w , c_1 , and c_2 . The functions of these terms and these factors are explained as follows:

1) *Momentum component* ($w \times vid(t)$): The inertial factor w characterizes inertia of masses, i.e., it controls the momentum of masses and how much mass remembers its previous velocity. Larger w causes the mass to have a better exploration ability, and smaller w values allow the mass to have a better exploitation ability.

2) *Cognitive component* ($c_1 \times rand_i \times acc_d^i(t)$): The first behavioral factor c_1 controls how much mass can be influenced by its acceleration at iteration t .

3) *Social component* ($c_2 \times rand_i \times (gbest - x_i(t))$): The second behavioral factor c_2 controls how much a mass can head toward the population’s best mass.

In the case of c_1 and c_2 , the larger values cause the mass to have a better exploitation ability. In fact, effective values can permit these three factors to achieve a good balance between exploration and exploitation.

VI. PROPOSED BINARY LOAD BALANCING PSO GSA

The Bin-LB-PSOGSA (Binary Load Balancing PSO GSA) works to distribute submitted application requests over VMs in an efficient, balanced distribution. At each time, requests of different users are submitted at different submission times to the cloud system. Then, a search process is performed by Bin-LB-PSOGSA to assign tasks of submitted application requests to VMs in a dynamic way. At the same time, a rescheduling of tasks that have already been submitted is re-applied, i.e., the task is bound to the submitted requests list, and a new search process is performed based on the new submitted requests list.

Unlike continuous search space, the search space is represented as a hypercube. Each mass moves over hypercube nodes (corners) by flipping one or more bits of the mass position matrix. Iteratively, the position matrix of a mass is binary-coded. But, the velocity matrix still consists of continuous values belonging to the real numbers. Each velocity element value holds the probability to flip or change the binary value of the corresponding position element. The process of flipping (or changing binary value) is performed by using a transfer function. In fact, transfer functions are used to determine the probability of the value of each bit in the mass position matrix (0 or 1). [18]

In the deep search process of Bin-LB-PSOGSA, each mass in the population represents one candidate solution or, in other words, a task-to-VM mapping. Each candidate solution has a fitness value, which is the value of the expected finish time of each submitted application request.

In pseudo code of the Bin-LB-PSOGSA (see Algorithm 1), first, the masses' population is initialized by the function initialize Masses (tasks, VMs) at line 2. In the initialization phase, tasks are assigned randomly to VMs. Then, for each iteration, global variables are updated that have to be changed iteratively, such as gravitational constant, best mass, and worst mass, by updateGlobalVariables(iteration) at line 4. Until the maximum iteration is reached and for each mass object in the population, the mass value of each mass object and gravitational force exerted by the population masses is calculated at lines 6 and 7, respectively. Then, at lines 8 to 10, mass position, velocity, and fitness are updated. As in Algorithm 2, the pseudo code that has clarified the way to update the velocity of each mass is presented.

After updating the mass's fitness, it is necessary to decide if the new fitness is better than what the mass object already found in its trip; this is done in line 12. If the new fitness is better, the personal best fitness and mass will be updated. Consequently, if the personal best mass is better than all solutions found by all mass objects, the global best mass is updated as in line 16. In fact, the global best mass is the promising best mass object that attracts most of the population due to its mass value (biggest mass value or heaviest mass). Based on the topology of the neighborhood that has been considered in this variant, Gbest topology (discussed in

subsection VIII), the population is influenced by the best global mass positions that are updated at line 18. In the next iteration, by updating the best global positions so far, other masses' objects take their new positions. Over time, most of the population comes increasingly closer to the best global mass, and finally, if the maximum iteration number is reached, the search process is terminated. Eventually, the best global mass is returned in the form of the best task-to-VM mapping at line 22.

VII. BIN-LB-PSOGSA COMPLEXITY

Let s be population size, v be VMs size, and c be submitted requests' tasks size. Initialization of masses is used to add random positions and velocities of each mass in the population. During initialization, the fitness of the current position of the mass is calculated. The time complexity of mass initialization is $O(v \times c)$. So, the time complexity for initialization of the whole population is $O(s \times v \times c)$.

In the iterations loop, first, global variables are being updated. The time complexity of that action is $O(s^3 + s \times c)$. The reason for such time complexity is that inside update Global Variables (iteration) there is a need for collecting the best and worst fitnesses from the whole swarm. The time complexity of those actions is $O(s)$. Inside this method, we also calculate the total forces that act on each mass and acceleration of the mass. The time complexity of those actions is $O(s^3 + s \times c)$. This is also the time complexity of the method update Global Variables (iteration).

Second, a loop is iterated for each mass in the population. Each mass's velocity and position are being updated. The time complexity of both of these updates is $O(v \times c)$. The next step is to update the fitness of the mass. The time complexity of that action is $O(v)$. The overall time complexity of the particle loop is $O(s \times v \times c)$.

For the iterations loop, the time complexity is therefore $O(s \times \text{MAX_ITERATION} \times (s^2 + v \times c))$.

The final step is to return the mapping from cloudlet to VM. The time complexity of this action is $O(v \times c)$. Eventually, if the time complexity of each step is combined, the final result is $O(s \times \text{MAX_ITERATION} \times (s^2 + v \times c))$.

VIII. NEIGHBORHOOD TOPOLOGY

The neighborhood topology adopted in the proposed Bin-LB-PSOGSA is the global neighborhood topology (Gbest) [19]. In other words, Gbest is a fully connected topology where all the masses are neighbors of each other and able to exchange information with each other. Further, the process of exchanging is fast due to the full connection between all population masses. Gbest topology makes the proposed Bin-LB-PSOGSA a fully informed strategy where every mass in the population learns from the same global best mass and is influenced by its positions.

Algorithm 1: Bin-LB-GSAPSO

```

1 function balance(tasks, vms)
  Input: tasks - list of submitted requests' tasks; vms - list of VMs
  Output: task to VM mappings
2 masses ← initializeMasses(tasks, vms);
3 for iteration ← 1 to MAX_ITERATION do
4   updateGlobalVariables(iteration);
5   foreach mass in masses do
6     calculateMass(mass, population);
7     calculateRelevantForce(mass, population);
8     updateVelocities(mass);
9     updatePositions(mass);
10    fitness ← updateFitness(mass);
11    personalBestFitness ← mass.getPersonalBestFitness();
12    if fitness < personalBestFitness then
13      personalBestFitness ← fitness;
14      personalBestPositions ← mass.getPositions();
15    end
16    if personalBestFitness < globalBestFitness then
17      globalBestFitness ← personalBestFitness;
18      globalBestPositions ← personalBestPositions;
19    end
20  end
21 end
22 return taskToVmMapping(globalBestPositions);
23 end

```

Algorithm 2: Update velocity

```

1 function updateVelocities(mass)
  Input: mass - current mass
2 for vm to vms do
3   vmVelocities[ ] ← mass.getVelocities()[vm.id];
4   vmPositions[ ] ← mass.getPositions()[vm.id];
5   for task to submittedTasks do
6     newVelocity[vm.id][task.id] = (w * vmVelocities[task.id])
7     + (c1 * Random.nextDouble() * currMass.getAcceleration()[task.id])
8     + (c2 * Random.nextDouble() * (vmGlobalbestPositions[task.id] -
9     vmPositions[task.id]));
10  end
11 end

```

IX. OBJECTIVE FUNCTION

The solution of proposed algorithms is to minimize the expected execution time of each task in submitted application requests or ET_{ij} of the task T_i that is running on VM_j . The calculation of ET_{ij} is as follows:

$$ET_{ij} = \frac{TL_i}{PS_{ij}} \quad (17)$$

The processing speed PS_{ij} related to T_i running on VM_j in the cloud depends on how many request tasks have been mapped to that VM (or n) as well as the total allocated MIPS of VM_j along all its processing elements or Pes (or Capacity $_j$). Calculation of the request processing speed PS_{ij} is shown in (18):

$$PS_{ij} = \frac{Capacity_j}{n} \quad (18)$$

X. BIN-LB-PSOGSA TECHNICAL PROCESSES

In this section, the mathematical notations and technical processes' steps are discussed. In the proposed algorithms, binary search space is considered. Therefore, the binary matrix encoding form in [13] is adopted to represent mass objects. Accordingly, each mass object has properties of mass position matrix and mass velocity matrix that will be decoded to a two-dimensional matrix. The first dimension represents the VM

number and the other tasks number. The position matrix takes binary values, while the velocity matrix keeps the continuous values. Additionally, each corresponds to a task-to-VM mapping as a candidate solution.

Here, the mathematical notation of the problem is described. The task set $T = \{ T_1, T_2, \dots, T_i, \dots, T_c \}$ will be mapped to VMs set $VM = \{ VM_1, VM_2, \dots, VM_j, \dots, VM_v \}$ using the relationship matrix representation in (18). Let s be the masses' population size, each mass object m represented by position matrix X_m of $c \times v$ position elements (c is the number of tasks and v is the number of VMs where $i = \{1, 2, 3, \dots, c\}$ and $j = \{1, 2, 3, \dots, v\}$ and $m = \{1, 2, 3, \dots, s\}$) as given in follows:

$$X_m = \begin{pmatrix} x_{11}^m & x_{12}^m & \dots & x_{1c}^m \\ x_{21}^m & x_{22}^m & \dots & x_{2c}^m \\ \dots & \dots & \dots & \dots \\ x_{v1}^m & x_{v2}^m & \dots & x_{vc}^m \end{pmatrix}$$

The element x_{ij}^m is the position of mass object m in row i and column j ; actually, it represents the distribution relationship between task T_i and virtual machine VM_i i.e., it explains whether T_i is mapped to VM_i or not. Position x_{ij}^m takes values of either 0 or 1. Namely, it indicates on which VM task T_i is working. So, if T_i is running on VM_i then position x_{ij}^m is equal to 1, but it equals 0 otherwise. Finally, positions x_{ij}^m that are equal to 1 are recorded composing solution position vector $P_m(t)$ at time t as in (19):

$$P_m(t) = \{p_{1,j}^m \quad p_{2,j}^m \quad \dots \quad p_{i,j}^m \quad \dots \quad p_{c,v}^m\} \quad (19)$$

where $p_{i,j}^m$ can take one of x_{ij}^m of the relation distribution matrix that has a value equal to 1 and $j = (1, 2, \dots, v)$.

Here, the technical steps of the proposed algorithms are explained in detail.

A. Population Initialization

Initially, position elements x_{ij}^m of each mass m position matrix X_m are initiated randomly by mapping each task to random VM, i. e. for each column at index i , one element is arbitrary assigned to the value 1 and the remaining elements (in that column) to 0. Iteratively, this process – initiating mass position matrix -- is repeated $c \times v$ times for each mass in the population.

For instance, assume that there are seven tasks and three VMs, and the population consists of 50 masses; therefore, the initial position matrix of mass 20 (X_{20}) will be as follows:

$$X_{20} = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

Namely, it is shown in the matrix (X_{20}) that T_1 is mapped to VM_2 , T_2 is mapped to VM_3 , T_3 is mapped to VM_2 , T_4 is mapped to VM_1 , T_5 is mapped to VM_1 , T_6 is mapped to VM_1 , and T_7 is mapped to VM_3 .

Finally, the tasks-to-VM vector or dimension vector D_k is defined in which indices of tasks' positions that hold value 1 are stored in it respectively as in (20):

$$D_k = \{d_1^k \quad d_2^k \quad \dots \quad d_i^k \quad \dots \quad d_c^k\} \quad (20)$$

where d_i^k is the index of assigned task T_i in position matrix X_k of the mass k , and vector indices are the VMs' IDs.

B. Mass Value, Relevance Force, and Acceleration Calculation

In this step, a mass value, relevance force, and finally the acceleration are calculated based on the gravity law equations in [14]. First, the mass value $M_m(t)$ of each mass object at iteration t is calculated based on (21) and (22):

$$mass_m(t) = \frac{fitness_m(t) \cdot worst(t)}{best(t) \cdot worst(t)} \quad (21)$$

$$M_m(t) = \frac{mass_m(t)}{\sum_1^{s-1} mass_b(t)} \quad (22)$$

where $fitness_m(t)$ is the fitness value of the mass m at iteration t , $worst(t)$ is the global worst fitness in iteration t , and $best(t)$ is the best one at that iteration. Further, $mass_m(t)$ is the current mass value where $mass_b^m(t)$ is a vector that holds the mass values of neighbor masses' objects at iteration t as shown in (23). This is because $\sum_1^{s-1} mass_b(t)$ is the summation of other masses' values. In fact, global best fitness represents the minimum expected finish time that can be achieved by the best global task-to-VM mapping since the global worst is the longest expected finish time.

$$mass_b^m(t) = \{mass_{b,1} \quad mass_{b,2} \quad \dots \quad mass_{b,(s-1)}\} \quad (23)$$

Second, the relevance gravitational forces exerted on mass m by another mass b is calculated based on (24):

$$F_d^{m,b}(t) = G(t) + \frac{M_m(t) \times M_b(t)}{R_{mb}(t) + \epsilon} (x_d^m - x_d^b) \quad (24)$$

In this equation, M_b is the value of the mass related to the active gravitational mass of object b , M_m is the value of the mass related to passive gravitational mass of object m , ϵ is a small constant, and $R_{mb}(t)$ is the straight-line distance in Euclidean space between mass object m and mass object b , and x_d^m and x_d^b are the corresponding positions at dimension d in both of the passive masses m and b . It is useful to mention that the active mass is the mass that generates the gravity, while the passive mass the mass responds to the gravity. The relevant gravitational force values are defined in a vector as in (25):

$$F_d^{m,b}(t) = \{f_1^{m,b} \quad f_2^{m,b} \quad \dots \quad f_i^{m,b} \quad \dots \quad f_c^{m,b}\} \quad (25)$$

where $f_i^{m,b}$ is the value of the exerted force on passive mass m from active mass b in the dimension d (number of dimensions d equals the number of tasks c).

Then, the total result gravitational forces exerted on mass m on the d^{th} direction at time (t) are as in (26):

$$F_d^m(t) = \sum_{k=1, k \neq m}^{s-1} rand_k \times F_d^{m,b}(t) \quad (26)$$

In this equation, $rand_m$ is a uniform random variable in the interval $[0, 1]$ generated for each mass m .

Based on Newton's law of gravity and Newton's law of motion, each mass object m moves toward the global best mass object by updating the acceleration vector of that object

iteratively. Under the concept of motion law, the acceleration of the mass object m on the d^{th} direction at time t is $Acc_{m,d}(t)$ as in the following equation:

$$Acc_d^m(t) = \frac{F_{m,d}(t)}{M_{mm}(t)} \quad (27)$$

where M_{mm} is the mass object m inertial mass.

C. Velocity Updating

The binary version of the velocity equation is as follows:

$$V_m(t+1) = (w \times v_{ij}^m(t)) + (c_1 \times rand_m \times acc_i^m(t)) + (c_2 \times rand_m \times (gbest_{ij}^m(t) - x_{ij}^m(t))) \quad (28)$$

where $V_m(t+1)$ is the velocity matrix of mass object m for the next iteration, and $v_{ij}^m(t)$ is the current velocity value of the element related to T_i and VM_i . The acceleration of mass m is $acc_i^m(t)$. The inertial weight (w) is calculated based on (29) where acceleration coefficients $C1$ and $C2$ are based on, respectively, (30) and (31). The $rand_m$ is a uniform random constant in the interval $[0, 1]$ generated for each mass object m . Its purpose is to give the search process a randomised characteristic.

For each iteration (t) , mass object m records the best global positions in its memory so that all masses can be increasingly closer until a maximum number of iterations is reached. Iteratively, the distance between masses and global best mass is decreased by subtracting the distance between positions $x_{ij}^m(t)$ and $gbest_{ii}^m(t)$, as stated in equation 28, term $(gbest_{ii}^m(t) - x_{ij}^m(t))$. Elements $Xgbest_m$ and current mass position matrix are subtracted one by one. In the case of the acceleration of masses' objects, a constant random value ($c_1 \times rand_m$) is multiplied with all elements in the acceleration vector, element by element. Also, in the current velocity matrix, the current value of the inertia weight is multiplied with all elements in the velocity matrix. The velocity matrix will be as follows:

$$V_m = \begin{pmatrix} v_{11}^m & v_{12}^m & \dots & v_{1c}^m \\ v_{21}^m & v_{22}^m & \dots & v_{2c}^m \\ v_{v1}^m & v_{v2}^m & \dots & v_{vc}^m \end{pmatrix}$$

To enhance the search process, we have considered a time-adaptive approach for the other controlling parameters such as inertia weight and acceleration coefficients (c_1 and c_2) as in (28), (29), and (30). For the inertia weight, we have adopted a time-varying inertia weight as introduced in [20] and acceleration coefficients in [21]. Here, w_{max} and w_{min} have constant values equal to 0.9 and 0.4, respectively, t is the current iteration, and t_{max} is the maximum iteration.

$$w = \frac{w_{max} - w_{min}}{t_{max}} \quad (29)$$

$$c_1 = 1 - \frac{t^3}{t_{max}^3} \quad (30)$$

$$c_2 = \frac{t^3}{t_{max}^3} \quad (31)$$

D. Position Updating

Each mass moves to the global best mass by updating positions and becomes increasingly closer to the global best mass over iterations.

Also, for the transfer function in the proposed Bin-LB-PSOGSA, we have used a time-adaptive approach as introduced in [22]. The time-varying transfer function is used here to enhance the exploration and exploitation processes. Moreover, to transform a real-valued velocity V_M to a binary value (0 or 1) in the process of updating positions (re-encoding) values of relation distribution matrix elements $x_{i,i}^k$ [13]. In fact, if the absolute value is large, the probability to flip a bit is higher. Updating of the position matrix elements is performed by applying the time-varying transfer function (TV_t) for each v_{ij}^m in mass' velocity matrix as stated in (32) and (33).

$$TV_t(t) = \frac{1}{1 + e^{-\frac{v_{ij}^m}{\varphi}}} \quad (32)$$

$$\varphi = \varphi_{max} - t \times \frac{(\varphi_{max} - \varphi_{min})}{t_{max}} \quad (33)$$

where φ_{max} and φ_{min} have constant values equal to 1.0 and 5.0, respectively.

E. Finding Global best Task-to-VM mapping

By the end of each iteration t , the best task-to-VM maps have been read from the global best mass position matrix. This process happens t times and for only the global best mass found during the searching process.

XI. EVALUATION

In this section, we show how to evaluate the proposed Bin-LB-PSOGSA to test its efficiency in achieving cloud balancing in term of the submitted load. The next subsections discuss the simulation tool and simulation setup we have used in the experiments. Additionally, we explain the algorithm meta-parameters, and finally, we conclude with the results of these experiments. In the last subsection, we assess the performance of the proposed algorithm in terms of load average and processing speed average against the Bin-LB-PSO algorithm.

A. Experimental Tool

The performance analysis of the proposed algorithm is carried out in a cloud simulator. The simulator CloudSim [23] is one of the best simulators for experimental purposes. This simulator is a generalized simulation framework that allows modeling, simulation, and experimenting with cloud computing infrastructure and application services.

In this section, we have analyzed the performance of our algorithm based on the results of simulation done using CloudSim. We have extended the classes of the CloudSim simulator to simulate our algorithm.

B. Simulation Setup

The simulation setup is detailed in Tables I and II. The experiment is carried out with 3 Datacenters each having two hosts, and the characteristics are 1024 MIPS Host processing power, 2 GB RAM, 1000 GB storage, 10240 Mbps (bandwidth), and 2 PEs (or cores). Each PE had the same processing power, as clarified in Table I.

In Table II, there are 5 VMs, and the characteristics are 128 MIPS (VM processing power) and 2 PEs (or cores).

In this experiment, the workload has been selected as introduced in [24]. Each task in the workload log, called a cloudlet by Cloudsim, was determined by the parameter PEs, or the number of processing elements (cores) required to perform each task. Each cloudlet required 4 to 256 PEs. The number of PEs is limited to powers of 2 due to the architecture of the supercomputer used in the log.

TABLE I. DATACENTER CONFIGURATION

Number of datacenters	1
Number of hosts per datacenter	4
Number of PEs per host	1
Number of MIPS per PE	1024 MI
RAM	2048 MB
Storage	1048576 MB
Bandwidth	10240 MB/s

TABLE II. VM CONFIGURATION

Number of VMs per host	5
Number of PEs per VM	2
Number of MIPS per PE	128

C. Algorithm Meta-Parameters

The algorithm meta-parameters, or in other words, the controlling parameter settings of Bin-LB-PSOGSA, are as mentioned in equations 28, 29, and 30. The maximum number of iterations is 500, and population size (number of masses) is 50. The acceleration constants $C1$ and $C2$ are set to 2 and 2, the inertial weight is linearly decreasing from 0.9 to 0.4. The initial gravitational constant ($G0$) is 1, descending constant (α) is 20, and small gravitational constant (ϵ) is e^{-1} . The search space bounds are in the range $[0, 100]$, and the velocity range is $[-8, 8]$.

D. Experimental Results

Here, the performance of the proposed algorithm in terms of average VM load and average VM processing speed is discussed. The next subsections explain the performance from different sides in detail.

1) Average VM load over time

Fig. 1 shows that in comparison with Bin-LB-PSO, the load of the proposed Bin-LB-PSOGSA is smaller than the load of Bin-LB-PSO in general. In particular, both Bin-LB-PSO and the proposed Bin-LB-PSOGSA have stable load values for a long time. This situation is due to the stability of the system because of the number of running application requests. After time passes, both as the time passes, the load increases due to the growth in the number of running application requests. Under the same environmental conditions, the proposed Bin-LB-PSOGSA outperformed Bin-LB-PSO in keeping the system balanced much longer. For instance, at moment 28, the load value of Bin-LB-PSO leaps due to the gap of some successes of application requests

before and after this moment (before it was 2 and after is 61). On the other hand, the load average in the proposed Bin-LB-PSOGSA has lower load average values than Bin-LB-PSO.

2) Average VM processing speed over time

Fig. 2 shows that over time, the average processing speed of application requests of both algorithms decreases in general. In particular, at moment 28, processing speed decreased dramatically due to the obvious increase in the running requests, which indicates the efficiency of both algorithms to utilize VMs. Therefore, it is clear that both Bin-LB-PSO and the proposed Bin-LB-PSOGSA have utilized VMs efficiently over time.

Both Fig. 1 and 2 shows that as the load increases over time, the processing speed of the submitted applications decreases, which proves that the proposed Bin-LB-PSOGSA is more efficient in keeping the load balanced over time as shown in Table III.

TABLE III. PERFORMANCE COMPARISON

Criteria	Bin-LB-PSOGSA	Bin-LB-PSO
VM processing speed (MIPS)	4376.84346	4376.84346
VM Load	389.0400819	451.3841267
Expected execution time (ms)	3900000	3800000
Performance efficiency	Better	Limited

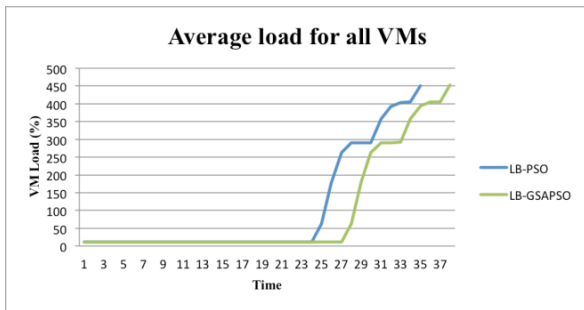


Fig. 1. Average VM load over time.

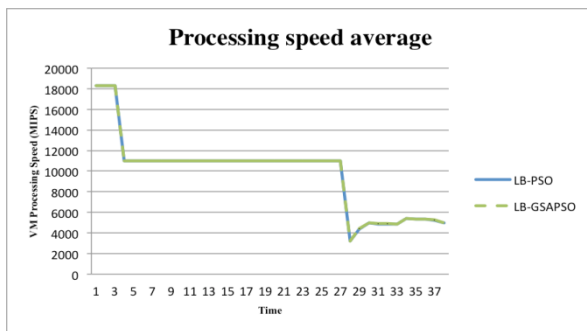


Fig. 2. Average VM processing speed over time.

XII. CONCLUSION

In this paper, we have proposed a load balancing task scheduling algorithm for cloud computing environments based on the binary hybrid gravitational search and particle swarm optimization strategy. It balances the load of application requests submitted from cloud users over virtual machines in

the cloud. The proposed algorithm enhances the overall VM utilization of the cloud system. We have compared our proposed hybrid algorithm with the pure Bin-LB-PSO. Results show that as the load increases over time, the processing speed of submitted applications decreases, which proves that the proposed Bin-LB-PSOGSA is more efficient in keeping the load balanced over time.

In the future, we plan to extend this kind of load balancing for workloads with dependent tasks. Also, we plan to improve this algorithm by considering other QoS factors, as well.

REFERENCES

- [1] Babu, K. R. Remesh, and P. Samuel. "Enhanced Bee Colony Algorithm for Efficient Load Balancing and Scheduling in Cloud." In *Innovations in Bio-Inspired Computing and Applications*, pp. 67-78. Springer International Publishing, 2016.
- [2] E. Pacini, C. Mateos, and C. G. Garino, "Distributed job scheduling based on Swarm Intelligence: A survey," *Computers & Electrical Engineering*, vol. 40, no. 1, pp. 252-269 2014.
- [3] E. Pacinia, C. Mateosb, and C. G. Garinoa, "Balancing Throughput and Response Time in Online Scientific Clouds via Ant Colony Optimization," *Advances in Engineering Software*, in press. Elsevier 2014.
- [4] S. A. Ludwig and A. Moallem, "Swarm intelligence approaches for grid load balancing," *Journal of Grid Computing*, vol. 9, no. 3, pp. 279-301, 2011.
- [5] S. Aslanzadeh and Z. Chaczko, "Load balancing optimization in cloud computing: Applying Endocrine-particulate swarm optimization," in *Electro/Information Technology (EIT)*, 2015 IEEE International Conference on, IEEE, 2015, pp. 165-169.
- [6] Awad, A. I., N. A. El-Hefnawy, and H. M. Abdel_kader. "Enhanced Particle Swarm Optimization for Task Scheduling in Cloud Computing Environments." *Procedia Computer Science* 65 (2015): 920-929.
- [7] R. K. Jena, "Multi Objective Task Scheduling in Cloud Environment Using Nested PSO Framework," *Procedia Computer Science*, vol. 57, pp.1219-1227, 2015.
- [8] K. Dasgupta, B. Mandal, P. Dutta, J. K. Mandal, and S. Dam, "A genetic algorithm (ga) based load balancing strategy for cloud computing," *Procedia Technology*. vol. 10, pp. 340-347, 2013.
- [9] X. Lu and Z. Gu, "A load-adaptive cloud resource scheduling model based on ant colony algorithm," in *Cloud Computing and Intelligence Systems (CCIS)*, 2011 IEEE International Conference on, IEEE, 2011, pp. 296-300.
- [10] S-S. Kim, J-H. Byeon, H. Liu, A. Abraham, and S. McLoone, "Optimal job scheduling in grid computing using efficient binary artificial bee colony optimization," *Soft Computing*, vol. 17, no. 5, pp. 867-882, 2013.
- [11] Z. Mousavinasab, R. Entezari-Maleki, and A. Movaghar, "A bee colony task scheduling algorithm in computational grids," in *Digital Information Processing and Communications*, Heidelberg, Springer Berlin, 2011, pp. 200-210.
- [12] D. Babu and P. V. Krishna, "Honey bee behavior inspired load balancing of tasks in cloud computing environments," *Applied Soft Computing*, vol. 13, no. 5, pp. 2292-2303, 2013.
- [13] L. Xu, K. Wang, Z. Ouyang, and X. Qi, "An improved binary PSO-based task scheduling algorithm in green cloud computing. In *Communications and Networking in China (CHINACOM)*," 2014 9th International Conference on, IEEE, 2014, pp. 126-131.
- [14] E. Rashedi, H. Nezamabadi-Pour, and S. Saryazdi, "GSA: a gravitational search algorithm," *Information Sciences*, vol. 179, no. 13, pp. 2232-2248, 2009.
- [15] Kennedy J, Eberhart R. "Particle swarm optimization." *Proceedings of the 4th IEEE International Conference on Neural Networks*, 1995, pp. 1942-1948.
- [16] J. Kennedy, "A discrete binary version of the particle swarm algorithm," *Proceedings of the 1997 IEEE International Conference on Systems, Man and Cybernetics*, vol. 5, IEEE press, New York, NY 1997, pp. 4104-4108.

- [17] Mirjalili, S., & Hashim, S. Z. M. (2010, December). A new hybrid PSOGSA algorithm for function optimization. In *Computer and information application (ICCIA), 2010 International Conference on* (pp. 374-377). IEEE.
- [18] F. Marini and B. Walczak, "Particle swarm optimization (PSO). A tutorial," *Chemometrics and Intelligent Laboratory Systems*, vol. 149, pp. 153-165, 2015.
- [19] D. Bratton, D. and J. Kennedy (2007, April). "Defining a standard for particle swarm optimization In *Swarm Intelligence Symposium, 2007. SIS 2007. IEEE* (pp. 120-127). IEEE.
- [20] C. Yang, W. Gao, N., Liu, and C. Song (2015). "Low-discrepancy sequence initialized particle swarm optimization algorithm with high-order nonlinear time-varying inertia weight." *Applied Soft Computing*, 29, 386-394.
- [21] G. Sun, A. Zhang, Z. Wang, Y. Yao, J., Ma, and G. D. Couples, (2016). Locally informed gravitational search algorithm. *Knowledge-Based Systems*, 104, 134-144.
- [22] J. Islam, X. Li, and Y. Mei, "A Time-Varying Transfer Function for Balancing the Exploration and Exploitation ability of a Binary PSO," *Applied Soft Computing*, vol. 59, pp. 182-196, 2017.
- [23] R. N. Calheiros, R. Ranjan, R., C. A. De Rose, and R. Buyya, , "CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms," *Software: Practice and Experience*, vol. 41. no. 1 23-50, 2011.
- [24] K. Windisch, V. Lo, R. Moore, D. Feitelson, and B. Nitzberg, "A comparison of workload traces from two production parallel machines," in *Frontiers of Massively Parallel Computing, Proceedings Frontiers' 96, Sixth Symposium on the IEEE*. 1996, pp. 319-326.

The P System Design Method based on the P Module

Ping Guo, Xixi Peng, Lian Ye
College of Computer Science
Chongqing University
Chongqing, 400044
China

Abstract—Membrane computing is a kind of biocomputing model. At present, the main research areas of membrane computing are computational models and P system design. With the expansion of the P system scale, how to rapidly construct the P system has become a prominent issue. Designing P system based on P module is a P system design method proposed in recent years. This method provides information hiding and can build P system through recursive combination. However, the current P module design lacks a unified design method and lacks the standard process of building P system from P module. This paper studies the structural characteristics of cell-like P systems, and proposes an improved P module design method and a process for assembling P systems through P modules. In order to fully expound the design method of P module, the P system for the square root of the large number was analyzed and designed. And the correctness of the P system based on the P module design method was verified by an instance.

Keywords—P module; P System; P system design; membrane computing; biocomputing models

I. INTRODUCTION

Membrane computing, also known as P system, is a branch of natural computing [1]. The models of P system are mainly divided into three types, namely, the cell-like P system [2], the tissue-like P system [3] and the neural-like P system [4]. They have been applied to solve the problems such as NP problems [5]–[8], image processing [9], [10], arithmetic operations [11]–[14] and so on.

In our previous work, most energy are put to implement the arithmetic operations in cell-like P systems: Ref. [11] firstly proposed an arithmetic P systems to implement the arithmetic operation in 2001; in [12] proposes an algorithm and builds expression P systems without priority rules for evaluating arithmetic expression; in [13] designed the P systems for addition, subtraction and multiplication; in [14] proposes a family of systems for solving Matrix-Vector Multiplication. Although we have obtained many excellent research results in the cell-like P system, the difficulty for constructing the P system has continued to increase due to the increasingly complex algorithm. As a result, some experts and scholars began to propose some new models of modular constructing the P system. In 2009, Romero-Campero et al. proposed a biology model for modular combination cells [15] and Serbanuta et al. proposed K systems embedded in P system which can develop new extensions of P system [16]. In 2010, Păun et al. proposed the dp system [17], which contains ideas for modularity.

Based on the above models, the modularized construction model, the P module [18], is proposed for simplifying the computing system structure and improving the reusability. At present, this model has been applied to some areas of research. In [19] proposed an improved generic version of P modules, an extensible framework for recursive composition of P systems. It proposed an solution to solve Byzantine agreement problem by P module. In [20] presented an improved deterministic solution for Flow-shop Scheduling problem. In [21] extended the P module theoretically and proposes the P module to solve the stereo matching problem in the application. Besides, it realized the discovering neighbors and Echo Algorithm. In [22] studied on the problem which aims to find out a point-disjoint and edge-disjoint path between source point and target point. All of these literature researches are related to the algorithms application of the P module. However, due to the lack of a unified design method, the P system based on the P module have low design efficiency and high error rate. In order to improve such problems, this paper designs a well-structured P module by combining the structural design methods in design methodology. The correctness of the dynamic execution of the P system is ensured with a good structure, making the P system easy to understand, easy to debug, and easy to maintain.

In this paper, the cell-like P system and the P module are introduced in Section 2. Section 3 improves the P module and proposes the design and assembly of the P module. With the method of structural design in design methodology, four methods for constructing P module are proposed to design well-structured P system based on P modules in Section 4. Section 5 gives an instance to show the working mechanism of P module by using the related definitions and design methods of the P module. Section 6 summarizes the research work and presents a deeper level of research in the future.

II. FOUNDATIONS AND RELATED WORKS

A. Cell-like P System

The cell-like P system is a class of P system constructed by biochemical reactions in abstract biological cells. In the cell-like P system, the substances in the cells are abstracted as computational objects and the biochemical reactions within the cells are abstracted as object evolutionary rules. A cell-like P system containing five membranes is shown by Fig. 1.

Fig. 1 is a schematic representation of a cell-like P system. A cell-like P system consists of the membranes (elementary membrane and combination membrane), the membrane regions

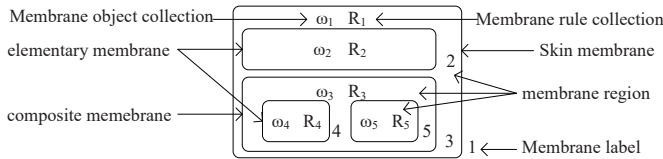


Fig. 1. The structure of cell-like P system.

surrounded by membranes, the membrane object collection in the regions and membrane rule collection. Formally, a cell-like P system (of degree $m \geq 1$) can be defined as form [23]: In the general model, the structure is in the form of nested membrane, which is not easy to be modularized, componentized and expanded. The rules in the cell-like P system can lead to high coupling degree.

$$\Pi = (O, \mu, \omega_1, \dots, \omega_m, R_1, \dots, R_m, i_o) \quad (1)$$

In the general model, the structure is in the form of nested membrane, which is not easy to be modularized, componentized and expanded. The rules in the cell-like P system can lead to high coupling degree.

B. The Related Works

The P module is a model that modularizes the biochemical reaction of a group of cells and supports information hiding. Formally, a P module can be defined as form:

$$\Pi = (O, K, \delta, P) \quad (2)$$

- 1) O is a finite non-empty alphabet of objects;
- 2) K is a finite set of cells, where each cell, $\delta \in K$, has the form $\delta = (Q, s_0, \omega_0, R)$ where,
 - a) Q is a finite set of states;
 - b) $s_0 \in Q$ is the initial state;
 - c) $\omega_0 \in O^*$ is the initial multiset of objects;
 - d) R is a finite ordered set of multiset rewriting rules of the general form:

$$sx \rightarrow_{\alpha} s'x'(u)_{\beta\gamma}; \quad (3)$$

where,

- (i) $s, s' \in Q$;
- (ii) $x, x' \in O^*, u \in O^*$;
- (iii) α is a rewriting operator, $\alpha \in \{min, max\}$, The rewriting operator $\alpha = min$ indicates that the rewriting is applied once, if the rule is applicable; and $\alpha = max$ indicates that the rewriting is applied as many times as possible, if the rule is applicable. When $\alpha = max$, α can be omitted in the rule.
- (iv) $\beta \in \{\uparrow, \downarrow, \updownarrow\}$;
- (v) $\gamma \in \{one, spread, repl\}$;
- 3) δ is a binary relation on K , i.e. a set of parent-child structural arcs, representing *duplex* or *simplex* communication channels between cells;
- 4) P is a subset of K , indicating the *port* cells, i.e. the only cells can be connected to other modules.

P module is a modular combination model of cells. It mainly uses the characteristics of its recursive combination to

realize the hidden functions of internal information and internal structure, so as to facilitate the construction of a complex P system.

III. DESIGN AND ASSEMBLY OF THE P MODULE

This section improves the P module with high encapsulation, information hiding, modular combination and high concurrency. Its special external definition, external reference and assembly mechanism make it highly independent, realize the reuse of modules and speed up the construction of P system.

A. P Module Improvement

The P module is a model of cell-like P system, which abstracts a cell into a P module. It has the characteristics of module encapsulation and the inheritance of rules and objects. Each module is independent and several P modules can be combined into the combination P module by a structured way. Formally, a P module (of degree $m \geq 1$) can be defined as form:

$$\Pi = (O, K, \delta, Q, D_{\uparrow}, D_{\downarrow}, R_{\uparrow}, R_{\downarrow}) \quad (4)$$

- 1) O is a finite non-empty alphabet of objects, $O=O_1 \cup O_2$. For each submodule, they contain the public objects from the parent module and their own objects.
 - a) O_1 is a subset of O , which represents private objects.
 - b) O_2 is a subset of O , disjoint of O_1 , which represents public objects.
- 2) K is a finite set of P modules.
- 3) δ is a subset of $(K \times K) \cup (K \times R_{\downarrow}) \cup (R_{\uparrow} \times K)$, i.e. a set of parent-child structural arcs, representing duplex or simplex communication channels, between two existing P modules or between an existing P modules and an external reference.
- 4) Q is a subset of O_2 , which is the generic synchronizing object set that P modules finally output.
- 5) D_{\uparrow} is a subset of K , representing *def_{\uparrow}* definitions, e.g. *def_{\uparrow\pi_i}* represents that the entrance module of this P module is Π_i ; D_{\downarrow} is a subset of K , representing *def_{\downarrow}* definitions, e.g. *def_{\downarrow\pi_i}* represents that the export module of this P module is Π_i .
- 6) R_{\uparrow} is a finite set, disjoint of K , representing *ref_{\uparrow}* references, e.g. *ref_{\uparrow a_i}* represents that the entrance arc of this P module is a_i ; R_{\downarrow} is a finite set, disjoint of K , representing *ref_{\downarrow}* references, e.g. *ref_{\uparrow b_i}* represents that the export arc of this P module is b_i .
- 7) Each cell, $\sigma \in K$, has the form $\sigma = (L, S, s_0, \omega_0, R)$. where,
 - a) L represents the inheritance rights of the rules; $L = \{\Gamma, \Delta, \Phi\}$; Γ represents this rule as a public rule, Δ represents this rule as a protected rule, Φ represents this rule as a private rule(this can be omitted).
 - b) S is a finite set of states;
 - c) $s_0 \in S$ is the initial state;
 - d) $\omega_0 \in O^*$ is the initial multiset of objects;
 - e) R is a finite ordered set of rules;

$$lsx \rightarrow_{\alpha} s'/x'; (id) \quad (5)$$

where,

- (i) $l \in L$;

- (ii) $s, s' \in S$;
- (iii) $x \in O^*$;
- (iv) α is a rewriting operator, $\alpha \in \{min, max\}$, The rewriting operator $\alpha = min$ indicates that the rewriting is applied once, if the rule is applicable; and $\alpha = max$ indicates that the rewriting is applied as many times as possible, if the rule is applicable. When $\alpha = max$, α can be omitted in the rule.
- (v) id is a number identified the sequence of rule execution. Prior and preference can be given high-ranking. If the rules in the same state are in the same priority, id can be omitted in the rule.

In this model, the connection relationship between P modules is a parent-child relationship, and their inheritance can be reflected by objects, O , and rules, R . Through the inheritance of the P module, we can organize system structure more effectively, clarify the relationship between modules, and make full use of existing modules to achieve more complex and deeper development.

B. P Module Assembly Mechanism

According to the definition of P module introduced above, a P system is a P module which is constructed by nested P modules. The nested P module is expressed by the combination P module which is constructed by P modules in the same layer. Given an arbitrary finite set of disjoint P modules, we can construct a combination P module by instantiating some of their external references to some of their external definitions, which implicitly instantiates the relationship of P modules in the same layer. When the parent P module is executed, the submodule will inherit the public objects and public rules of the parent P module to further initialize the internal structure of the module and start to work. The siblings can be executed in parallel, this shows the powerful computing power of the whole system. The combination P module can encapsulate the details of the interior, users only pay attention to their input and output.

Considering of a finite family of n P modules, $\Psi = \{\Pi_i | i \in [1, n]\}$, where $\Pi_i = (O_i, K_i, \delta_i, Q_i, D_{\uparrow_i}, D_{\downarrow_i}, R_{\uparrow_i}, R_{\downarrow_i}) (i \in [1, n])$, the result of a composition P module depends on one kind of actual instantiation that the external reference and the definition are matched. The external reference is matched to external definition by two partial mappings, $\rho_{\uparrow} : \cup_{i \in [1, n]} R_{\uparrow_i} \rightarrow \cup_{i \in [1, n]} D_{\uparrow_i}$, $\rho_{\downarrow} : \cup_{i \in [1, n]} R_{\downarrow_i} \rightarrow \cup_{i \in [1, n]} D_{\downarrow_i}$. A previously uninstantiated arc (σ, r) , where $(\sigma \in K_i, r \in R_{\downarrow_i} | i \in [1, n])$, is instantiated as $(\sigma, \rho_{\downarrow(\sigma)})$, and a previously uninstantiated arc (r, σ) , where $(\sigma \in K_i, r \in R_{\uparrow_i} | i \in [1, n])$, is instantiated as $(\rho_{\uparrow(\sigma)}, \sigma)$.

Based on what has been described above, the P module family Ψ can be expressed as the form, $\Pi = (O, K, \delta, Q, D_{\uparrow}, D_{\downarrow}, R_{\uparrow}, R_{\downarrow})$, when $\rho_{\uparrow}, \rho_{\downarrow}$ are the partial mappings that define the instantiation (as previously introduced), if:

- 1) Ψ is cell-disjoint;
- 2) $O = \cup_{i \in [1, n]} O_i$;
- 3) $K = \cup_{i \in [1, n]} K_i$;
- 4) $\delta = \{(\rho_{\uparrow(\sigma)}, \tilde{\rho}_{\downarrow(\sigma)}) | \cup_{i \in [1, n]} \sigma_i\}$, where $\tilde{\rho}_{\uparrow(\sigma)} = \sigma \in Dom(\rho_{\uparrow})? \rho_{\uparrow(\sigma)} : \sigma, \tilde{\rho}_{\downarrow(\sigma)} = \sigma \in Dom(\rho_{\downarrow})? \rho_{\downarrow(\sigma)} : \sigma$;

- 5) $Q = \cup_{i \in [1, n]} Q_i / \cup_{i \in [1, n-1]} Q_i = Q_n; (Q_n \text{ is the output objects as the exit of the combination P module})$
- 6) $D_{\uparrow} \subseteq \cup_{i \in [1, n]} D_{\uparrow_i}, D_{\downarrow} \subseteq \cup_{i \in [1, n]} D_{\downarrow_i}$;
- 7) $R_{\uparrow} = \cup_{i \in [1, n]} R_{\uparrow_i} \setminus Dom(\rho_{\uparrow}), R_{\downarrow} = \cup_{i \in [1, n]} R_{\downarrow_i} \setminus Dom(\rho_{\downarrow})$;

As described above, we can know the concrete the construction and assembly mechanism of P modules in P system, which includes the nested combination principle of the P module in the same layer and the perfect encapsulation mechanism. The construction and assembly mechanism also make a detailed definition of δ, D, R as a way of communication. The P modules construction and assembly mechanism facilitates the design of P system for complex algorithms, where every P module provides encapsulation and information hiding to other P modules.

IV. BASIC STRUCTURE OF THE P SYSTEM

The P system is constructed by layer upon layer encapsulation using P module. P modules in the same layer are assembled by the construction and assembly mechanism and encapsulated into a combination P module. The construction and assembly methods of a combination P module include four ways, i.e. the sequence method, the branch method, the cycle method and the parallel method, which show the four structures of the P module, respectively.

A. Sequential Method

Since the specific implementation rules of each P module will be determined by the function to be performed, the definition of the rules in each P module need to be abstracted to be a form, which is shown below through the two rules. The general design definition of a elementary P module perform a series of calculations on the initial object set, and finally output the result set. As shown below, here are two rules to represent this process, r_1 represents a series of operations on the initial object set, which are a series of operations except the output of the result set, and the evolution from the initial set of objects x to y is accomplished by multiple rules in the specific implementation. r_2 represents the calculated set of objects is evolved into the set of public objects required by the submodule, so that the submodule can inherit from it to obtain a complete initial set of objects.

$$r_1: \Gamma / \Delta / \Phi S_0 x_0, \dots, x_{i_0} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_0}$$

$$r_2: \Gamma / \Delta / \Phi S_1 y_0, \dots, y_{j_0} \rightarrow_{min/max} S_0 z_0, \dots, z_{k_0}$$

Fig. 2 illustrates a combined P module through the sequential modular composition of two elementary P modules.

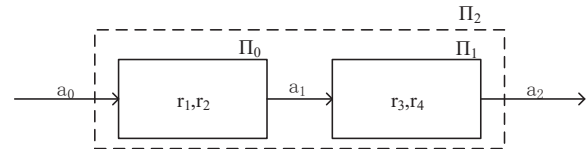


Fig. 2. The sequential structure based on the P module.

In $\Pi_0 < def_{\downarrow \Pi_0}, ref_{\downarrow a_1} >$, using the ruleset following this paragraph objects $\alpha_i (i \in [1, k_0])$ can be obtained by inputting objects, $x_i (i \in [1, i_0])$.

$$r_1: \Gamma / \Delta / \Phi S_0 x_0, \dots, x_{i_0} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_0}$$

$$r_2: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_0} \rightarrow_{min/max} S_0 \alpha_0, \dots, \alpha_{k_0}$$

In $\Pi_1 < def_{\downarrow \Pi_1}, ref_{\downarrow a_2} >$ following this paragraph, using the ruleset following this paragraph, objects $\beta_i (i \in [1, k_1])$ can be obtained by inputting objects, $\alpha_i (i \in [1, k_0])$.

$$r_3: \Gamma/\Delta/\Phi S_0 \alpha_0, \dots, \alpha_{k_0} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_1}$$

$$r_4: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} \rightarrow_{min/max} S_0 \beta_0, \dots, \beta_{k_1}$$

The combination P module, $\Pi_2 < def_{\downarrow \Pi_0}, ref_{\downarrow a_2} >$, contains two P modules, Π_0 and Π_1 , which also appears as an external def_{\downarrow} definition, and makes external ref_{\downarrow} references to a unspecified P module, a_2 . There is a definition of Π_2 following this paragraph.

$$\Pi_2 = (O, K, \delta, Q, D_{\uparrow}, D_{\downarrow}, R_{\uparrow}, R_{\downarrow})$$

where,

$$1) O = O_1 \cup O_2$$

$$\begin{aligned} &= \{ x_0, \dots, x_{i_0}, y_0, \dots, y_{j_0}, x_0, \dots, x_{i_1}, y_0, \dots, y_{j_1} \\ &\quad \} \cup \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1} \} \\ &= \{ x_0, \dots, x_{i_0}, y_0, \dots, y_{j_0}, x_0, \dots, x_{i_1}, y_0, \dots, y_{j_1}, \\ &\quad \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1} \} \end{aligned}$$

$$2) K = \{ \Pi_0, \Pi_1 \}$$

$$3) \delta = \{ (\Pi_0, ref_{\downarrow a_1} \rightarrow def_{\downarrow \Pi_1}) \}$$

$$4) Q = \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1} \} / \{ \alpha_0, \dots, \alpha_{k_0} \} = \{ \beta_0, \dots, \beta_{k_1} \}$$

$$5) D_{\downarrow} = \{ def_{\downarrow \Pi_0} \}, D_{\uparrow} = \{ \}$$

$$6) R_{\downarrow} = \{ ref_{\downarrow a_2} \}, R_{\uparrow} = \{ \}$$

We can connect Π_0 and Π_1 by the generic instantiation: $(\Pi_0, ref_{\downarrow a_1} \rightarrow def_{\downarrow \Pi_1})$.

B. Branch Method

Fig. 3 illustrates a combined P module through the branch modular composition of four P modules.

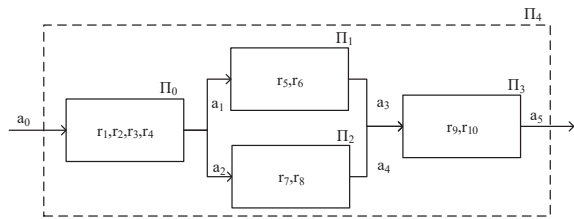


Fig. 3. The branch structure based on the P module.

In $\Pi_0 < def_{\downarrow \Pi_0}, ref_{\downarrow a_5} >$, using the ruleset following this paragraph, there is no material input to determine the object m (m is a set of objects), if there is a material object, get $\alpha_i (i \in [1, k_0])$; otherwise, get $\beta_i (i \in [1, k_0])$.

$$r_1: \Gamma/\Delta/\Phi S_0 x_0, \dots, x_{i_0}, m \rightarrow_{min/max} S_1 y_0, \dots, y_{j_0}; 1$$

$$r_2: \Gamma/\Delta/\Phi S_0 x_0, \dots, x_{i_0}, \rightarrow_{min/max} S_1 y_0, \dots, y_{j_0}; 2$$

$$r_3: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_0} \rightarrow_{min/max} S_0 \alpha_0, \dots, \alpha_{k_0}$$

$$r_4: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_0} \rightarrow_{min/max} S_0 \beta_0, \dots, \beta_{k_0}$$

In $\Pi_1 < def_{\downarrow \Pi_1}, ref_{\downarrow a_3} >$, using the ruleset following this paragraph, objects $\phi_i (i \in [1, k_1])$ can be obtained by inputting objects, $\alpha_i (i \in [1, k_0])$.

$$r_5: \Gamma/\Delta/\Phi S_0 \alpha_0, \dots, \alpha_{k_0} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_1}$$

$$r_6: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} \rightarrow_{min/max} S_0 \phi_0, \dots, \phi_{k_1}$$

In $\Pi_2 < def_{\downarrow \Pi_2}, ref_{\downarrow a_4} >$, using the ruleset following this paragraph, objects $\phi_i (i \in [1, k_1])$ can be obtained by inputting objects, $\beta_i (i \in [1, k_0])$.

$$r_7: \Gamma/\Delta/\Phi S_0 \beta_0, \dots, \beta_{k_0} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_2}$$

$$r_8: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_2} \rightarrow_{min/max} S_0 \phi_0, \dots, \phi_{k_1}$$

In $\Pi_3 < def_{\downarrow \Pi_3}, ref_{\downarrow a_5} >$, using the ruleset following this paragraph, objects $\gamma_i (i \in [1, k_2])$ can be obtained by inputting objects, $\phi_i (i \in [1, k_1])$.

$$r_7: \Gamma/\Delta/\Phi S_0 \phi_0, \dots, \phi_{k_1} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_2}$$

$$r_8: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_2} \rightarrow_{min/max} S_0 \gamma_0, \dots, \gamma_{k_2}$$

The combined P module, $\Pi_4 < def_{\downarrow \Pi_0}, ref_{\downarrow a_5} >$, contains four P modules, Π_0, Π_1, Π_2 and Π_3 , which also appears as an external def_{\downarrow} definition, and makes one external ref_{\downarrow} references to one unspecified P module, a_5 . There is a definition of Π_4 following this paragraph.

$$\Pi_4 = (O, K, \delta, Q, D_{\uparrow}, D_{\downarrow}, R_{\uparrow}, R_{\downarrow})$$

where,

$$1) O = O_1 \cup O_2$$

$$\begin{aligned} &= \{ x_0, \dots, x_{i_0}, y_0, \dots, y_{j_0}, m, x_0, \dots, x_{i_1}, y_0, \\ &\quad \dots, y_{j_1} \} \cup \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_0}, \phi_0, \dots, \\ &\quad \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \} \\ &= \{ x_0, \dots, x_{i_0}, y_0, \dots, y_{j_0}, m, x_0, \dots, x_{i_1}, y_0, \\ &\quad \dots, y_{j_1}, \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_0}, \phi_0, \dots, \phi_{k_1}, \\ &\quad \gamma_0, \dots, \gamma_{k_2} \} \end{aligned}$$

$$2) K = \{ \Pi_0, \Pi_1, \Pi_2, \Pi_3 \}$$

$$3) \delta = \{ (\Pi_0, ref_{\downarrow a_1} \rightarrow def_{\downarrow \Pi_1}), (\Pi_0, ref_{\downarrow a_2} \rightarrow def_{\downarrow \Pi_2}), (\Pi_1, ref_{\downarrow a_3} \rightarrow def_{\downarrow \Pi_3}), (\Pi_2, ref_{\downarrow a_4} \rightarrow def_{\downarrow \Pi_3}) \}$$

$$4) Q = \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_0}, \phi_0, \dots, \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \} / \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_0}, \phi_0, \dots, \phi_{k_1} \} = \{ \gamma_0, \dots, \gamma_{k_2} \}$$

$$5) D_{\downarrow} = \{ def_{\downarrow \Pi_0} \}, D_{\uparrow} = \{ \}$$

$$6) R_{\downarrow} = \{ ref_{\downarrow a_5} \}, R_{\uparrow} = \{ \}$$

We can connect Π_0 and Π_1 by the generic instantiation: $(\Pi_0, ref_{\downarrow a_1} \rightarrow def_{\downarrow \Pi_1})$, Π_0 and Π_2 by the generic instantiation: $(\Pi_0, ref_{\downarrow a_2} \rightarrow def_{\downarrow \Pi_2})$, Π_1 and Π_3 by the generic instantiation: $(\Pi_1, ref_{\downarrow a_3} \rightarrow def_{\downarrow \Pi_3})$ and connect Π_2 and Π_3 by the generic instantiation: $(\Pi_2, ref_{\downarrow a_4} \rightarrow def_{\downarrow \Pi_3})$.

C. Cycle Method

Fig. 4 illustrates a combined P module through the cycle modular composition of two P modules which construct do-while model.

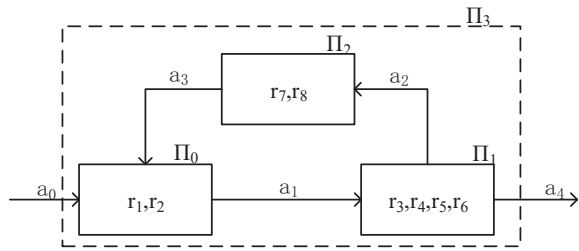


Fig. 4. The cycle structure based on the P module.

In $\Pi_0 < def_{\downarrow\Pi_0}, ref_{\downarrow a_1} >$, using the ruleset following this paragraph, objects $\alpha_i (i \in [1, k_0])$ can be obtained by inputting objects, $x_i (i \in [1, i_0])$.

$$r_1: \Gamma/\Delta/\Phi S_0 x_0, \dots, x_{i_0} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_0}$$

$$r_2: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} \rightarrow_{min/max} S_0 \alpha_0, \dots, \alpha_{k_0}$$

In $\Pi_1 < def_{\downarrow\Pi_1}, ref_{\downarrow a_2}, ref_{\downarrow a_4} >$, there is no material input to determine the object m (m is a set of objects), if there is a material object, get $\beta_i (i \in [1, k_1])$; otherwise, get $\phi_i (i \in [1, k_0])$.

$$r_3: \Gamma/\Delta/\Phi S_0 \alpha_0, \dots, \alpha_{k_0}, m \rightarrow_{min/max} S_1 y_0, \dots, y_{j_1}; 1$$

$$r_4: \Gamma/\Delta/\Phi S_0 \alpha_0, \dots, \alpha_{k_0}, \rightarrow_{min/max} S_1 y_0, \dots, y_{j_1}; 2$$

$$r_5: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} \rightarrow_{min/max} S_0 \beta_0, \dots, \beta_{k_1}$$

$$r_6: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} \rightarrow_{min/max} S_0 \phi_0, \dots, \phi_{k_1}$$

In $\Pi_2 < def_{\downarrow\Pi_2}, ref_{\downarrow a_3} >$, using the ruleset following this paragraph, objects $\gamma_i (i \in [1, k_2])$ can be obtained by inputting objects, $\beta_i (i \in [1, k_1])$.

$$r_7: \Gamma/\Delta/\Phi S_0 \beta_0, \dots, \beta_{k_1} \rightarrow_{min/max} S_1 y_0, \dots, y_{j_2}$$

$$r_8: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_2} \rightarrow_{min/max} S_0 \gamma_0, \dots, \gamma_{k_2}$$

The combined P module, $\Pi_3 < def_{\downarrow\Pi_0}, ref_{\downarrow a_4} >$, contains three P modules, Π_0 , Π_1 and Π_2 , which also appears as an external def_{\downarrow} definition, and makes two external ref_{\downarrow} references to one unspecified P modules, a_4 . There is a definition of Π_3 following this paragraph.

$$\Pi_3 = (O, K, \delta, Q, D_{\uparrow}, D_{\downarrow}, R_{\uparrow}, R_{\downarrow})$$

where,

- 1) $O = O_1 \cup O_2$

$$= \{ x_0, \dots, x_{i_0}, m, y_0, \dots, y_{j_0}, y_0, \dots, y_{j_1}, y_0, \dots, y_{j_2} \} \cup \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1}, \phi_0, \dots, \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \}$$

$$= \{ x_0, \dots, x_{i_0}, m, y_0, \dots, y_{j_0}, y_0, \dots, y_{j_1}, y_0, \dots, y_{j_2}, \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1}, \phi_0, \dots, \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \}$$
- 2) $O = O_1 \cup O_2 = \{ x_0, \dots, x_{i_0}, m, y_0, \dots, y_{j_0}, y_0, \dots, y_{j_1}, y_0, \dots, y_{j_2} \} \cup \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1}, \phi_0, \dots, \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \} = \{ x_0, \dots, x_{i_0}, m, y_0, \dots, y_{j_0}, y_0, \dots, y_{j_1}, y_0, \dots, y_{j_2}, \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1}, \phi_0, \dots, \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \}$
- 3) $K = \{ \Pi_0, \Pi_1, \Pi_2 \}$

- 4) $\delta = \{ (\Pi_0, ref_{\downarrow a_1} \rightarrow def_{\downarrow\Pi_1}), (\Pi_1, ref_{\downarrow a_2} \rightarrow def_{\downarrow\Pi_2}), (\Pi_2, ref_{\downarrow a_3} \rightarrow def_{\downarrow\Pi_0}) \}$
- 5) $Q = \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1}, \phi_0, \dots, \phi_{k_1}, \gamma_0, \dots, \gamma_{k_2} \} / \{ \alpha_0, \dots, \alpha_{k_0}, \beta_0, \dots, \beta_{k_1}, \phi_0, \dots, \phi_{k_1} \} = \{ \gamma_0, \dots, \gamma_{k_2} \}$
- 6) $D_{\downarrow} = \{ def_{\downarrow\Pi_0} \}, D_{\uparrow} = \{ \}$
- 7) $R_{\downarrow} = \{ ref_{\downarrow a_3}, ref_{\downarrow a_4} \}, R_{\uparrow} = \{ \}$

We can connect Π_0 and Π_1 by the generic instantiation: $(\Pi_0, ref_{\downarrow a_1} \rightarrow def_{\downarrow\Pi_1})$ and connect Π_1 and Π_2 by the generic instantiation: $(\Pi_1, ref_{\downarrow a_2} \rightarrow def_{\downarrow\Pi_2})$ and connect Π_2 and Π_0 by the generic instantiation: $(\Pi_2, ref_{\downarrow a_3} \rightarrow def_{\downarrow\Pi_0})$.

D. Parallel Method

The parallel structure can be seen as a variant of the branch structure, but the operation rules in this structure are very different from the branch structure due to the large number of P modules in parallel computation and its unique parallelism. Fig. 5 shows the generic parallel structure of parallel computing modules of degree m (m is a variable).

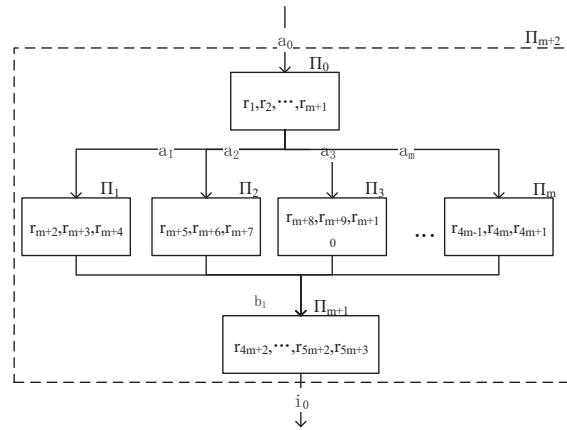


Fig. 5. The parallel structure based on the P module.

The parallelism of the parallel structure is embodied in the $1 - m$ P module. They equally inherit the public object set from Π_0 , and Π_{m+1} inherits the object set of m parallel P modules.

For the parallel structure, each parallel module first inherits the same object set from Π_0 , then uses some of the inheriting object sets to execute the respective calculation rules, and finally outputs the result set and remaining inherited objects to Π_{m+1} . Π_{m+1} inherits the object set of all parallel P modules, Π_{m+1} module needs to process these inheriting objects in order to run correctly.

Due to the existence of multiple inheritance and the existence of a special case of 1 pair n and n pair 1, it becomes more complicated to maintain the consistency of the data. Some inheriting object set still remain the submodule, because the parallel execution P modules use only part of the inheriting object set. In the remaining inherited objects, one is the object set as the public global variables that need pass to the child module, and the other is the redundant object set. For the first case, the submodules of the parallel modules inherits the object sets of the m P modules and results in a multiple of the

number of object sets. Therefore, it is necessary to divide the public global variables in the submodule (i.e. the public global variables/the number of the parallel P modules). In the second case, since it is useless data, each parallel submodule needs to destroy the redundant object set (because the redundant object set does not have a general purpose, so specific problems need to be specifically designed). To make the above situation clear, set the public global variable to be ω for the rule design. The following is the general structure design of the parallel structure.

In $\Pi_0 < def_{\downarrow \Pi_0}, ref_{\downarrow a_1}, ref_{\downarrow a_2}, ref_{\downarrow a_3}, \dots, ref_{\downarrow a_m} >$, the original objects include x_0, \dots, x_{i_0} and ω (it is the public global variable), using the ruleset following this paragraph, objects $\alpha_i, \beta_i, \dots, \phi_i (i \in [1, k_0])$ can be obtained by inputting objects, $x_i (i \in [1, i_0])$.

$$\begin{aligned} r_1: \Gamma/\Delta/\Phi S_0 x_0, \dots, x_{i_0} &\rightarrow_{min/max} S_1 y_0, \dots, y_{j_0} \\ r_2: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_0} &\rightarrow_{min/max} S_0 \alpha_0, \dots, \alpha_{k_0} \\ r_3: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_0} &\rightarrow_{min/max} S_0 \beta_0, \dots, \beta_{k_0} \\ \dots \\ r_{m+1}: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_0} &\rightarrow_{min/max} S_0 \phi_0, \dots, \phi_{k_0} \end{aligned}$$

In $\Pi_1 < def_{\downarrow \Pi_1}, ref_{\downarrow b_1} >$, using the ruleset following this paragraph, objects $\alpha'_i (i \in [1, p_0])$ can be obtained by inputting objects, $\alpha_i (i \in [1, k_0])$. For the remaining objects, i.e. $(\beta_0, \dots, \beta_{k_0})$, they need to be removed by the rule named r_{m+4} .

$$\begin{aligned} r_{m+2}: \Gamma/\Delta/\Phi S_0 \alpha_0, \dots, \alpha_{k_0} &\rightarrow_{min/max} S_1 y_0, \dots, y_{j_1} \\ r_{m+3}: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} &\rightarrow_{min/max} S_0 \alpha'_0, \dots, \alpha'_{p_0} \\ r_{m+4}: \Gamma/\Delta/\Phi S_1 \beta_0, \dots, \beta_{k_0}, \dots, \phi_0, \dots, \phi_{k_0} &\rightarrow_{max} S_0 \end{aligned}$$

In $\Pi_2 < def_{\downarrow \Pi_2}, ref_{\downarrow b_2} >$, using the ruleset following this paragraph, objects $\beta'_i (i \in [1, p_0])$ can be obtained by inputting objects, $\beta_i (i \in [1, k_0])$.

$$\begin{aligned} r_{m+5}: \Gamma/\Delta/\Phi S_0 \beta_0, \dots, \beta_{k_0} &\rightarrow_{min/max} S_1 y_0, \dots, y_{j_1} \\ r_{m+6}: \Gamma/\Delta/\Phi S_1 y_0, \dots, y_{j_1} &\rightarrow_{min/max} S_0 \beta'_0, \dots, \beta'_{p_0} \\ r_{m+7}: \Gamma/\Delta/\Phi S_1 \alpha_0, \dots, \alpha_{k_0}, \dots, \phi_0, \dots, \phi_{k_0} &\rightarrow_{max} S_0 \end{aligned}$$

Due to the uncertainty in the number of the parallel P modules, we will not list rules of other parallel P modules here. Their difference is that the use of different initialization object sets, the output set of different object sets and the remaining object sets needed to be removed.

$\Pi_{m+1} < def_{\downarrow \Pi_{m+1}}, ref_{\downarrow i_0} >$ can inherit the same object set, ω , from n P modules and result in the error of ω , so need to get the correct object set, ω , by $\omega = \omega/m$ (rule: r_{4m+2}). Then perform corresponding calculations on different kinds of object sets. Its rule set is as listed below.

$$\begin{aligned} r_{4m+2}: \Gamma/\Delta/\Phi S_0 \omega^m &\rightarrow_{max} S_1 \omega; \\ r_{4m+3}: \Gamma/\Delta/\Phi S_0 \alpha'_0, \dots, \alpha'_{p_0} &\rightarrow_{min/max} S_1 y_0, \dots, y_{j_2}; \\ r_{4m+4}: \Gamma/\Delta/\Phi S_0 \beta'_0, \dots, \beta'_{p_0} &\rightarrow_{min/max} S_1 y_0, \dots, y_{j_2}; \\ \dots \\ r_{5m+2}: \Gamma/\Delta/\Phi S_0 \phi'_0, \dots, \phi'_{p_0} &\rightarrow_{min/max} S_1 y_0, \dots, y_{j_2}; \\ r_{5m+3}: \Gamma/\Delta/\Phi S_0 y_0, \dots, y_{j_2} &\rightarrow_{min/max} S_1 \gamma_0, \dots, \gamma_{k_1}; \end{aligned}$$

The combined P module, $\Pi_{m+2} < def_{\downarrow \Pi_0}, ref_{\downarrow i_0} >$, contains three kinds of P modules, i.e. the control P module, the parallel P modules and the merge P module. The control P module is Π_0 , the parallel P module are $\Pi_i (i \in [1, m])$, the merge P module is Π_{m+1} . There is a definition of Π_{m+2} following this paragraph.

$$\Pi_{m+2} = (O, K, \delta, Q, D_{\uparrow}, D_{\downarrow}, R_{\uparrow}, R_{\downarrow})$$

where,

- 1) $O = O_1 \cup O_2 = \{ x_0, \dots, x_{i_0}, y_0, \dots, y_{j_0}, y_0, \dots, y_{j_1}, y_0, \dots, y_{j_2} \} \cup \{ \omega, (\alpha_0, \dots, \alpha_{k_0}), (\beta_0, \dots, \beta_{k_1}), \dots, (\phi_0, \dots, \phi_{k_1}), (\gamma_0, \dots, \gamma_{k_2}), (\alpha'_0, \dots, \alpha'_{p_0}), (\beta'_0, \dots, \beta'_{p_0}), \dots, (\phi'_0, \dots, \phi'_{p_0}) \} = \{ x_0, \dots, x_{i_0}, y_0, \dots, y_{j_0}, y_0, \dots, y_{j_1}, y_0, \dots, y_{j_2}, (\alpha_0, \dots, \alpha_{k_0}), (\beta_0, \dots, \beta_{k_1}), \dots, (\phi_0, \dots, \phi_{k_1}), (\gamma_0, \dots, \gamma_{k_2}), (\alpha'_0, \dots, \alpha'_{p_0}), (\beta'_0, \dots, \beta'_{p_0}), \dots, (\phi'_0, \dots, \phi'_{p_0}), \omega \}$
- 2) $K = \{ \Pi_i (i \in [1, m]) \}$
- 3) $\delta = \{ (\Pi_0, ref_{\downarrow a_i} \rightarrow def_{\downarrow \Pi_i}) (i \in [1, m]), (\Pi_i, ref_{\downarrow b_i} \rightarrow def_{\downarrow \Pi_{m+1}}) (i \in [1, m]) \}$
- 4) $Q = \{ \omega, (\alpha_0, \dots, \alpha_{k_0}), (\beta_0, \dots, \beta_{k_1}), \dots, (\phi_0, \dots, \phi_{k_1}), (\gamma_0, \dots, \gamma_{k_2}), (\alpha'_0, \dots, \alpha'_{p_0}), (\beta'_0, \dots, \beta'_{p_0}), \dots, (\phi'_0, \dots, \phi'_{p_0}) \} / \{ (\alpha_0, \dots, \alpha_{k_0}), (\beta_0, \dots, \beta_{k_1}), \dots, (\phi_0, \dots, \phi_{k_1}), (\alpha'_0, \dots, \alpha'_{p_0}), (\beta'_0, \dots, \beta'_{p_0}), \dots, (\phi'_0, \dots, \phi'_{p_0}) \} = \{ \omega, \gamma_0, \dots, \gamma_{k_2} \}$
- 5) $D_{\downarrow} = \{ def_{\downarrow \Pi_0} \}, D_{\uparrow} = \{ \}$
- 6) $R_{\downarrow} = \{ ref_{\downarrow i_0} \}, R_{\uparrow} = \{ \}$

We can connect Π_0 and the parallel P modules, $\Pi_i (i \in [1, m])$ by the generic instantiation: $(\Pi_0, ref_{\downarrow a_i} \rightarrow def_{\downarrow \Pi_i}) (i \in [1, m])$ and connect $\Pi_i (i \in [1, m])$ and Π_{m+1} by the generic instantiation: $(\Pi_i, ref_{\downarrow b_i} \rightarrow def_{\downarrow \Pi_{m+1}}) (i \in [1, m])$.

V. AN EXAMPLE: P SYSTEM DESIGN BASED ON P MODULE

Based on the high computational complexity of calculating the arithmetic square root of a large number, this section proposes an efficient algorithm to reduce its computational complexity, and implement the algorithm in the P system by using the P module and four P module construction methods.

A. Square Root Algorithm of a Large Number

Two algorithms for calculating the square root of a large number are introduced here. One is a square root estimation algorithm for estimating the scope of the square root. This algorithm is illustrated by Table I.

The algorithm is applied to estimate the square root, including four steps and the time complexity of $Sqrte(n)$ is about $O(d)$ (d is the number of digits of n).

The other is a square root algorithm through m bisection calculation algorithm. This algorithm is an improved algorithm of 2-points, which is illustrated by Table II.

According Table II, the complexity of $Mbisection(b, n, a, m)$ is about $O(\log_m n)$ (m is the number of the interval splitted, n is the interval size).

TABLE I. ALGORITHM: *EstimateSqr(x)*

Input: x
Output: the left interval point of estimated square root and the interval size of estimated square root
Steps:
(1) Circularly dividing 100 with no remainder, then get the high-value of the inputting number expressed as the numbers, y , and the number of cycles.
(2) Select the square root of the high-value according to the formula, i.e. the high-value $> [81, 64, 49, 36, 25, 16, 9, 4, 1]$, then the high-value square root result is $[10, 9, 8, 7, 6, 5, 4, 3, 2]$.
(3) Combine the square root of the high-value and the cycles.
(4) Output: the left interval point of estimated square root and the size of estimated square root.
End

TABLE II. *Mbisection(b, l, x, m)*

Input: b, l, x
x : the original number to be squared;
b : the left interval point of \sqrt{x} ;
l : the size of the interval of \sqrt{x} ;
m : the quantity of equidistant intervals
original interval: $[b, b + l]$
Output: the square root
Steps:
(1) Split the interval into m equidistant intervals and obtain $m + 1$ copies of endpoint number in the interval, x_0, x_1, \dots, x_m .
(2) Parallely calculate $f(x_i) = x_i^2 - x$ ($i \in [0, m]$), then output these key value pairs $(x_i, f(x_i))$ ($i \in [0, m]$).
(3) Filter these key value pairs, if there is the number, $f(x_i) = 0$, then output x_i as the final result; otherwise, output two adjacent numbers, x_i, x_{i+1} when $f(x_i) < 0$ and $f(x_{i+1}) > 0$.
(4) If $x_{i+1} - x_i > 2, x_{i+1}$ and x_i are as a new round of inputting and enter (1); otherwise, output $x_i + 1$ as the final result.
End

B. Square Root Algorithm of a Large Number based on the P Module

The previous chapter proposed two algorithms for solving the square root of a large number, but the computational efficiency of each algorithm is not optimized. To reduce the computational complexity, the two algorithms can be combined to form an efficient algorithm named *Bigsplite*, which integrates the square root estimation algorithm, *EstimateSqr(x)*, and the square root algorithm through m bisection calculation algorithm, *Mbisection(b, l, x, m)*. By using P modules, the P system, *Bigsplite*, has high powerful parallel execution capabilities, high reusability and low coupling.

Provided that input α , the square root algorithm of a large number based on the P module, *Bigsplite*(α), is illuminated by Fig. 6. While Fig. 6 shows that how to construct the P system by using P modules. In the beginning of the system construction, the initial objects structure of the elementary modules only contains one object, c , except for $\Pi_1 \dots \Pi_9$. The objects that corresponds to $\Pi_1 \dots \Pi_9$ are $b^4 c^4 q^2 s, b^9 c^5 q^3 s, b^{16} c^7 q^4 s, b^{25} c^9 q^5 s, b^{36} c^{11} q^6 s, b^{49} c^{13} q^7 s, b^{64} c^{15} q^8 s, b^{81} c^{17} q^9 s, b^{100} c^{19} q^{10} s$. The specific rules of this algorithm in the P system are shown in the appendix. According to the flow chart and the rule sets, the detailed implementation process is described in detail below.

- 1) Firstly, copy α to τ for saving global parameters, α . Through circularly dividing the data, τ , by 100, the quotient of the cycle calculation value- δ and the number

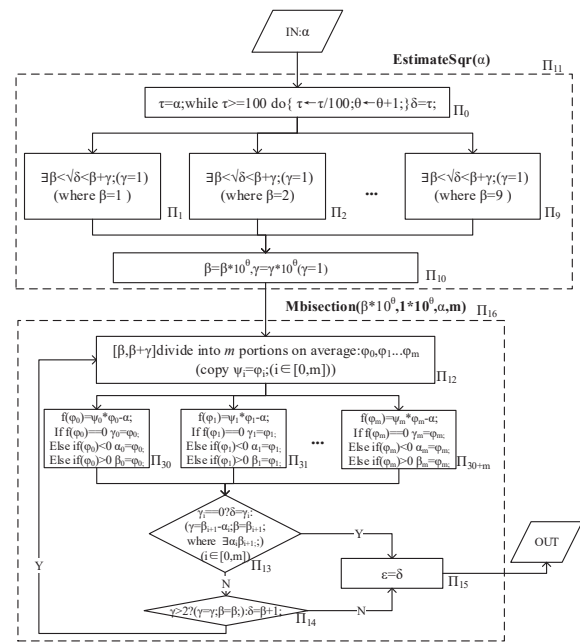


Fig. 6. The square root algorithm of large number based on the P module.

- of cycle calculations- θ are obtained after the cycle ends. At present, there are public objects in Π_0 , i.e. the original objects, α , the highest-segment value of α , δ , and the bits of $\sqrt{\alpha}$ except the highest-bit, θ . Go to the step 2.
- 2) In Π_1, \dots, Π_9 , these P modules parallel compute the square root of the high-segment value, δ , inheriting from Π_0 . Through finding a P module of Π_1, \dots, Π_9 that make the value, q , satisfying the arithmetic formula: $\sqrt{\delta} < q < \sqrt{\delta} + 1$ and converting q to be $\beta + \gamma$ ($\gamma = 1$), $\sqrt{\delta}$ can be expressed by $\beta < \sqrt{\delta} < \beta + \gamma$ ($\gamma = 1$). For other P modules, they need to remove inherited public objects, α and θ . Finally, there are public objects in Π_i (Π_i is the selected cell), i.e. the original objects, α , the bits of $\sqrt{\alpha}$ except the highest-bit, θ , the left interval point of $\sqrt{\delta}$, β , and the size of $\sqrt{\delta}$, γ . Go to step (3). (Parallel Computing)
- 3) In Π_{10} , combine of digits, θ , and the valuation interval of the square root of high-segment value, $[\beta, \beta + \gamma]$, then get the valuation range of α , i.e. $[\beta * 10^\theta, (\beta + \gamma) * 10^\theta]$ which is replaced by $[\beta, \beta + \gamma]$ ($\gamma = 1$). Finally, public objects have the original objects, α , the left valuation interval point, β , and the valuation interval size, γ . Go to step (4).
- 4) In Π_{12} , the valuation interval, $[\beta, \beta + \gamma]$ ($\gamma = 1$), is divided into m intervals on average, which are showed as $m + 1$ endpoints, $\varphi_0, \varphi_1, \dots, \varphi_m$ (they are also $\psi_0, \psi_1, \dots, \psi_m$). Public objects include the original objects, $\alpha, \varphi_0, \varphi_1, \dots, \varphi_m$ and $\psi_0, \psi_1, \dots, \psi_m$. Go to step (5).
- 5) In Π_{30+i} ($i \in [0, m]$), they respectively calculate $f(\varphi_i, \psi_i) = \varphi_i * \psi_i - \alpha$ by the objects they want to use, i.e. φ_i, ψ_i (i = the id of P module-30), and respectively destroy other objects, i.e. φ_i, ψ_i ($i \neq$ the id of cell-30). If $f(\varphi_i, \psi_i) = 0$, producing γ_i ($\gamma_i = \varphi_i$); else if $f(\varphi_i, \psi_i) < 0$, producing α_i ($\alpha_i = \varphi_i$); else if $f(\varphi_i, \psi_i) > 0$, producing β_i ($\beta_i = \varphi_i$). There must be

two kinds of objects for each module, one is one of the $\alpha_i, \beta_i, \gamma_i$ (i =the id of cell-30), the second is that each module has a common global variable, α . Go to step (6).(Parallel Computing)

- 6) In Π_{13} , determine whether there is an accurate result from the results of the calculation, γ_i , if so, then the square root value, $\delta = \gamma_i$, inherited by Π_{15} and go to step(7); otherwise, find the interval of the square root, i.e. $[\beta, \beta + \gamma]=[\alpha_i, \beta_{i+1}]$. At present, there is either δ and α or β, γ and α . Go to step (8).
- 7) In Π_{15} ,output the final result $\varepsilon(\varepsilon = \delta)$.
- 8) In the cell Π_{14} , determine whether the size of the interval, i.e. γ , is more than 2, if so, the result, $\delta = \beta + 1$ is given to Π_{15} and go to step (7); otherwise,pass β, γ, α to Π_{12} and go to step (4).

The square root algorithm of a large number based on the P module, *Bigsplite*(α), can pass through two phases, *EstimateSqr*(α) and *Mbisection*(β, γ, α, m), when the result interval of the square root is $[\beta, \beta + \gamma]$ by *EstimateSqr*(α). The size of the second parameter of the algorithm named *Mbisection* is much smaller than the original input number, α , so the total number of the recursive computing is decreased, but the time complexity of *Mbisection*(β, γ, α, m) is still $O(\log_m n)$. So the time complexity of the square root algorithm of large number based on the P module is $O(\log_m n)$.

C. Structure of P Module for Calculating the Square Root of a Large Number

Fig. 7 is a P module flow chart illustrating the modular combination of P modules that calculates the square root of a large number.

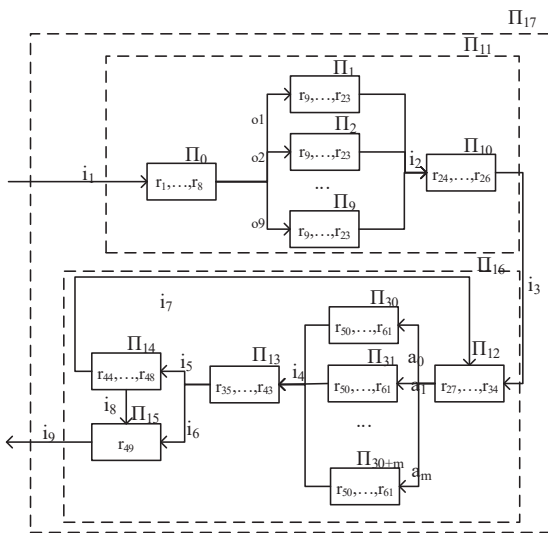


Fig. 7. A P module flow of a parallel algorithm.

As is shown Fig. 7, the number of P modules in the P system are $(19 + m)$. The P system can be expressed as the biggest combination P module $\Pi_{17} < def_{\downarrow \Pi_0}, ref_{\downarrow i_9} >$. The entire definition and operation mechanism of the system is mainly represented by public object sets, rule sets, and construction and assembly mechanism. The description of the

public object set can not only show the evolutionary direction and flow of the entire system calculation, but also more easily incorporate the rule set in the appendix to ensure the correctness of the rule design of P system. The following is a simple description of the public and private objects in the evolution of the rules of each module.

- 1) In Π_0 , $O_1=\{ \tau, b, c \}$ and $O_2=\{ \alpha, \theta, \delta \}$;
- 2) In the parallel P modules of Π_{10} , $O_1=\{ b, c, q, s, e \}$ and $O_2=\{ \alpha, \theta, \delta, \beta, \gamma \}$;
- 3) In Π_{11} , $O_1=\{ \theta, c \}$ and $O_2=\{ \alpha, \beta, \gamma \}$;
- 4) In Π_{13} , $O_1=\{ c, \gamma, \beta, d_i, b_i \}$ and $O_2=\{ \alpha, \varphi_i, \psi_i \}$; ($i \in [0, n]$)
- 5) In the parallel P modules of Π_{14} , $O_1=\{ c, v, r, \phi_i, \psi_i \}$ and $O_2=\{ \alpha, \alpha_i, \beta_i, \gamma_i \}$; ($i \in [0, n]$)
- 6) In Π_{15} , $O_1=\{ \alpha_i, \beta_i, \gamma_i, c, p_i \}$ and $O_2=\{ \alpha, \delta, \gamma, \beta \}$; ($i \in [0, n]$)
- 7) In Π_{16} , $O_1=\{ c \}$ and $O_2=\{ \alpha, \delta, \gamma, \beta \}$;
- 8) In Π_{17} , $O_1=\{ c, \delta \}$ and $O_2=\{ \alpha, \varepsilon \}$;

As a whole, the biggest P module, $\Pi_{17} < def_{\downarrow \Pi_0}, ref_{\downarrow i_9} >$, includes a closely related combination of two functional combination P modules, $\Pi_{11} < def_{\downarrow \Pi_0}, ref_{\downarrow i_3} >$ and $\Pi_{16} < def_{\downarrow \Pi_{12}}, ref_{\downarrow i_9} >$. Π_{11} implements the algorithm *EstimateSqr*(x) and Π_{19} implements the algorithm *Mbisection*(b, x, a, m). These two P modules form a sequential structure which describes that narrow the computing scope in the first step and accurately calculate to obtain the final result in the next step.

The combination P module $\Pi_{11} < def_{\downarrow \Pi_0}, ref_{\downarrow i_3} >$ mainly describes a square root estimation algorithm for estimating the scope of the square root. Π_{11} is a parallel structure which is combined by 11 P modules $\Pi_i (i \in [0, 10])$. $\Pi_i (i \in [1, 9])$ are the core parallel modules.

In the elementary P module $\Pi_0 < def_{\downarrow \Pi_0}, ref_{\downarrow o_1}, ref_{\downarrow o_2}, ref_{\downarrow o_3}, ref_{\downarrow o_4}, ref_{\downarrow o_5}, ref_{\downarrow o_6}, ref_{\downarrow o_7}, ref_{\downarrow o_8}, ref_{\downarrow o_9} >$, mainly obtain the high-value and the number of digits other than the high-value by processing the original data in parallel. The rules in $\Pi_i < def_{\downarrow \Pi_i}, ref_{\downarrow i_2} > (i \in [1, 9])$ are exactly the same. $\Pi_{10} < def_{\downarrow \Pi_{10}}, ref_{\downarrow i_3} >$ outputs the estimated square root result.

The combination P module $\Pi_{16} < def_{\downarrow \Pi_{12}}, ref_{\downarrow i_9} >$ mainly describes a square root algorithm through *m* bisection calculation approach. In a whole, Π_{16} is a cycle structure, do-while. Two ways which include ($\Pi_{13} \rightarrow \Pi_{14} \rightarrow \Pi_{15}$) and ($\Pi_{13} \rightarrow \Pi_{15}$) can end the cycle. ($\Pi_{13} \rightarrow \Pi_{14} \rightarrow \Pi_{15}$) is executed when the size of the interval is not greater than 2. ($\Pi_{13} \rightarrow \Pi_{15}$) is executed when the exact square root value are obtained. If the size of interval is more than 2, the cycle continue. $\Pi_{12} < def_{\downarrow \Pi_{12}}, ref_{\downarrow a_0}, ref_{\downarrow a_1}, \dots, ref_{\downarrow a_m} >$ split the interval into *m* equidistant intervals and obtain *m* + 1 copies of endpoint number in the interval, $\varphi_0, \varphi_1, \dots, \varphi_m (\psi_0, \psi_1, \dots, \psi_m)$ Also, Π_{16} is a parallel structure. $\Pi_{30+i} (i \in [0, m])$ is the parallel P modules in the combination P module Π_{16} , which calculate a formula, $f(\varphi_i, \psi_i) = \varphi_i * \psi_i - \alpha (i \in [0, m])$.

D. Calculate Instance

We assume that the inputting data is 69399 and the number of parallel computing P module is 11, namely, $m=11$. So 69399 copies of objects into the membrane system to evolve.

- 1) The membrane structure after the original data, Π_0 includes $O_2 = \{ \alpha^{69399} \}$ and $O_1 = \{ c \}$. Then obtaining the set of public objects, $O_2 = \{ \delta^6, \theta^2, \alpha^{69399} \}$.
- 2) The parallel P modules in Π_{10} , namely, Π_1, \dots, Π_9 , inherit the public objects from Π_0 and start the core calculation of *EstimateSqr*(69399). After the parallel calculation, Π_{11} inherits $\beta^2, \gamma^1, \theta^2, \alpha^{69399}$ from Π_{10} . Then Π_{11} gets $\beta^{200}, \gamma^{100}$ and α^{69399} after calculation.
- 3) Start *Mbisection*(200, 69399, 100, 11). Π_{13} inherits the objects from Π_{11} , so the inputting objects in Π_{19} are $\beta^{200}, \gamma^{100}$ and α^{69399} . By the rules, Π_{13} gets objects, $O_2 = \{ \phi_i^{200+i*10} (i \in [0, 10]), \psi_i^{200+i*10} (i \in [0, 10]), \alpha^{69399} \}$. Then $\Pi_{30+i} (i \in [0, 10])$ correspondingly inherit O_2 in Π_{13} . Namely, in $\Pi_{30+i} (i \in [0, 10])$, $O_2 = \{ \phi_i^{200+i*10} (i \in [0, 10]), \psi_i^{200+i*10} (i \in [0, 10]), \alpha^{69399} \}$.
- 4) After the first cycle of the algorithm, *Mbisection*(200, 69399, 100, 11), Π_{16} gets $\alpha_6^{260}, \beta_7^{270}$ and α^{69399} from Π_{36} and Π_{37} . Then Π_{16} gets β^{260}, γ^{10} and α^{69399} from Π_{15} .
- 5) After the second cycle of the algorithm, *Mbisection*(260, 69399, 10, 11), Π_{16} gets $\alpha_1^{262}, \beta_2^{264}$ and α^{69399} from Π_{31} and Π_{32} . Then Π_{16} gets β^{262}, γ^2 and α^{69399} from Π_{15} and ends, while Π_{17} inherits δ^{263} from $\Pi_{15} (\delta^{263} = \beta^{262} + c)$ and saved as ε^{263} . The objects ε^{263} represents the square root of 69399. That means the final result is the number, 263.

VI. CONCLUSION

In this paper, the design and assembly of the P module is formed to provide a framework for constructing a P system by recursive combined P modules, so that the P module combines the three characteristics of packaging, information hiding and modular combination. By designing a well-structured P module, the correctness of the dynamic execution of the P system is ensured with a good structure, the efficiency of software development is improved, and the error rate is reduced.

As the application of this paper, we use the definition and design method of P module to solve the large number square root problem. By designing a P system for solving the square root of a large number, the correctness and high efficiency of the P system design method based on the P module are clarified.

The design method in this paper is mainly aimed at the cell-like P system and further work can apply it to other models of P systems, such as the tissue-like P system and neural-like P system.

REFERENCES

- [1] C.Martín-Vide, Gh.Păun, G.Rozenberg. "Membrane systems with carriers". Theoretical Computer Science, 2002, 270 (1): 779-796
- [2] A.Vitale, G.Mauri, C.Zandron. "Simulation of abounded symport/antiport P system with Brane calculi". Biosystems, 2008, 91(3): 558-571.
- [3] R.Freund, Gh.Păun, M.J.Pérez-Jiménez. "Tissue P systems with channel states". Theoretical Computer Science, 2005, 330(1): 101-116.
- [4] O.H.Ibarra, A.Păun, et al. "Normal forms for spiking neural P systems". Theoretical Computer Science, 2007, 372(2): 196-217.
- [5] Guo P, Dai Y, Chen H. "A p system for Hamiltonian cycle problem". Optik - International Journal for Light and Electron Optics, 2016, 127(20):8461-8468.

- [6] Guo P, Zhu J, Zhou M. "A Family of Uniform P Systems for All-Satisfiability Problem". Journal of Computational & Theoretical Nanoscience, 2016, 13(1):135-142.
- [7] Ping Guo, Yuwen Zhong, Haizhu Chen, Mingqiang Zhou. "A P system for finding all solutions of the degree-constrained spanning tree problem". Pre-Proceedings of the Second Asian Conference on Membrane Computing (ACMC2015), 2015.
- [8] Ping Guo, Junqi Xiang, Jingya Xie, Jinhang Zheng, "A P System for Solving All-Solutions of TSP", International Journal of Advanced Computer Science and Applications, Vol. 8, No. 9, 2017, 357-364
- [9] Yahya, R.I., Shamsuddin, S.M., Hasan, S., Yahya, S.I., "Tissue-like P system for segmentation of 2D hexagonal images". ARO- Sci. J. Koya Univ. 4(1), 3542 (2016).
- [10] Isawasan, P., Venkat, I., Subramanian, K., Khader, A., Osman, O., Christinal, H., "Region-based segmentation of hexagonal digital images using membrane computing". In: 2014 Asian Conference on Membrane Computing (ACMC). IEEE (2014).
- [11] A.Atanasiu. "Arithmetic with membranes". Romanian Journal of Information Science and Technology, 2001, 4(1): 5-20.
- [12] P.Guo, J.Chen. "Arithmetic operation in membrane system". In:Proceedings of International Conference on BioMedical Engineering and Informatics. Sanya, Hainan, China, 231-234(2008).
- [13] Ruilong Yang, Ping Guo, Jia Li, Ping Gu, "Arithmetic operations with membranes based on arithmetic formula tables", Chinese Journal of Electronics, 2015, 24(2).
- [14] Ping Guo, Lijiao Wei, Haizhu Chen. "A P Systems for Matrix-Vector Multiplication", Journal of Computational and Theoretical Nanoscience, Vol.12, No. 11, pp. 4279-4288, 2015.
- [15] Francisco J. Romerocampero, Jamie Twycross, Miguel Camara, et al. "Modular assembly of cell systems biology models using P systems". International Journal of Foundations of Computer Science, 2009, 20(03):427-442.
- [16] Șerbănuță T, Ștefănescu G, Roșu G. "Defining and Executing P Systems with Structured Data in K". Membrane Computing. Springer-Verlag, 2009:374-393.
- [17] Păun G, Pérez-Jiménez M D J. "Solving Problems in a Distributed Way in Membrane Computing: dP Systems". International Journal of Computers Communications & Control, 2010, 5(2):238-250.
- [18] Dinneen M J, Kim YB, Nicolescu R. "P systems and the Byzantine agreement". Journal of Logic & Algebraic Programming, 2010, 79(6):334-349.
- [19] Dinneen M J, Kim Y B, Nicolescu R. "A Faster P Solution for the Byzantine Agreement Problem". International Conference on Membrane Computing. Springer-Verlag, 2010:175-197.
- [20] Dinneen M.J., Kim YB, Nicolescu R. "Synchronization in P Modules". International Conference on Unconventional Computation. Springer-Verlag, 2010:32-44.
- [21] Nicolescu R. "Parallel and Distributed Algorithms in P Systems". International Conference on Membrane Computing. Springer-Verlag, 2011:35-50.
- [22] Nicolescu R, Wu H. "BFS Solution for Disjoint Paths in P Systems". UC'11 Proceedings of the 10th international conference on Unconventional computation, Springer-Verlag, 2011:164-176.
- [23] Guo P, Zhang H, Chen H, et al. "Fraction Reduction in Membrane Systems". The Scientific World Journal,2014(2):858527.

APPENDIX: RULESET

- 1 The rules in Π_0

$$r_1: S_0\alpha \rightarrow_{max} S_1\alpha\tau;$$

$$r_2: S_1\tau^{100}c \rightarrow_{min} S_2\tau^{100}c; 1$$

$$r_3: S_1\tau c \rightarrow_{min} S_5\tau c; 2$$

$$r_4: S_2\tau^{100} \rightarrow_{max} S_3b;$$

$$r_5: S_3\tau \rightarrow_{max} S_4;$$

$$r_6: S_3b \rightarrow_{min} S_4b\theta;$$

$$r_7: S_4b \rightarrow_{min} S_1\tau;$$

$$r_8: S_5\tau \rightarrow_{min} S_0\delta;$$
- 2 The rules in Π_1 to Π_9

- $r_9: S_0\delta b \rightarrow_{max} S_1d;$
- $r_{10}: S_1bs \rightarrow_{min} S_2bs; 1$
- $r_{11}: S_1s \rightarrow_{min} S_4s; 2$
- $r_{12}: S_2bc \rightarrow_{min} S_3e;$
- $r_{13}: S_3bs \rightarrow_{min} S_4bs; 1$
- $r_{14}: S_3s \rightarrow_{min} S_5s; 2$
- $r_{15}: S_4\delta \rightarrow_{max} S_0;$
- $r_{16}: S_4e \rightarrow_{max} S_0bc; 1$
- $r_{17}: S_4d \rightarrow_{max} S_0b; 2$
- $r_{18}: S_4\alpha \rightarrow_{max} S_0;$
- $r_{19}: S_4\theta \rightarrow_{max} S_0;$
- $r_{20}: S_5qs \rightarrow_{max} S_0\gamma; 1$
- $r_{21}: S_5q \rightarrow_{max} S_0\beta; 2$
- $r_{22}: S_5e \rightarrow_{max} S_0bc; 2$
- $r_{23}: S_5d \rightarrow_{max} S_0b; 2$

3 The rules in Π_{10}

- $r_{24}: S_0\theta c \rightarrow_{max} S_1c$
- $r_{25}: S_1\beta \rightarrow_{max} S_0\beta^{10}$
- $r_{26}: S_1\gamma \rightarrow_{max} S_0\gamma^{10};$

4 The rules in Π_{12}

- $r_{27}: S_0\gamma^n \rightarrow_{max} S_1d_0d_1 \dots d_n; 1$
- $r_{28}: S_0\beta \rightarrow_{max} S_1b_0b_1 \dots b_n; 1$
- $r_{29}: S_0\gamma c \rightarrow_{max} S_1\gamma c; 2$
- $r_{30}: S_1\gamma c \rightarrow_{max} S_2cd_0d_1 \dots d_n; 1$
- $r_{31}: S_1c \rightarrow_{max} S_2c; 2$
- $r_{32}: S_1\gamma \rightarrow_{max} S_2; 3$
- $r_{33}: S_2d_i \rightarrow_{max} S_0\varphi_i^i\psi_i^i;$
- $r_{34}: S_2b_i \rightarrow_{max} S_0\varphi_i\psi_i;$

5 The rules in Π_{13}

- $r_{35}: S_0\alpha_m \rightarrow_{max} S_1\alpha; 1$
- $r_{36}: S_0\gamma_i c \rightarrow_{min} S_1c\gamma_i; 1$
- $r_{37}: S_0c \rightarrow_{min} S_2c; 2$
- $r_{38}: S_1\gamma_i \rightarrow_{max} S_3\delta;$
- $r_{39}: S_2\alpha_i\beta_{i+1} \rightarrow_{max} S_3p_{i+1}\beta_{i+1}\beta;$
- $r_{40}: S_3\beta_i p_i \rightarrow_{max} S_3\gamma; 1$
- $r_{41}: S_3p_i \rightarrow_{max} S_0; 2$
- $r_{42}: S_3\alpha_i \rightarrow_{max} S_0; 2$
- $r_{43}: S_3\beta_i \rightarrow_{max} S_0; 2$

6 The rules in Π_{14}

- $r_{44}: S_0c\gamma^3 \rightarrow_{min} S_0c\gamma^3; 1$
- $r_{45}: S_0c \rightarrow_{min} S_1c; 2$
- $r_{46}: S_1\beta \rightarrow_{max} S_0\delta;$
- $r_{47}: S_1c \rightarrow_{max} S_0\delta c;$
- $r_{48}: S_1\gamma \rightarrow_{max} S_0;$

7 The rules in Π_{15}

- $r_{49}: S_0\delta \rightarrow_{max} S_0\epsilon;$

8 The rules in Π_{30} to Π_{30+m}

- $r_{50}: S_0\psi_i c \rightarrow_{min} S_1c; 1$
- $r_{51}: S_0c \rightarrow_{min} S_2c; 2$
- $r_{52}: S_1\varphi_i \rightarrow_{max} S_0\varphi_i r;$
- $r_{53}: S_2\alpha r \rightarrow_{max} S_3v;$
- $r_{54}: S_3\alpha v \rightarrow_{max} S_4\alpha v; 1$
- $r_{55}: S_3r v \rightarrow_{max} S_5v; 2$
- $r_{56}: S_3v \rightarrow_{max} S_6v; 3$
- $r_{57}: S_4\varphi_i \rightarrow_{max} S_7\alpha_i;$
- $r_{58}: S_5\varphi_i \rightarrow_{max} S_7\beta_i;$
- $r_{59}: S_6\varphi_i \rightarrow_{max} S_7\gamma_i;$
- $r_{60}: S_7v \rightarrow_{max} S_0\alpha;$
- $r_{61}: S_7r \rightarrow_{max} S_0;$

Cascades Neural Network based Segmentation of Fluorescence Microscopy Cell Nuclei

Sofyan M. A. Hayajneh*, Mohammad H. Alomari†, Bassam Al-Shargabi‡

*Department of Electrical and Computer Engineering
Isra University, Amman, Jordan

†Department of Electrical Engineering
Applied Science Private University, Amman, Jordan

‡Department of Computer Information System
Isra University, Amman, Jordan

Abstract—The visual extraction of cellular, nuclear and tissue components from medical images is very vital in the diagnosis routine of different health related abnormalities and diseases. The objective of this work is to modify and efficiently combine different image processing methods supported by cascaded artificial neural networks in an automated system to perform segmentation analysis of medical microscopy images to extract nuclei located in either simple or complex clusters. The proposed system is applied on a publicly available data sets of microscopy nuclei cells. A GUI is designed and presented in this work to ease the analysis and screening of these images. The proposed system shows promising performance and reduced computational time cost. It is hoped that thus system and the corresponding GUI will construct platform base for several biomedical studies in the field of cellular imaging where further complex investigations and modelling of microscopy images could take place.

Keywords—Artificial neural networks; machine learning, DSP, fluorescence microscopy; biomedical imaging; cell nuclei; image segmentation

I. INTRODUCTION

A. Digital Image Processing

Digital image processing and analysis is among rapidly growing technologies. It encompasses a wide-ranging field of applications in our everyday life. Medical, industrial, text recognition, biometrics and graphics, are just examples of hundreds of possible applications of image processing. Although every single application needs a well-designed approach to parse and extract the required useful information and, most of these approaches can be categorized under a single or multiple major aspects that include but are not limited to: image visualization, sharpening, enhancement, recognition, verification, retrieval, segmentation and /or restoration, etc. [1]-[4].

Most of the image processing-based applications are assembled and designed by applying a schematic classical methodology that includes one or more of the following phases [1]. Image pre-processing and enhancement, Objects segmentation (extraction), Statistical features extraction, Objects (Candidates) selection and pruning, Objects post-processing and (or) Features classification.

Segmentation is an important phase in image analysis where the image is divided into meaningful disjoint regions with similar properties, such as gray level, color, texture, brightness, contrast, etc. It is often one of the first and most difficult phases in image analysis. Due to its importance, a great variety of segmentation algorithms have been proposed to tackle a wide range of applications such as microscopy, biomedical engineering, biomedical imaging, bioinformatics and pattern recognition [5]-[9]. The segmentation process will be the focus of our attention in this work and more details about segmentation process will be highlighted in Section II.

B. Medical Imaging and Computer Aided Diagnosis (CAD)

Medical imaging techniques are usually implemented to greatly enhance the extraction process of numerical features that provide efficient and as sufficient as possible representation of medical signs features, activities or general regions of interest in the different types of medical images acquired by different types of medical instruments. These images along with the possible integration of advanced image processing techniques and medical and physiology concepts can positively improve the ability of knowledge extraction and increase the ability to screen and (or) predict diseases that may have serious impacts on our daily life [10]-[12].

The difficult medical diagnostic routine (time consuming and tedious process) can be improved by providing the specialists (e.g. radiologists, pathologists, biologists) with quantitative data and statistical measurements which are extracted from such medical images so that the visualized version is much more informative [1], [13]. The development in computational power (e.g. processor speed, RAM, graphical hardware) has driven the development of several image processing algorithms that have had a significant impact in several medical research and applications.

Recently, computer aided diagnosis (CAD) has quickly become a widespread and unescapable useful tool for diagnostic examinations in many daily routine works across a wide range of medical and clinical areas such as microscopy imaging, tissue culturing, cancer research, Confocal microscopy, heart diseases, brain tumors, blood diseases screening, etc. [14]-[18]. Such CAD (Fig. 1) systems have drawn a lot of attention because they can represents a second opinion for the specialist and they allow a large scale

statistical evaluation besides the classical human screening evaluation [19], [20].

C. Fluorescence Microscopy and Nuclei Image Segmentation

In fluorescence microscopy, the object under consideration is itself the light source rather than a light reflector. This is due the fact that certain material absorbs light at one wavelength (excitation) and emits a detectable visible (i.e. longer specific wavelength) light. This physical phenomenon (fluorescence) is utilized to visually and separately observe different components of the specimen such as cell membrane, cytoplasm, nuclei, gene, tissue, or proteins (Fig. 2).

The component of interest is specifically stained in the specimen using particular fluorescent dyes and then the fluorescence microscopes make it possible to visualize and store digital images of this component [21]. In fluorescence microscopy applications, researchers are typically interested in as much as accurate localization of the boundaries of the observed fluorescently labelled structural and functional units (cell nuclei, genes, etc.).

Heterogeneity of different biological features can be an issue arising from the use of fluorescence imaging data. However, to a large extent this may be overcome by specially designed CADs, which can correctly take into account the physical variations seen between cells and therefore across a set of fluorescence microscopy images.



Fig. 1. Example of a CAD microscope (Image courtesy: Carl Zeiss- ELYRA Super resolution Microscopy).

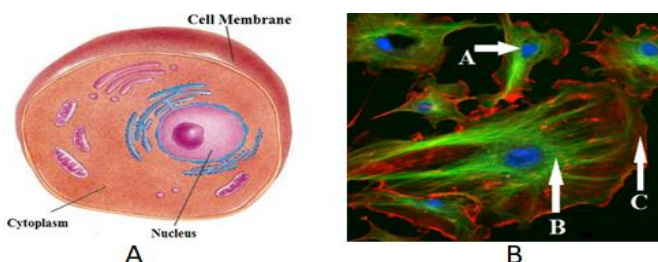


Fig. 2. (A) The main cell structure, (B) A specimen showing Endothelial cells stained with different dyes (Image courtesy: http://en.wikipedia.org/wiki/Fluorescence_microscopy).

II. REALTED WORKS

A. Cellular and Nuclear Image Segmentation Literature Background

Semantic image segmentation in fluorescence microscopy analysis refers to the separation process of cell components from the surrounding background by finding the boundaries of cellular, nuclear or histological structures with an adequate accuracy from images of stained tissues with different

markers, Fig. 2B. Nuclear segmentation is an important step in the pipeline of many cytometry analyses because it forms the basis of many other operations (e.g. cell counting, cell-cycle assignment, cell tracing) and is often the first step in the overall cell segmentation [22]. An increasing number of screenings and investigations are done using different types of fluorescence microscopy images either on individual or sequence (i.e. live-cells imaging) images.

In fluorescence labelled images of blood and bone marrow, high degrees of nuclei segmentation accuracy is reported by applying a classical image processing techniques such as shading correction and background (grayscale opening) followed by Otsu's method and watershed algorithm based on inverse distance transform [23]. In [24] a modified algorithm using the watershed algorithm based on morphological filtering operations is applied to choose the seeds of cell nuclei in tissue sections (i.e. foreground) and background as well. In this case, the merging of touching and overlapping regions is used to solve the over-seeded situations. In [24] method, it is required to manually choose and set specific values of certain parameters based on test images and then use them on images of the same dataset or images taken under the same conditions.

In [25], the problem of touching cells is addressed and treated by detecting the concave points from the polygonal approximations. After applying morphological filtering and adaptive thresholding to detect contour, this contours are segmented using the concave points. This approach is combined with a customized ellipse shape fitting such that each segment of the contour has a fitted ellipse.

A fully automated approach based on graph cut model is also used for segmenting the touching cell nuclei [26]. The background and the foreground separation is achieved using a minimal geodesic length, then the individual nuclei are found by a graph cut which include image gradient information and a priori knowledge about the shape of the nuclei. The graph-cut is also used for cells segmentation for the tracking problem in microscopy images [27].

The advantages of active contour method (flexibility), multi-resolution methods (low processing time), multi-scale methods (smoothing) and region-growing methods (statistical modelling) are combined in [28] to construct an accurate and fast cell nuclei segmentation algorithm.

A new method for leucocytes segmentation based on nuclei classification is presented in [29]. The overlapping and isolated configuration situations are classified based on Bayesian networks and stepwise merging strategy. Some morphological features of the nuclei, such as the compactness, smoothness and moments are used followed by a watershed algorithm to find the proper nuclei boundaries. The overlapping nuclei are segmented into isolated nuclei using an intensity gradient transform and watershed algorithm.

Some artificial neural network based approaches such as bidirectional associative memory (BAM) [9] is effectively used in medical segmentation application because it has some preferable features such as: supervision is required only for selecting texture primitives, no training is required and it is

robust for the presence of noise and distortion. More reported segmentation algorithms in the field of fluorescence cell nuclei and histological structures can be found in [19]. In [30], a supervised learning-based system is proposed for segmenting different types of biomedical images where the focus was to describe a general purpose system that can, with few modifications, be used in a variety of image segmentation applications as long as enough labeled data is available for training. The system used the intensity neighborhoods as nonparametric feature vectors for pixel classification.

Although Fluorescence microscopy is a rapid expanding technique and it has made it possible to identify cellular components with a high degree of specificity, more attention is required to make its analysis fully automated and as meaningful as possible. Such work is challenging due to the large variations of features of the cellular components (size, shape, orientation, texture, etc.).

A totally automatic system is still not a reality; so much work remains, mainly in the early steps, which may involve segmentation, recognition of nuclei from the background followed by refinement, counting, and statistics calculations. From another side, many works still focus in certain regions (ROI) inside the image under processing rather than processing the whole image (i.e. inner and boundaries).

The main aim of this work and its further consequence modules is to focus on the specifics of images acquired by fluorescence microscopes (in particular, images cell nuclei) and design successful and efficient fully automated segmentation system that can be used to overcome these specifics simultaneously using tuned image processing and artificial intelligence techniques.

B. Clinical Nuclei Related Work

Many variations of the basic and CAD microscope instrument are now available and used in great success applications, allowing us to look into spaces much too small to be seen with the naked eye. The processing of digital microscope images includes the utilization of digital image processing technologies to analyze, model and visualize these images. Medicine, chemistry, cancer research, pharmacology, biological research and numerous related fields are common places for this type of microscope image processing. The processing of such images is improved by designing a direct special interfacing of microscope imaging instruments with image processing systems and interfaces (Fig. 1) [20], [21].

In many applications, it is very important to achieve accurate and efficient segmentation, classification and grouping of nuclei and cells in fluorescence microscopy images. This importance comes to enhance the understanding of cells functions [31] and enters the processing workflow of pathological diagnosis [32]. Examples of this include the immune-histochemically staining estimation and morphological grading where the detection of cell nuclei on histological slides is required.

From another hand, the precise quantitative statistics about nuclear structure and morphology along with their visualization can uncover important clues for the diagnosis of benign, pre-neoplastic, and neoplastic (cancer) lesions. Also, this type of quantification and classification should ease the understanding of the anatomical variation of different organs by the analysis of their corresponding tissues [33].

The nuclei segmentation were used as a basis to investigate the subcellular localization of proteins at a proteome-wide scale [34]. The cell nuclei were considered as seeds to perform protein segmentation using the watershed method and hence be able to identify the subcellular localization patterns even of complex ones.

C. Dataset

In this paper, the publicly available U2OS dataset described in [22] will be considered. Fig. 3 shows three images from the datasets under consideration. Although these images depict different levels of segmentation difficulty, most of them follow a similar histogram general profile where double-peaks (bimodal) are clear. It can be distinguished that in many occasions clustered touching nuclei are exist. Also, it is easy to figure out that the histograms of these images do not occupy the full dynamic range of the gray scale (i.e. 0-255) and they are concentrated on the low (dark) side of the scale. Different issues such as touching objects also takes place; these issues will be considered in the design of the different stages of the proposed segmentation system as will be described later.

III. PROPOSED METHODOLOGY

The main outline of the proposed methodology can be summarized as follow (Fig. 4) Stage-1) Preprocessing. Stage-2) Detecting the nuclei candidates using NN, Stage-3) Classification of candidates into single isolated nucleus or clustered (overlapped or touching) nuclei, Stage-4) Separation of individual nucleus within clusters using modified watershed algorithm combined with NN iteratively, and 5) Refining each separated nucleus. The following subsections describe the details of these stages, followed by qualitative and quantitative evaluations.

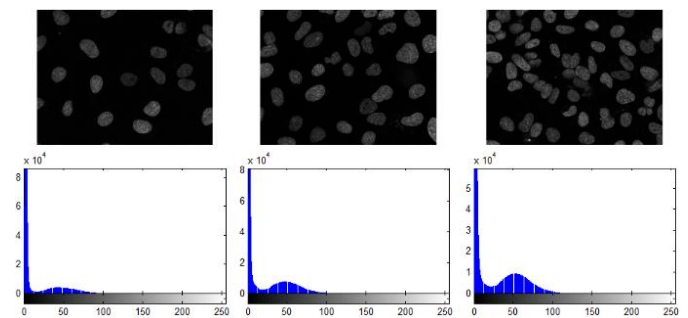


Fig. 3. (Top) Three examples of different cell nuclei images. (Left) dna-20, well separated with low number of nuclei, (Middle) dna-19, moderate level. (Right) dna-41, difficult clustered nuclei [22]. (Bottom) The corresponding histograms.

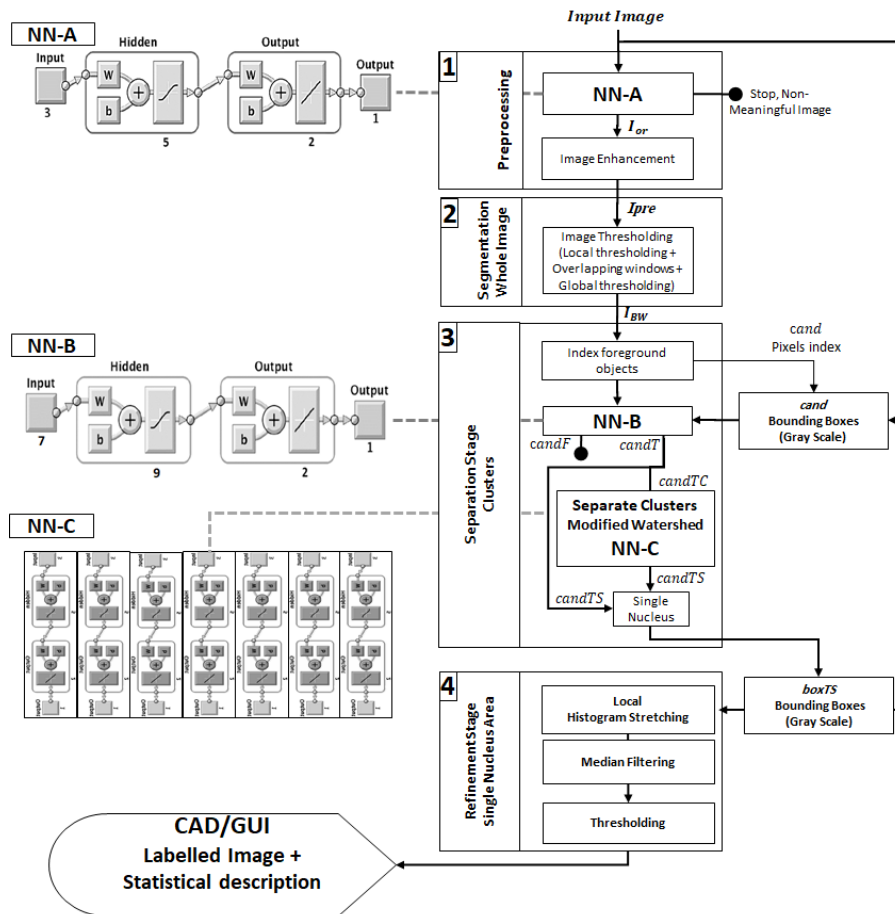


Fig. 4. The proposed cell nuclei segmentation system showing the multi stages of detection and verification.

A. Preprocessing Stage

Some input images contain no meaningful visual data, hence no nuclei can be identified or some noisy bright areas may be wrongly identified as nuclei. In the preprocessing stage, a neural network (NN-A) with back propagation algorithm (BPA) is used to identify this type of images. Part of the images were manually investigated and each image is assigned a flag (F1, for meaningful images; F2, for non-meaningful images). Then, three statistical features [(1)-(3)] were calculated for each image. The NN-A structure is: three inputs nodes, one hidden layer of 5 nodes and two output nodes. The training vector of NN-A is $[\mu, \sigma, E]$ and the neuron for output layer indicates F1 or F2 cases.

$$\mu = \frac{\sum_{r=1}^R \sum_{c=1}^C IM_{r,c}}{R * C} \quad (1)$$

$$\sigma = \sqrt{\frac{\sum_{r=1}^R \sum_{c=1}^C (IM_{r,c} - \mu)^2}{(R * C) - 1}} \quad (2)$$

$$E = - \sum_{n=0}^{255} p_n \log_2 p_n \quad (3)$$

where μ, σ, E are the mean, standard deviation and entropy of the input image (IM), respectively. R and C are the number of rows and columns of IM, respectively, and p contains the 255-bins histogram count of IM. It is found that NN-A can

exclude all non-meaningful images. This step will be useful in the generalization of the proposed method for other types of medical images.

Original image (hereinafter, I_{or}) that has high information content is then preprocessed to reduce the noise and enhance the quality. Although the nuclei in such images are usually brighter than the dark background of the surrounding tissue, it is usually difficult to identify nuclei consistently over the whole image because some of these images are subject to non-uniform illumination and noise. To smooth out the possibility of uneven illumination, a gray morphological top-hat operator is applied.

$$I_{un} = I_{or} - ((I_{or} \ominus StrEl) \oplus StrEl) \quad (4)$$

where \ominus and \oplus denote erosion and dilation, respectively and $StrEl$ is the structure element used to perform the opening operation. It is found that the optimum size of $StrEl$ is correlated to the mean value of the image (the greater the area of nuclei in the image, the lesser the background area and hence the mean is greater).

Then, a common adaptive histogram stretching approach is applied on I_{un} to enhance the image contrast. In this approach, (18 pixels \times 18 pixels) sub images are processed individually to achieve contrast stretching of 95% (centered) of the sub image histogram, then the processed sub images are combined

through a bilinear interpolation to compensate for the possible induced boundaries between these processed portions. Then a 9×9 median filter is applied to compensate for the increase in the noise amplitude which comes as a result of the above enhancement steps [1]. The result of this step is hereinafter called (I_{pre}), an example is shown in Fig. 8B.

B. Segmentation Stage

Generally, segmentation process can be achieved based on the separation of regions that have similar properties, such as gray level, color, texture, brightness, contrast, etc. In this stage, an adaptive gray level based thresholding approach is employed as follows:

- 1) Calculate TH_G ; the Otsu's [35] threshold of the global preprocessed image (I_{pre}).
- 2) Divide I_{pre} into square blocks (64 pixels × 64 pixels) containing the gray levels (BLK_g).
- 3) For each block i (BLK_{g_i}), do the following:

a) Calculate

$$TH_{avg} = w \cdot TH_G + (1 - w) \cdot TH_L \quad (5)$$

where TH_{L_i} is the Otsu's threshold of BLK_{g_i} . And w is weighting factor, it is found experimentally that $w = .85$ is an optimum choice.

b) Convert each block BLK_{g_i} individually to a binary block (BLK_{BW_i}) using TH_{avg} threshold calculated in step 0.

c) Combine the resulted binary blocks together to form the whole black-white image I_{BW1} .

4) Repeat steps 2) and 3) above using a block size of (32 pixels × 32 pixels), this ends up with I_{BW2} .

5) Calculate $I_{BW} = I_{BW1} \cdot I_{BW2}$; the bitwise logical Anding operation.

6) A series of morphological opening (using a disk shaped structure elements of 2, 3, 4 pixels radius) are performed on I_{BW} image to refine the shapes of the detected objects and remove noisy pixels and tiny white objects. Also, this operation removes small objects (i.e. dark pixels) inside the foreground of an image.

7) Then, a morphological flood filling operation is performed to remove small holes in the foreground.

C. Single Nucleus and Clustered Nuclei Separation

In this stage, the foreground objects in the binary I_{BW} image are isolated and indexed (labelled) individually, each object forms a region of possible nuclei candidates ($cand$). Some $cand$ regions contain single nucleus, some regions contain nuclei cluster. These two types are called ($candT$). From the other hand, some regions contain faulty $cand$ ($candF$) (i.e. Noise detected as nucleus), (Fig. 5).

D. Non-nucleus Candidates Pruning

It is recognized that some gray intensity properties of $candF$ regions are not highly changeable. As an example, it is found by deep investigation that the Otsu's threshold of nuclei blocks (i.e. blocks contains true nuclei) is almost close to the mean of the pixels within the block, while non-nuclei blocks

don't satisfy this. Also, for the nuclei blocks, it is found that Kurtosis (peakedness) [36] of the part of the histogram right to the Otsu's threshold is higher than that of non-nuclei blocks, Fig. 6.

To prune $candF$, a NN (hereinafter NN-B) with BPA training is used. Several NN topologies were designed and tested to determine the optimum NN-B. It is found that the optimum NN-B contains (7 input nodes, one hidden layer with 9 nodes, 2 output nodes).

A number of ($candF$) and ($candT$) regions (similar to Fig. 5) were manually cropped. Then, the features described in (6)-(12) were calculated for these regions. The vector [$\mu_R, \sigma_R, skew_R, kurtosis_R, DR_R, BP, ROB_R$] after normalization represent the input training vector to NN-B. The neuron of the output layer indicates whether the block is ($candF$) or ($candT$).

$$\mu_R = \frac{\sum_{r=1}^{R_R} \sum_{c=1}^{C_R} candR_{r,c}}{R_R * C_R} \quad (6)$$

$$\sigma_R = \sqrt{\frac{\sum_{r=1}^{R_R} \sum_{c=1}^{C_R} (candR_{r,c} - \mu_R)^2}{(R_R * C_R) - 1}} \quad (7)$$

$$skew_R = \frac{1}{R_R * C_R} \sum_{r=1}^{R_R} \sum_{c=1}^{C_R} \left(\frac{candR_{r,c} - \mu_R}{\sigma_R} \right)^3 \quad (8)$$

$$kurtosis_R = \frac{\sum_{r=1}^{R_R} \sum_{c=1}^{C_R} (candR_{r,c} - \mu_R)^4}{R_R * C_R * \sigma_R^4} - 3 \quad (9)$$

$$DR_R = \max_{r,c} candR_{r,c} - \min_{r,c} candR_{r,c} \quad (10)$$

$$BP_R = \frac{\forall (candR_{r,c} \geq 1.25 * \mu_R), \sum_{r=1}^{R_R} \sum_{c=1}^{C_R} (1)}{R * C} \quad (11)$$

$$ROB_R = \frac{BP}{R_R * C_R} \quad (12)$$

where $\mu_R, \sigma_R, skew_R, kurtosis_R$ are the mean, standard deviation, skew, kurtosis of the cropped region ($candR$), respectively. R_R and C_R are the number of rows and columns of $candR$, respectively. DR_R represent the dynamic range of intensities of $candR$, BP represents the number of bright pixels in $candR$ nprmalized to the full image size, and ROB_R represents the ratio of bright pixels to the total area of $candR$. R and C as defined in (1).

By training and testing, it is found that NN-B is capable to prune $candF$ and accept $candT$ for further processes with higher degree of accuracy. Numerical evaluation of the NN-B goodness is included in section IV.

E. Separation of Clustered Nuclei

Some $candT$ contain more than one nuclei, this is caused by touching or overlapping nuclei Fig. 5-B. To get accurate count and statistics of nuclei, it is important to split $candT$ regions that contain clustered nuclei into single nucleus ($candTS$). To achieve this, a region based segmentation using watershed transform is applied on $candT$. The idea behind the basic watershed transform is to define the catchment basins and the watershed lines between them. Generally, the watershed transform is applied to the image gradient.

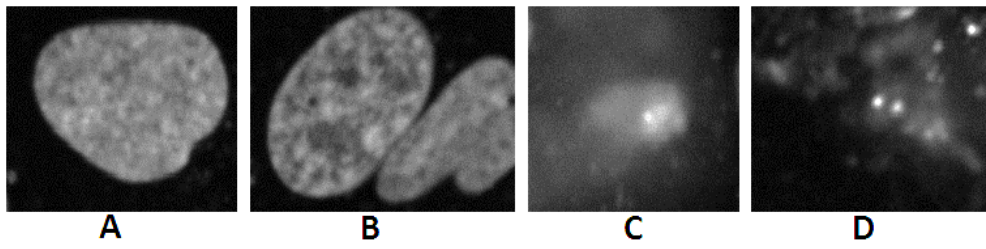


Fig. 5. Samples of different regions, (A) single nucleus (**candT**), (B) clustered touching nuclei (**candT**), (C and D) non-nuclei (**candF**).

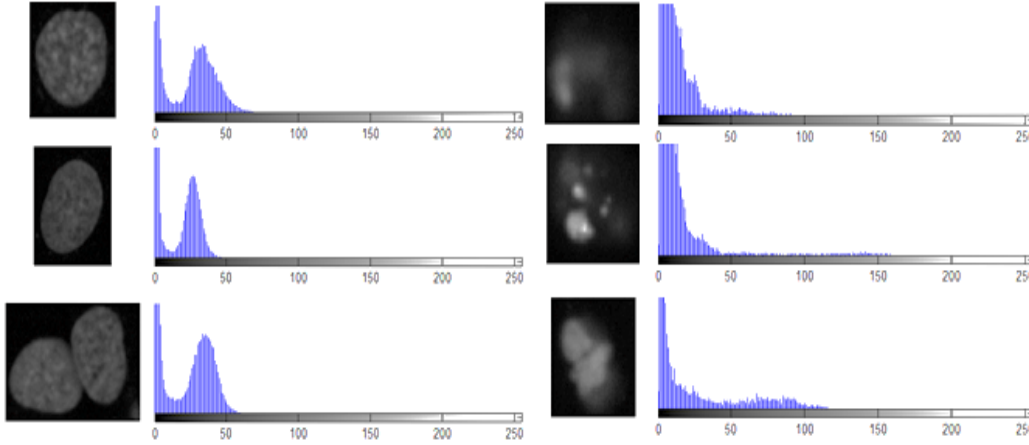


Fig. 6. Samples of different sub regions: (Left) nuclei and their corresponding histograms (right) non-nuclei region and their corresponding histograms.

Direct application of the watershed transform results in a large number of tiny regions. To solve this problem, the proposed system presents a modified iterative watershed algorithm to separate the touching nuclei in nuclei clusters. The iterative scenario is controlled and tuned using a neural network (hereinafter NN-C). To control the number of separated nuclei within each candT, candT is scaled using a variable scaling parameter (SP). The separation algorithms are as follows:

- a) Initiate SP to 1. (all watershed regions are detected).
- b) Define ScandT; the scaled version of candT using SP.
- c) Apply watershed transform on ScandT
- d) Calculate the areas of the detected regions in step C, and choose the 10 largest ($Reg_i ; i \in \{1: 10\}$).
- e) For each Reg_i , using (13)-(19), calculate the following features' vector:

$$FV_i = [Area_i, Per_i, Ratio_i, CenX_i, CenY_i, MinAx_i, MaxAx_i]$$
- f) For each Reg_i , use the trained NN-C to estimate another version of FV_i (called FV_{iNNC})
- g) Calculate the squared error (SErr) between FV_i and FV_{iNNC} .
- h) If SErr is reducing, then reduce SP linearly and perform a new iteration starting from step B. If the reduction in SErr is below a predefined threshold, then stop iterations and choose the regions candTS that have small SErr less than a predefined threshold that were calculated during training.

$$Area_i = \frac{A_i}{R_i * C_i} \quad (13)$$

$$Per_i = \frac{P_i}{R_i * C_i} \quad (14)$$

$$Ratio_i = \frac{Per_i}{Area_i} \quad (15)$$

$$CenX_i = \frac{\sum_{j=1}^{A_i} X_{ij}}{R_i * A_i} \quad (16)$$

$$CenY_i = \frac{\sum_{j=1}^{A_i} Y_{ij}}{C_i * A_i} \quad (17)$$

$$MinAx_i = \text{length of minor axis of } Reg_i \quad (18)$$

$$MaxAx_i = \text{length of major axis of } Reg_i \quad (19)$$

where A_i and P_i are the number of pixels in the area and in the perimeter of the region i , respectively. $Area_i$, Per_i , $CenX_i$ and $CenY_i$ are the area, perimeter, center of mass in horizontal direction and center of mass in vertical direction, respectively.

FV_{iNNC} vector represents the neuron of the output layer of the NN-C. candT and ScandT are normalized to a fixed dimension ($Rn * Cn$) and the mean value for each row and the mean value for each column of both candT and ScandT are calculated and saved. These constructed vector have a length ($2 * Rn + 2 * Cn$) represent the input training and testing vector for NN-C. Fig. 7. shows examples' results of the separation step.

F. Refinement Stage

In this stage, every single nuclei **candTS** is processed alone to refine its borders. This stage will take out local non-nuclei pixels that were misclassified as nuclei members. From the other hand, the local nuclei pixels that were misclassified as non-nuclei will be added to their corresponding nucleus. This is done as follows:

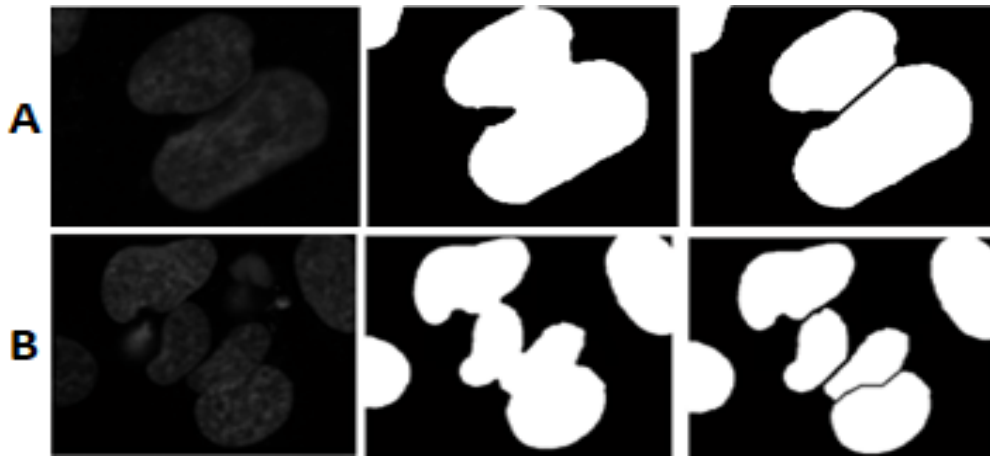


Fig. 7. Two examples of Separation of clustered nuclei. (Left) Enhanced cropped regions. (Center) Initial segmentation (CandT). (Right) Split nucleus (CandTS) using modified watershed and NN-C.

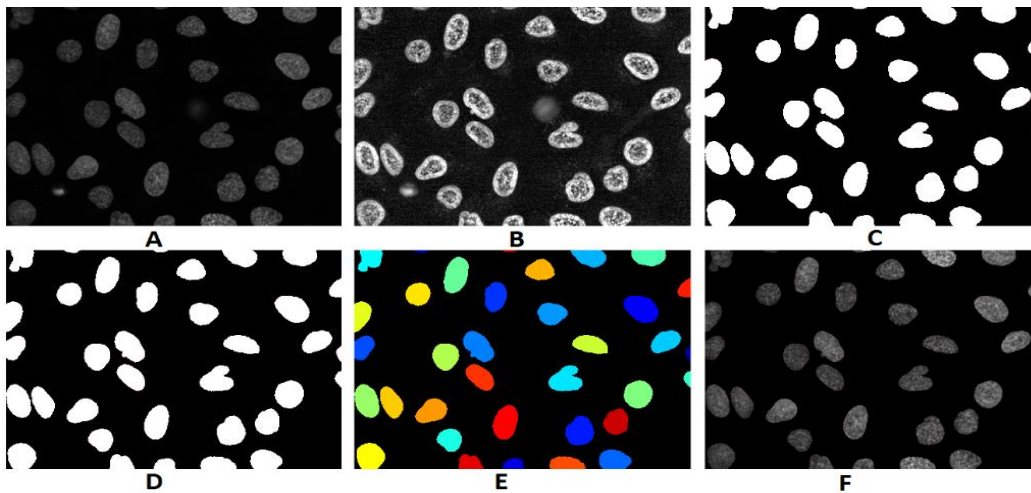


Fig. 8. (A) Input image. (B) Enhanced image. (C) Segmented image. (D) Refined result. (E) Color labeling. (F) Original gray scale profile of the segmented nuclei.

a) Calculate the centroid of the *candTS*.

b) Define the bounding rectangular box *boxTS* that contains all the pixels of *candTS* and 10% more pixels from each side (Top, Bottom, Left, and Right).

c) From the original image I_{or} ; crop out the sub image $I_{or,box}$ contained within the boundaries of *boxTS*.

d) Enhance $I_{or,box}$ individually by applying a histogram stretching transform over 95% of the range and a median filter of size 8*8. This step will increase the independency of each nuclei boundary.

e) Use the calculated centroid in step 1 as a seed to start a region growing based thresholding within *boxTS*. In this step, the region growing is not allowed to expand in the directions where other *candTS* are located. This vital to avoid rejoining of separated nuclei.

It is found that this enhancement stage highly improves the local region of each individual nuclei and hence provide more fine and accurate details of the nuclei borders.

An intensive empirical study based on several cross validation runs were carried out on different nuclei images of

different noise levels and different segmentation difficulty to choose some features' values in the above stages. Some parameters were fixed and some of them are automatically changing (tuned) according to other statistics calculated during the processing of each individual image. No assumptions were assumed and no human interaction is required to choose certain areas in the image under processing (i.e. the proposed system operates on the whole image rather than a specific region of interest ROI).

Along with the proposed system, an interactive graphical user interface (GUI) has been developed. This GUI allows the user to select the data set and then all the associated images are displayed. Fig. 9. shows the developed GUI including all the display options.

IV. PERFORMANCE EVALUATION AND DISCUSSION

For evaluation, the proposed methodology is applied on the (U20S) dataset described above. In order to evaluate the NN-B performance, the false acceptance rate (FDR) is calculated. FDR is an error measure that shows the probability that a *candF* is detected as a *candT*. Samples of the calculated FDR for randomly selected images are shown in Table I.

Row 3 shows the FDR without using the NN-B, while row 4 shows the FDR using the NN-B. It is apparent that the FDR is reduced using the NN-B. Using NN-B reduces the average FDR from 12.8% to 2.1%.

Fig. 10 shows examples of the segmentation obtained by the proposed methodology compared to the provided hand-labelled (ground truth) images. For finer details comparison, a zoomed version for some of the nuclei shown in Fig. 10. are also shown in Fig. 11.

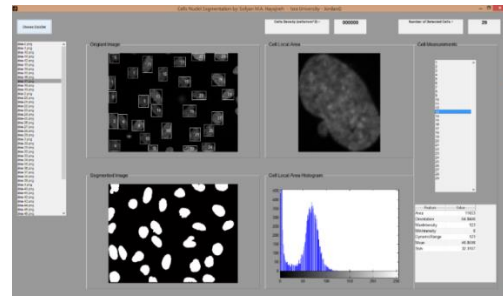


Fig. 9. The proposed GUI.

TABLE I. FDR FOR RANDOMLY SELECTED IMAGES WITH AND WITHOUT USING NN-B

Image number (dna-)	0	10	14	18	23	29	3	30	35	39	4	40	45	48	8	
Number of <i>candT</i>	41	42	47	44	43	37	28	33	29	35	27	34	28	32	28	Avg.
FDR (No NN-B) (%)	4.90	4.80	8.50	20.5	4.70	8.10	21.4	9.10	20.7	14.3	11.1	17.6	3.60	6.30	21.4	12.8
FDR (NN-B) (%)	0	0	2.1	2.3	0	2.7	3.6	0	6.9	2.9	0	2.9	0	3.1	0	2.1

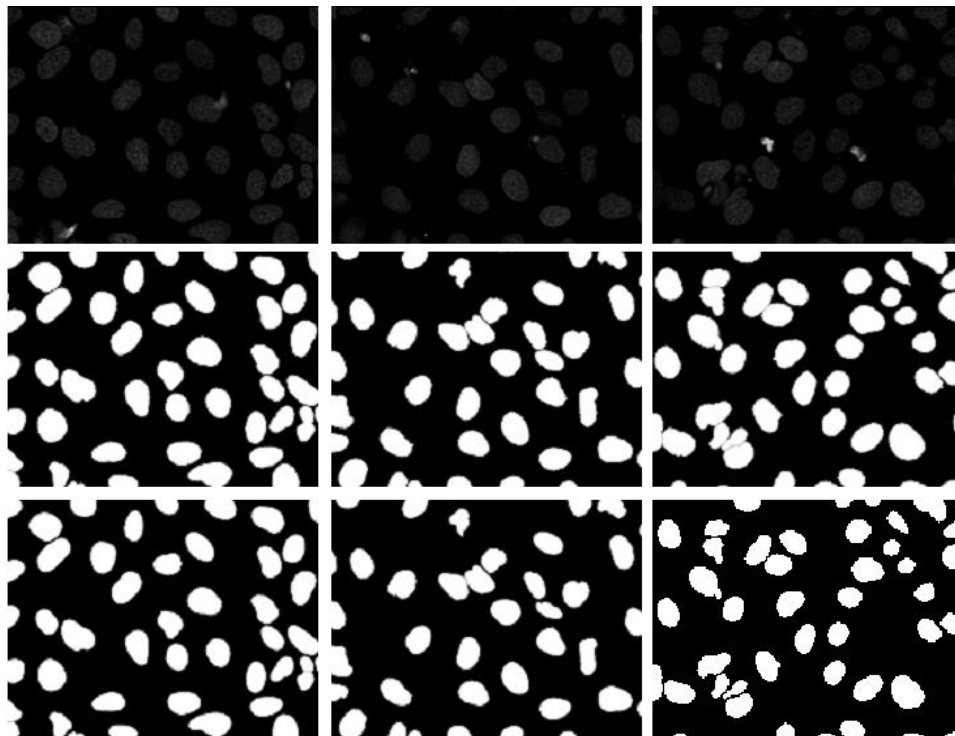


Fig. 10. Three examples of nuclei segmentation. (top) original, (middle) ground truth, (bottom) proposed.

Almost, as Fig. 11 depicts, it is found that ground truth nuclei borders are bigger (outer) than the segmented ones. It is thought that this is due to an oversize segmentation in the hand labelling. Anyway, a simple morphological dilation process could overcome this issue. However, we purposely prefer not to apply this post processing in the refinement stage since the main focus is to describe and evaluate the proposed method. Fig. 12 depicts this fact clearly, it is clear that the total area of all segmented nuclei - using the proposed method - in each image is less than that in the ground truth reference images. This ensures that a simple dilation process could

lead to more area convergence. Again, we purposely prefer not to apply this step.

Although Fig. 10 and 11 show a very qualitative reasonable result, quantitative evaluation should provide another proof of satisfaction. To this end, the proposed method is compared to other methods in literature [22] that were applied on the (U20S) dataset.

Let (IP) represents the proposed binary image output and (IGT) represents the ground truth binary image, [22] described different quantitative performance metrics that include: (1) Hausdroff (HD) metric and normalized sum of distances

(NSD), both can be considered as a spatially based metrics, (2) (Split) error: two IP nuclei are assigned to a single IGT nucleus, (Merge): single IP nucleus embody two IGT nuclei, (Add): an IP nucleus is assigned to the IGT background, and (Miss): an IGT nucleus is assigned to the IP background. Table II shows the result of the proposed methodology compared to other methods.

As stated above, the proposed method provides slight less area in most of the cases. However, this does not affect the error count metrics as they depend on counting objects numbers and comparing pairs of objects (i.e. one to one object in IP and IGT). From the other hand, this slightly affects the values of the HD and NSD as they are aware of the spatial locations of nuclei and their contours so that good values of HD and NSD metrics can still apparently be obtained.

For the same reason, we got lower Dice metric value (0.88), this value can reach an average of 0.94 when a dilation operation is applied on the final segmented binary nuclei. Again, we purposively prefer not apply such preprocessing step.

In general, using the watershed approach leads to less merge errors while increasing the split error [22]. In the

proposed algorithm, it is apparent (Table II) that we can achieve less merge errors using the watershed while keeping split errors at low levels. This is due to applying a modified region growing algorithms that can merge regions that are wrongly separated (split) in the watershed step.

Also, a reasonable (add error) metric were obtained due to applying the non-nucleus pruning operation based on an NN training. From the other hand, the (miss error) were kept at low levels due to applying an adaptive thresholding criteria that takes into account the local and global intensity variations through the image.

It is important to compare the computational cost of the proposed method to some previous methods. It can be shown that the proposed method presents a low computational approach while providing reasonable results. In average, the proposed method takes ~23 seconds to label the input image. This represents a significant reduction in computational cost compared to many other methods that include [30] which takes 49.2 minutes using a supervised learning approach, while in [41] it takes 30 minutes using the template matching approach.

TABLE II. QUANTITATIVE COMPARISON OF THE PROPOSED METHOD AND OTHER METHODS. EACH VALUE REPRESENTS AN AVERAGED PERFORMANCE METRIC OVER ALL THE IMAGES IN THE DATASET

Method	Reference	Pixel	Spatially Based		Error Count			
		Dice	HD	NSD	Split	Merge	Add	Miss
Otsu's	[35]	0.88	30.6	0.11	1.10	2.40	0.3	5.60
(RC) Iterative Thresholding	[37]	0.92	34.8	0.12	1.10	2.40	0.3	5.50
Watershed (direct)	[24]	0.82	25.9	0.34	13.2	1.20	1.8	3.20
Watershed (gradient)	[38]	0.85	34.6	0.30	7.70	2.00	2.0	2.90
Merging	[39]	0.91	13.1	0.07	1.80	2.10	1.0	3.30
Active Mask	[40]	-	148.3	0.55	10.50	2.10	0.4	10.8
Template matching	[41]	-	77.8	0.06	0.58	1.45	0.9	3.48
Level set	[42]	-	96.6	0.09	1.10	0.35	2.75	0.85
K-means	[43]	-	94.6	0.11	1.56	0.30	2.6	1.60
Proposed		0.88	23.3	0.07	0.46	0.30	0.35	1.10

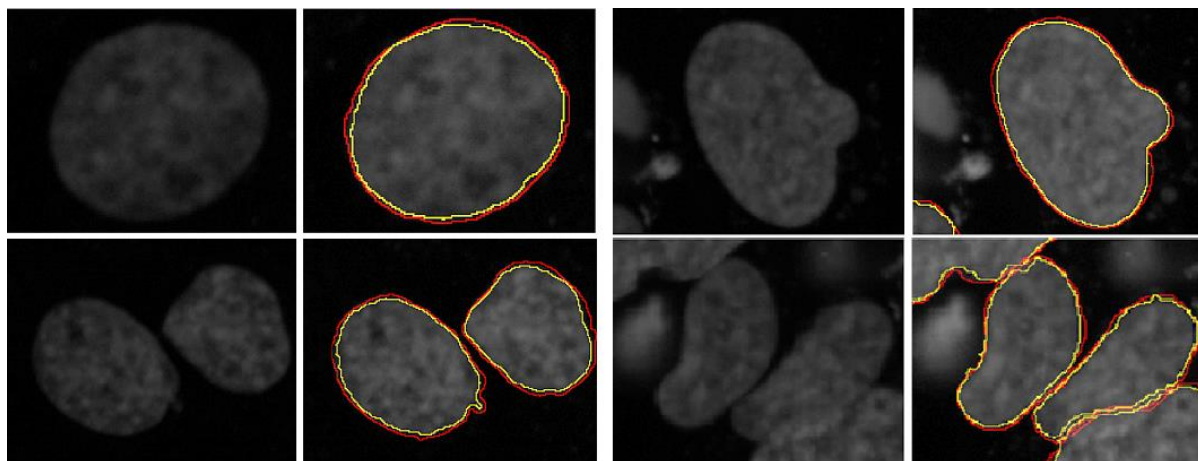


Fig. 11. Zoomed versions of Fig. 10. (left) original (right-red) ground truth (right-yellow) proposed.

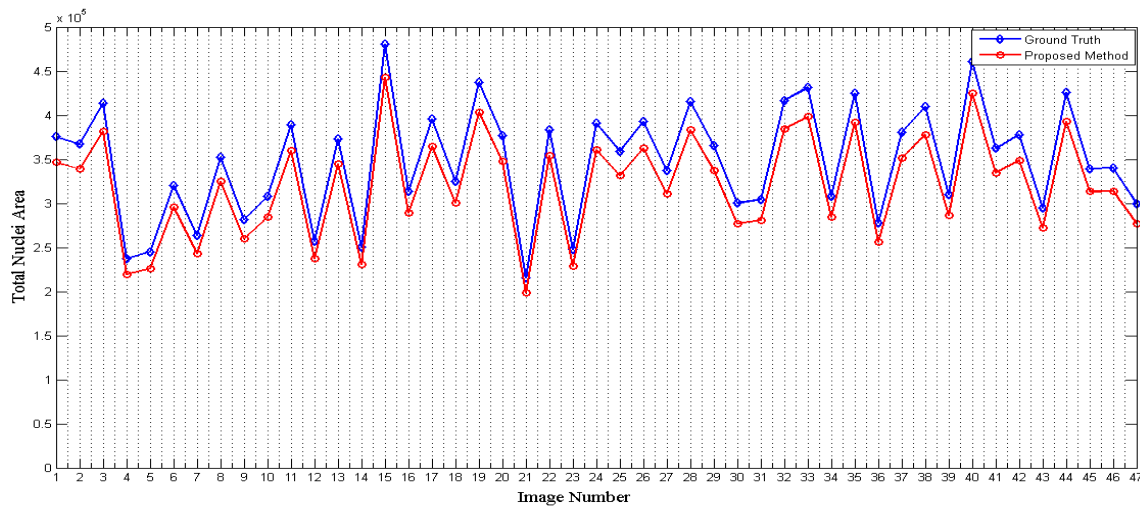


Fig. 12. Total area of nuclei detected in both ground truth reference images and the proposed methodology.

Good computational time cost of 15 seconds were achieved by [31], but using a (Linux system of 48GB RAM). It is still thought that the proposed algorithm dramatically outperforms this time because a (Windows system of 2G RAM) is used compared to 48GB RAM.

V. CONCLUSIONS AND FURTHER WORK

The proposed system in this work combined many image processing techniques in a single automated platform to automate processing of microscopy cell images. Different stages of the proposed system are supported by cascaded neural networks used to extract features and tune other processes. A GUI is also designed and provided to make the proposed system more user friendly and helpful in the forthcoming works. It was shown that the proposed system is capable of providing reasonable and very good promising results in the segmentation stages. It is hoped that the successful of this work and its subsequent development will pave the way for our vision of advanced levels of processing that includes real time processing of living cells and nuclei, and the three dimensional modelling of cells and histological structures. Future work to improve the outcome of the current work should include more accurate and efficient techniques for improving the nuclei segmentation. Another types of analysis such as automated cell detection, counting, classification, and tracking could then be built into a toolbox that would facilitate automation analysis of stem cell behavior as an example.

The proposed system is designed to deal with the whole image region without focusing into certain (ROI) regions which means more accurate and meaningful results. The proposed system depends on auto tuning of some related parameters which means that it can be extended to be applied on other datasets without the need for new methodologies. The proposed system shows a promising results compared to other systems. It also shows a rescannable reduction in processing time which makes it applicable in near real time diagnosis systems.

The ground truth images of the database under consideration are hand-segmented to separate touching nuclei

without labelling overlapped regions. So that, the separation stage described in the methodology section is designed to get accurate count of nuclei without taking into account the full area of each nucleus. As future work, it is hoped to enhance the refinement stage by increasing the accuracy of separation of clustered nuclei so that overlapped region is associated more accurately and more meaningfully to the separated nuclei. From the other hand, it is of the future planes to provide such system as a software as a service (SaaS) [44] and to allow the integration of this system with other related systems.

REFERENCES

- [1] R. C. Gonzalez and R. E. Woods, Digital image processing, 2002.
- [2] M. J. Aschwanden, "Image processing techniques and feature recognition in solar physics," Solar Physics, vol. 262, pp. 235-275, 2010.
- [3] M. Sonka, V. Hlavac, and R. Boyle, Image processing, analysis, and machine vision: Cengage Learning, 2014.
- [4] C. Solomon and T. Breckon, Fundamentals of Digital Image Processing: A practical approach with examples in Matlab: John Wiley & Sons, 2011.
- [5] T. W. Nattkemper, "Automatic segmentation of digital micrographs: A survey," Medical Informatics, vol. 11, pp. 847-851, 2004.
- [6] M. M. Fraz, P. Remagnino, A. Hoppe, B. Uyyanovvara, A. R. Rudnicka, C. G. Owen, and S. A. Barman, "Blood vessel segmentation methodologies in retinal images-A survey," Computer methods and programs in biomedicine, vol. 108, pp. 407-433, 2012.
- [7] S. Wienert, D. Heim, K. Saeger, A. Stenzinger, M. Beil, P. Hufnagl, M. Dietel, C. Denkert, and F. Klauschen, "Detection and Segmentation of Cell Nuclei in Virtual Microscopy Images: A Minimum-Model Approach," Sci. Rep., vol. 2, 07/11/online 2012.
- [8] N. Sharma and L. M. Aggarwal, "Automated medical image segmentation techniques," Journal of medical physics/Association of Medical Physicists of India, vol. 35, p. 3, 2010.
- [9] N. Sharma, A. K. Ray, S. Sharma, K. K. Shukla, S. Pradhan, and L. M. Aggarwal, "Segmentation and classification of medical images using texture-primitive features: Application of BAM-type artificial neural network," Journal of Medical Physics / Association of Medical Physicists of India, vol. 33, pp. 119-126, 2008.
- [10] M. J. McAuliffe, F. M. Lalonde, D. McGarry, W. Gandler, K. Csaky, and B. L. Trus, "Medical image processing, analysis and visualization in clinical research," in Computer-Based Medical Systems, 2001. CBMS 2001. Proceedings. 14th IEEE Symposium on, 2001, pp. 381-386.
- [11] I. Bankman, Handbook of medical image processing and analysis: academic press, 2008.

- [12] W. Birkfellner, Applied medical image processing: a basic course: Taylor & Francis, 2014.
- [13] J. Peti-Peterdi and P. D. Bell, "Confocal and two-photon microscopy," in Renal Disease, ed: Springer, 2003, pp. 129-138.
- [14] F. Gorunescu, "Data mining techniques in computer-aided diagnosis: non-invasive cancer detection," Journal of Gastrointestinal and liver diseases, vol. 16, pp. 427-430, 2007.
- [15] S. Kara and F. Dirgenali, "A system to diagnose atherosclerosis via wavelet transforms, principal component analysis and artificial neural networks," Expert Systems with Applications, vol. 32, pp. 632-640, 2007.
- [16] M. S. Sharif, R. Qahwaji, S. Hayajneh, S. Ipson, R. Alzubaidi, and A. Brahma, "An efficient system for preprocessing confocal corneal images for subsequent analysis," in Computational Intelligence (UKCI), 2014 14th UK Workshop on, 2014, pp. 1-8.
- [17] B. Salah, M. Alshraideh, R. Beidas, and F. Hayajneh, "Skin cancer recognition by using a neuro-fuzzy system," Cancer informatics, vol. 10, p. 1, 2011.
- [18] K. Doi, "Computer-aided diagnosis in medical imaging: historical review, current status and future potential," Computerized medical imaging and graphics, vol. 31, pp. 198-211, 2007.
- [19] C. SMOCHINĂ, P. HERGHELEGIU, and V. MANTA, "Image Processing Techniques used in Microscopic Image Segmentation," Gheorghe Asachi Technical University of Iași, 2011.
- [20] Q. Wu, F. Merchant, and K. Castleman, Microscope image processing: Academic press, 2010.
- [21] H. Lodish, A. Berk, S. L. Zipursky, P. Matsudaira, D. Baltimore, and J. Darnell, Microscopy and Cell Architecture. New York: W.H. Freeman & Company, 2000.
- [22] L. P. Coelho, A. Shariff, and R. F. Murphy, "Nuclear segmentation in microscope cell images: a hand-segmented dataset and comparison of algorithms," in Biomedical Imaging: From Nano to Macro, 2009. ISBI'09. IEEE International Symposium on, 2009, pp. 518-521.
- [23] V.-T. Ta, O. Lézoray, A. Elmoataz, and S. Schüpp, "Graph-based tools for microscopic cellular image segmentation," Pattern Recognition, vol. 42, pp. 1113-1125, 2009.
- [24] C. Wählby, I. M. SINTORN, F. Erlandsson, G. Borgefors, and E. Bengtsson, "Combining intensity, edge and shape information for 2D and 3D segmentation of cell nuclei in tissue sections," Journal of Microscopy, vol. 215, pp. 67-76, 2004.
- [25] X. Bai, C. Sun, and F. Zhou, "Splitting touching cells based on concave points and ellipse fitting," Pattern recognition, vol. 42, pp. 2434-2446, 2009.
- [26] O. Daněk, P. Matula, C. Ortiz-de-Solórzano, A. Muñoz-Barrutia, M. Maška, and M. Kozubek, "Segmentation of touching cell nuclei using a two-stage graph cut model," in Image Analysis, ed: Springer, 2009, pp. 410-419.
- [27] H. F. Yang and Y. Choe, "Cell tracking and segmentation in electron microscopy images using graph cuts," in IEEE International Symposium on Biomedical Imaging: From Nano to Macro, 2009. ISBI'09., 2009, pp. 306-309.
- [28] G. Srinivasa, M. C. Fickus, Y. Guo, A. D. Linstead, and J. Kovacevic, "Active mask segmentation of fluorescence microscope images," Image Processing, IEEE Transactions on, vol. 18, pp. 1817-1829, 2009.
- [29] M. R. JEONG, B. Ko, and J. Y. NAM, "Overlapping nuclei segmentation based on Bayesian networks and stepwise merging strategy," Journal of microscopy, vol. 235, pp. 188-198, 2009.
- [30] C. Chen, J. A. Ozolek, W. Wang, and G. K. Rohde, "A general system for automatic biomedical image segmentation using intensity neighborhoods," Journal of Biomedical Imaging, vol. 2011, p. 8, 2011.
- [31] J.-P. Bergeest and K. Rohr, "Efficient globally optimal segmentation of cells in fluorescence microscopy images using level sets and convex energy functionals," Medical Image Analysis, vol. 16, pp. 1436-1444, 10// 2012.
- [32] T. J. Fuchs and J. M. Buhmann, "Computational pathology: Challenges and promises for tissue analysis," Computerized Medical Imaging and Graphics, vol. 35, pp. 515-530, 10// 2011.
- [33] W. Wei, J. A. Ozolek, Slepč, x030C, D. ev, A. B. Lee, C. Cheng, and G. K. Rohde, "An Optimal Transportation Approach for Nuclear Structure-Based Pathology," Medical Imaging, IEEE Transactions on, vol. 30, pp. 621-631, 2011.
- [34] L. P. Coelho, T. Peng, and R. F. Murphy, "Quantifying the distribution of probes between subcellular locations using unsupervised pattern unmixing," Bioinformatics, vol. 26, pp. i7-i12, 2010.
- [35] N. Otsu, "A Threshold Selection Method from Gray-Level Histograms," IEEE Transactions on Systems, Man and Cybernetics, vol. 9, pp. 62-66, 1979.
- [36] J. S. Suri, Advances in diagnostic and therapeutic ultrasound imaging: Artech House, 2008.
- [37] T. Ridler and S. Calvard, "Picture thresholding using an iterative selection method," IEEE transactions on Systems, Man and Cybernetics, vol. 8, pp. 630-632, 1978.
- [38] L. Vincent and P. Soille, "Watersheds in digital spaces: an efficient algorithm based on immersion simulations," IEEE Transactions on Pattern Analysis & Machine Intelligence, pp. 583-598, 1991.
- [39] G. Lin, U. Adiga, K. Olson, J. F. Guzowski, C. A. Barnes, and B. Roysam, "A hybrid 3D watershed algorithm incorporating gradient cues and object models for automatic segmentation of nuclei in confocal image stacks," Cytometry Part A, vol. 56, pp. 23-36, 2003.
- [40] G. Srinivasa, M. Fickus, M. N. G. Rivero, S. Y. Hsieh, Y. Guo, A. D. Linstead, and J. Kovacevic, "Active mask segmentation for the cell-volume computation and Golgi-body segmentation of HeLa cell images," in Biomedical Imaging: From Nano to Macro, 2008. ISBI 2008. 5th IEEE International Symposium on, 2008, pp. 348-351.
- [41] C. Chen, W. Wang, J. A. Ozolek, and G. K. Rohde, "A flexible and robust approach for segmenting cell nuclei from 2D microscopy images using supervised learning and template matching," Cytometry Part A, vol. 83, pp. 495-507, 2013.
- [42] T. F. Chan and L. Vese, "Active contours without edges," Image processing, IEEE transactions on, vol. 10, pp. 266-277, 2001.
- [43] K. S. Ravichandran and B. Ananthi, "Color skin segmentation using K-means cluster," February 22 2016 2009.
- [44] S. Hayajneh, "Cloud Computing SaaS Paradigm for Efficient Modelling of Solar Features and Activities," International Journal of Cloud Applications and Computing (IJCAC), vol. 5, pp. 20-34, 2015.

An Automatic Segmentation Algorithm for Solar Filaments in H-Alpha Images using a Context-based Sliding Window

Ibrahim A. Atoum

College of Applied Sciences
Al Maarefa Colleges for Science and Technology
Riyadh, Saudi Arabia

Abstract—There are many features which appear on the surface of the sun. One of these features that appear clearly are the dark threads in the Hydrogen alpha ($H\alpha$) spectrum solar images. These ‘filaments’ are found to have a definite correlation with Coronal Mass Ejections (CMEs). A CME is a large release of plasma into space. It can be hazardous to astronauts and the spacecraft if it is being ejected towards the Earth. Knowing the exact attributes of solar filaments may open the way towards predicting the occurrence of CMEs. In this paper, an efficient and fully automated algorithm for solar filament segmentation without compromising accuracy is proposed. The algorithm uses some statistical measures to design the thresholding equations and it is written in the C++ programming language. The square root of the range as a measure of variability of image intensity values is used to determine the size of the sliding window at run time. There are many previous studies in this area, but no single segmentation method that could precisely claim to be fully automatic exists. Samples were taken from several representative regions in low-contrast and high-contrast solar images to verify the viability and efficacy of the method.

Keywords—Solar image processing; solar filament; segmentation; sliding window; Coronal mass ejections

I. INTRODUCTION

Solar filaments are huge regions of very thick relatively cool plasma compared to the surface of the sun. These solar features appear as elongated and dark filaments in $H\alpha$ solar images. CMEs are enormous bubbles of hot plasma (billions of tonnes of magnetized plasma) that is propagating away from the solar corona to the interplanetary medium at a very high velocity [1]. The importance of studying solar filaments comes from considering its disappearances as a significant indicator for the possible occurrence of CMEs, which is considered as the major cause of geomagnetic storms. It is now almost certain that there is a close association between CME and filament disappearances ([1]-[12]). Most of the filament detection techniques depend on sliding windows over the image region of interest. This sliding window is the number of pixels that are shifted while scanning the image. The scanning speed of images greatly depends on the size of the sliding window and the computational power of the image processing system. A context-based sliding window is proposed in this study to fully automate the process of solar filament segmentation in $H\alpha$ Images. The use of often user selectable, non-adaptable and non-automatic different sizes of

structuring elements in filament detection techniques led to them being defined as non-automated. According to this hypothesis, no single filament segmentation method thus could accurately claim to be fully automatic. Additionally, most of the previous studies ([1], [5], [7], [8], [13]-[20]) follow an image scanning process using a fixed size sliding window according to a fixed step size. Also the sliding window size affects the results of the different filtering techniques like the median filter [21]. The paper is organized as follows. The concept of the context-sensitive sliding window is presented in the next section. Section III overviews the algorithm of the context-based sliding window. Section Four is the conclusion.

Fig. 1 shows the unevenly illuminated solar images. These images were enhanced using an image enhancement technique introduced by [16].

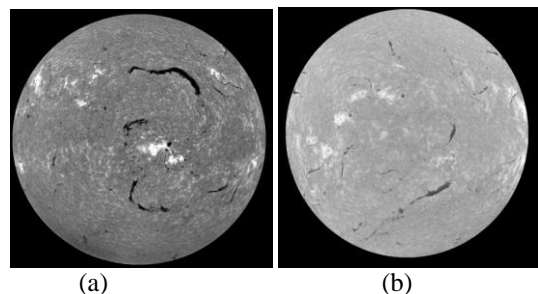


Fig. 1. (a) Solar image observed by Big Bear Solar Observatory (BBSO) on 2nd January 2001; (b) Solar image observed by Meudon Observatory on 1st April 2002.

II. CONTEXT-BASED SLIDING WINDOW

The common purpose of a sliding window is to give the possibility of checking the presence of the object in a rectangular box at different positions in the image. Subsequently, it is possible to use any of the many proposed methods to check the presence of an object in the sliding window or not. All previous methods of segmenting solar filaments rely entirely on customizing the statistical measures using a fixed sized sliding window. In general, the values of the statistical parameters in a local window are affected by the size of it. This research strongly proffers to adapt the size of the sliding window according to the window contents by using the Range (R) as a measure of variability, that is, the

difference between the lowest and highest image intensity values contained in the sliding window. To ascertain the effectiveness of using R and to optimize the image segmentation method, two samples; the first window (W_1) and the second window (W_2) were taken from two solar images as shown in the figures.

III. HOMOGENEITY AND HETEROGENEITY OF IMAGE REGIONS

To give an adequate opportunity for segmenting solar images and hence judge between the homogeneous and

heterogeneous regions, four values were calculated from different representative image samples, these being the value of R , the window side length; the window size (W_{SIZE}) as the integer square root of R ; the mean value (W_{μ}) and the standard deviation (W_{σ}). These values are shown in Fig. 2 and 3. Samples were taken to represent from non-filament regions (in order to represent them) as shown in Fig. 2(a) and Fig. 3(a). Other samples were taken from filament regions. Finally, the last samples were taken to represent the filament edges as shown in Fig. 2(c) and Fig. 3(c). Fig. 4 displays the same parameters for the low contrast image filaments.

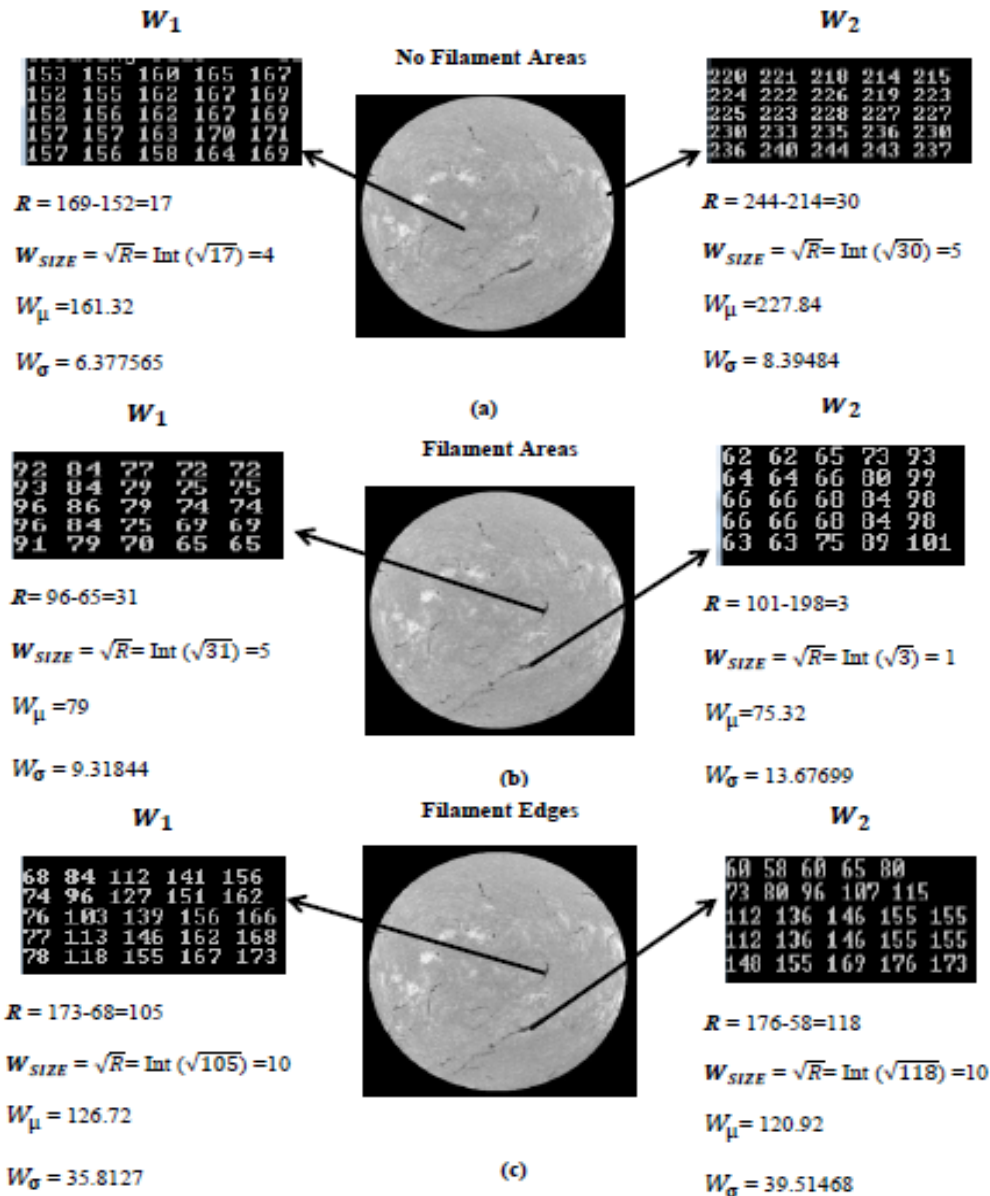


Fig. 2. Samples taken from solar image observed at Meudon Observatory on 1st April 2002: (a) Two samples from no filament pixels. (b) Two samples from filament pixels. (c) Two samples from Filament Edges.

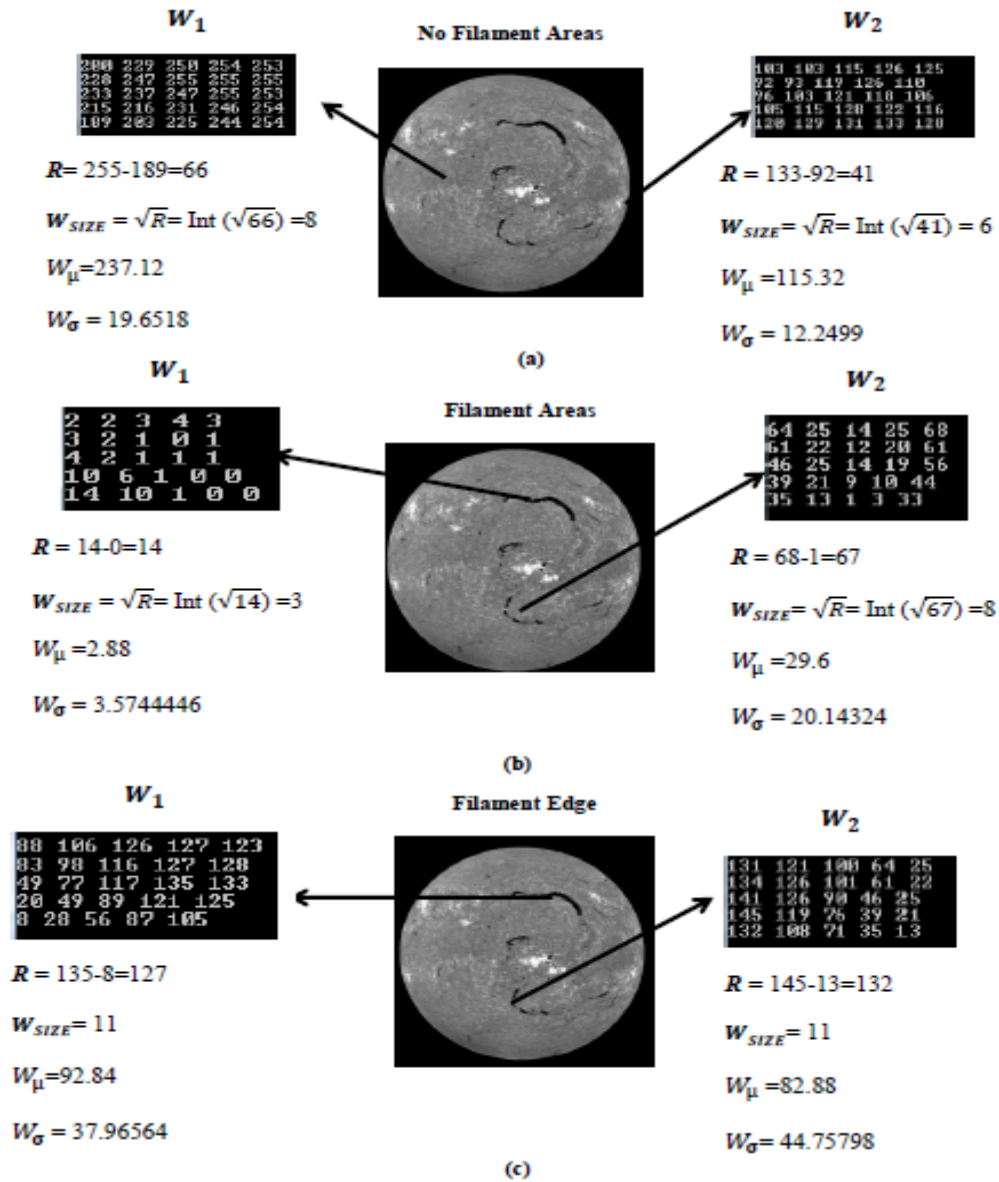


Fig. 3. Samples taken from solar image observed at Big Bear Solar Observatory 2nd January 2001: (a) Two samples from no-filament pixels. (b) Two samples from filament pixels. (c) Two samples from Filament Edges.

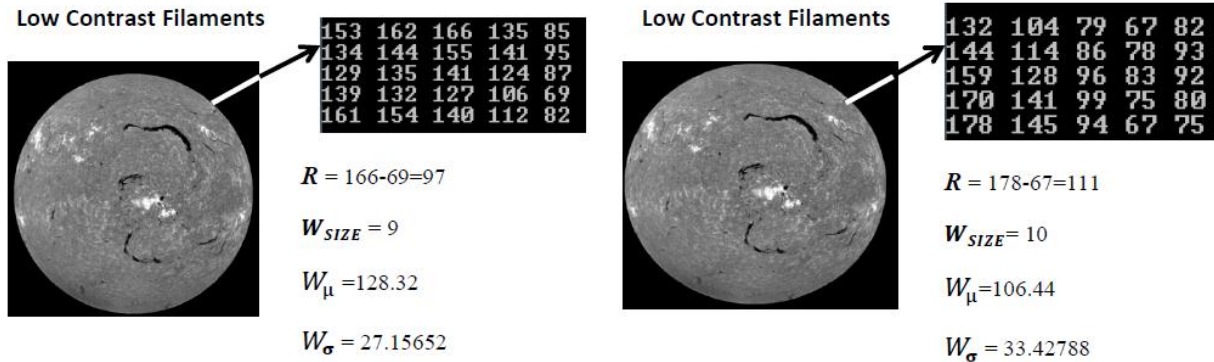


Fig. 4. Two samples from Low Contrast Filaments.

Where,

W_μ : the average for the window which is computed as $W_\mu = \frac{\sum_{i=1}^m W_i}{m}$, where W_i is the set of pixels inside the window and m being the number of pixels inside the window.

W_σ : the window standard deviation which is computed as $W_\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (W_i - W_\mu)^2}$, where W_i are each pixel value in the window and n being the number of pixels inside the window.

R : the window image intensity range value which is computed as $R = W_{max} - W_{min}$, where W_{max} and W_{min} are the maximum and minimum pixel image intensity values inside the window.

It can be observed from the previous values that the range values become high only in the regions with heterogeneous and non-convergent intensity values, while it will be low in the regions with relatively uniform (low spread) and convergent intensity values.

IV. ADAPTIVE FILAMENT SEGMENTATION

The process of image segmentation is the first and most important phase in analyzing and interpreting these images. The success or failure of the segmentation process may be considered as the success or failure of the subsequent image classes. There are many features that appear on Ha images but our goal through this algorithm is to separate the image into white foreground filaments and the black background.

Fig. 5 shows the results after applying the segmentation code as shown in Fig. 7 on the solar image observed by Big Bear Solar Observatory (BBSO) on 2nd January 2001. Fig. 6 shows the results after applying the segmentation code on the solar image observed by Meudon Observatory (MO) on 1st April 2002.

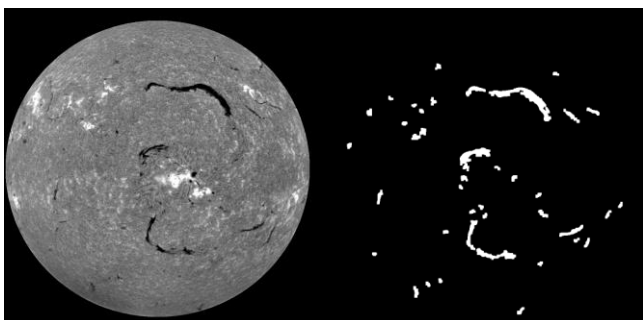


Fig. 5. The output after applying the algorithm on the BBSO image.

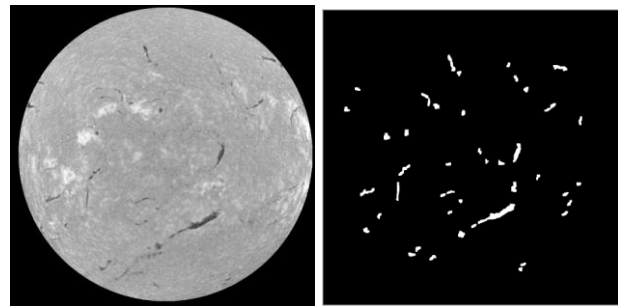


Fig. 6. The output after applying the algorithm on the MO image.

Setting the window size to any initial value

```
For i=1 to m do // m is the number of image rows
For j=1 to n do // n is the number of image columns
Calculate  $I_\mu$  and  $I_\sigma$ 
For n=1 to W do // W: window rows
For m=1 to W do
Calculate  $W_\mu$ ,  $W_R$  and  $W_\sigma$ 
If ((( $I_\mu > (W_\mu + I_\sigma)$ ) OR ( $W_\mu < I_\sigma$ ))) then
Filaments
Else
No Filaments
End if
End for m
End for n
End for j
End for i
```

Fig. 7. The segmentation pseudo code.

V. EVALUATION AND RESULTS

The size of the sliding window will change automatically after calculating the range in the first iteration based on the sliding window intensity values. The highest value that may exist in any sliding window is 255 (white) and the lowest one is 0 (black), which means that the maximum possible can only be, $255 - 0 = 255$. A low window size indicates a homogeneous region (such as no-filament or filament region) and a high value indicates a non-homogenous region (such as the regions of filament edges). In *Ha* solar images; the filaments appear as dark features. This means if the sliding window is around the edges points of the filament pixels then the range value will be high otherwise, it will be low. The values calculated for each sample in all the regions are summarized in Table I.

TABLE I. STATISTICAL DETAILS OF THE SOLAR IMAGE AND THE SLIDING WINDOW

		R		W _{SIZE}		W _μ		W _σ		I _μ	I _σ
		W ₁	W ₂	W ₁	W ₂	W ₁	W ₂	W ₁	W ₂		
Solar Image with High Intensity values	No Filament	17	30	4	5	161.32	227.84	6.377565	8.39484	172	15.7931
	High contrast Filament	3	31	1	5	75.32	79	13.67699	9.31844		
	Filament edge	118	105	10	10	120.92	126.72	39.51468	35.8127		
Solar Image with Low Intensity values	No Filament	41	66	6	8	115.32	237.12	12.2499	19.6518	124	24.6386
	Low contrast Filament	67	14	8	3	29.6	2.88	20.14324	3.5744446		
	Filament edge	132	127	11	11	82.88	92.84	44.75798	37.96564		
	High contrast Filament	97	111	9	10	128.32	106.44	27.15652	33.42788		

Where,

I_{μ} : The mean value for the whole image which is computed as $I_{\mu} = \frac{\sum_{i=1}^n P_i}{n}$, where P_i is the set of pixels and n being the number of pixels inside the solar image.

I_{σ} : The standard deviation for the whole image which is computed as

$$I_{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - I_{\mu})^2}$$

where n the number of pixels inside the whole image and x_i the pixel value.

Tables II and III contain the values of the different operators of (1) and (2) for all the regions. Tables I, II and III verify that (1) is correct for harmonic (relatively uniform) regions (with filaments or without filaments).

$$I_{\mu} > (W_{\mu} + I_{\sigma}) \tag{1}$$

However, it was observed that (1) does not work well with the areas that have low contrast filaments; so another condition is adapted to segment these filaments, as shown in (2) and as it is clearly seen in Table III.

$$W_{\mu} < (R + W_{\sigma}) \tag{2}$$

TABLE II. FIRST WINDOW: FIRST EQUATION VERIFICATION

I _μ	R	W _μ	W _σ	W _μ + W _σ	W ₁		R + W _σ	I _σ	W _μ + I _σ
					No Filament	Filament			
172	17	161.32	6.377565	167.697565	No Filament	Filament	23.377565	15.7931	177.1131
	3	75.32	13.67699	88.99699	Filament	Filament	16.67699		75.32
	118	120.92	39.51468	160.43468	Filament Edge	Filament	157.51468		120.92
124	41	115.32	12.2499	127.5699	No Filament	No Filament	53.2499	24.6386	139.9586
	67	29.6	20.14324	49.74324	Low contrast filament	Low contrast filament	87.14324		29.6
	132	82.88	44.75798	127.63798	Filament Edge	Filament Edge	176.75798		82.88
	97	128.32	27.15652	155.47652	High Contrast Filament	High Contrast Filament	124.15652		152.9586

TABLE III. SECOND WINDOW: FIRST EQUATION VERIFICATION

I _μ	R	W _μ	W _σ	W _μ + W _σ	W ₂		R + W _σ	I _σ	W _μ + I _σ
					No Filament	Filament			
172	30	227.84	8.39484	236.23484	No Filament	No Filament	38.39484	15.7931	243.6331
	31	79	9.31844	88.31844	Filament	Filament	40.31844		79
	105	126.72	35.8127	162.5327	Filament Edge	Filament Edge	140.8127		126.72
124	66	237.12	19.6518	256.7718	No Filament	No Filament	85.6518	24.6386	261.7586
	14	2.88	3.5744446	6.4544446	Low Contrast Filament	Low Contrast Filament	17.5744446		2.88
	127	92.84	37.96564	130.80564	Filament Edge	Filament Edge	164.96564		92.84
	111	106.44	33.42788	139.86788	High Contrast Filament	High Contrast Filament	144.42788		

VI. CONCLUSION

An accurate, robust and novel algorithm for solar filament segmentation has been introduced in this paper. The algorithm adopted an adaptive, automated sliding window size according to the extent of the heterogeneity of the window pixels instead of using a fixed value. The algorithm uses extracted windows from two different solar images, one high contrast image and one low contrast image. Some statistical calculations were used to judge the presence of the solar filament or not. The range as a statistical measure played an important role in

automating the segmentation process. The mean value and standard deviation of the image intensity values in addition to the range, mean value and standard deviation of the adaptive sliding window pixels have been used to make up the filament segmentation equations. This technique can be better evaluated in the subsequent solar image classes which include merging broken structures; filaments characterizations and detecting filament disappearances; additionally it can be the basis for all these subsequent automatic operations.

REFERENCES

- [1] L. Alejandro, "The source region of coronal mass ejections," *The Astrophysical Journal*, vol. 688, pp. 647-655, November 2008.
- [2] N. Gopalswamy, R. Kundu, K. Manoharan, A. Raouf, N. Nitta, and P. Zarka, "X-ray and radio studies of a coronal eruption: shock wave, plasmoid, and coronal mass ejection," *The Astrophysical Journal*, vol. 486, pp. 1036-1044, September 1997.
- [3] P. Subramanian and K. Dere, "Source regions of coronal mass ejections," *The Astrophysical Journal*, vol. 561, pp. 372-395, July 2001.
- [4] M. Al-Omari, R. Qahwaji, T. Colak, and S. Ipson, "Machine learning-based investigation of the associations between CMEs and filaments," *Solar Phys.*, vol. 262, pp. 511-539, April 2010.
- [5] N. Fuller, N., J. Aboudarham, and D. Bentley, "Filament recognition and image cleaning on Meudon H α spectroheliograms," *Solar Phys.*, vol. 227, pp.61-73, March 2005.
- [6] N. Gopalswamy, "Coronal mass ejections of solar cycle 23," *J. Astrophys. Astr.*, vol. 27, pp. 243-254, June 2006.
- [7] N. Gopalswamy, et. al., "The *SOHO/LASCO CME catalog*," *Earth, Moon, and Planets*, vol. 104, pp. 295-313, April 2009
- [8] J. Jing, B. Yurchyshyn, G. Yang, Y. Xu and H. Wang, "On the relation between filament eruptions, flares and coronal mass ejections," *Astrophysical Journal*, vol. 614, pp. 1054-1062, October 2004.
- [9] S. Pojoga and S. Huang, "On the sudden disappearances of solar filaments and their relationship with coronal mass ejections," *Adv. Space Res.*, vol. 32, pp. 2641-2646, December 2003.
- [10] E. Robbrecht, Berghmans, D. and R. Van der Linden, "Automated LASCO CME catalog for solar cycle 23: sre CMEs scale invariant," *The Astrophysical Journal*, vol. 691, pp. 1222-1234, February 2009.
- [11] B. Schmieder, "Magnetic source regions of coronal mass ejections," *J. Astrophys. Astr.* Vol. 27, pp. 139-149, June 2006.
- [12] P. Zhou, X. Wang, J. Zhang, "Large-scale source regions of earth-directed coronal mass ejections," *Astronomy and Astrophysics*, vol.445, pp.1133-1141, January 2006.
- [13] J. Aboudarham, et. al., "Automatic detection and tracking of filaments for a solar feature database," *Annales Geophysicae*. Vol. 26, pp. 243-248. February 2008.
- [14] A. Vourlidas, R. A. Howard, E. Esfandiari, S. Patsourakos, S. Yashiro, and G. Michalek, "Comprehensive analysis of coronal mass ejection mass and energy properties over a full solar cycle," *The Astrophysical Journal*, vol. 722, pp. 1522-1538, September 2010.
- [15] P. Bernasconi, D. Rust and D. Hakim, "Advanced automated solar filament detection and characterization code: description, performance and results," *Solar Phys.*, vol. 228, pp. 97-119, January 2005
- [16] H. Fang and P.Chen, "Developing an advanced automated method for solar filament recognition and its scientific application to a solar cycle of MLSO H α data," *Sol. Phys.*, vol.286, pp. 385-404, March 2013.
- [17] A. Joshi, N. Srivastava and S. Mathew, "Automated detection of filaments and their disappearance using full-disc H α images," *Sol. Phys.*, vol. 262, pp. 425-436, April 2010.
- [18] E. Robbrecht and D. Berghmans, "A broad perspective on automated CME tracking: towards higher level space weather forecasting," *Geophysical Monograph Series*, vol. 165, pp. 33-42, October 2006.
- [19] Y. Yuan Y., F. Shih, J. Jing, H. Wang H. and J. Chae, "Automatic solar filament segmentation and characterization," *Sol. Phys.*, vol. 272, pp. 101-117, August 2011.
- [20] F. Shih and J. Kowalski, "Automatic extraction of filaments in H-alpha solar images," *Solar Phys.* vol. 218, pp. 99-122, December 2003.
- [21] E. Robbrecht and D. Berghmans, "Automated recognition of coronal mass ejections (CMEs) in near-real-time data," *Astronomy and Astrophysics*, vol. 425, pp. 1097-1106, June 2004.

Relative Humidity Profile Estimation Method with AIRS (Atmospheric Infrared Sounder) Data by Means of SDM (Steepest Descend Method) with the Initial Value Derived from Linear Estimation

Kohei Arai¹

¹Graduate School of Science and Engineering
Saga University
Saga City, Japan

Abstract—Relative humidity profile estimation method with AIRS (Atmospheric Infrared Sounder) data by means of SDM (Steepest Descend Method) with the initial value derived from LED: Linear Estimation Method is also proposed. Through experiments, it is found that there is almost 15 (%) of relative humidity estimation error. Therefore, it can be said that the relative humidity is still tough issue for retrieval. It is also found that the estimation error does not depend on the designated atmospheric models, Mid-Latitude Summer/Winter, Tropic. Even if the assigned atmospheric model is not correct, the proposed SDM based method allows almost same estimated relative humidity. In other word, it is robust against atmospheric model.

Keywords—Atmospheric Infrared Sounder (AIRS); Steepest Descend Method (SDM); LED; MODTRAN; relative humidity; atmospheric model; Infrared sounder

I. INTRODUCTION

Vertical profiles of air-temperature and relative humidity (Water vapor) can be estimated with satellite based Infrared Sounder data [1]. Retrieval accuracy, however, is not good enough for Earth system science. Estimation accuracy of air-temperature and relative humidity at tropopause¹ altitude is not good enough, in particular, because there are gradient changes of air-temperature and relative humidity profile in the tropopause. Therefore, observed radiance at the specific channels is not changed for the altitude.

There is least square based conventional method for estimation of vertical profiles of air-temperature and relative humidity. In the method, Root Mean Square (RMS) difference between observed radiance and calculated radiance (model based radiance) with the designated physical parameters in the model is minimized. Then, the designated physical parameters including air-temperature and relative humidity at the minimum RMS difference are to be solutions.

The most typical least square method is Newton-Raphson method² which gives one of local minima. Newton-Raphson method needs the first and the second order derivatives,

Jacobian and Hessian at around the current solution. It is not easy to formularize these derivatives analytically. The proposed method is based on Levenberg Marquardt (LM)³ of non-linear least square method. It uses numerically calculated the first and the second order derivatives instead of analytical based derivatives, namely, these derivatives can be calculated with radiative transfer model based radiance calculations. At around the current solution in the solution space, directional derivatives are calculated with the radiative transfer model such as MODTRAN.

The proposed method is validated for air-temperature and relative humidity profile retrievals with Infrared: IR sounder⁴ data obtained from AQUA/AIRS (AIRS instrument onboard AQUA satellite [2]-[7]) by mean of Steepest Descent Method: SDM. Although SDM based optimization method shows relatively good accuracy, it takes a comparatively large computer resources and it is falling in local minima, not in the global optimum. To avoid this situation of which the solution is falling in local minima and accelerate the convergence through giving initial value which is derived from a Linear Estimation Method (LEM). Due to the fact that retrieval accuracy depends on the initial value, the most appropriate initial value is given from the LEM in the proposed method.

From the previous research works, a comparison of retrieving accuracy between Newton-Raphson method and the proposed method based on LM method [8] is made to demonstrate an effectiveness of the proposed method in terms of estimation accuracy for the altitude of tropopause [9]. Global Data Assimilation System: GDAS⁵ data of assimilation model derived 1-degree mesh data is used as truth data of air-temperature and relative humidity profiles. The experimental data show that the proposed method is superior to the conventional Newton-Raphson method.

The next section describes the proposed method together with typical conventional Newton-Raphson method for retrieving vertical profiles followed by experiments. Then

¹ <http://en.wikipedia.org/wiki/Tropopause>

² http://en.wikipedia.org/wiki/Newton's_method

³ http://en.wikipedia.org/wiki/Levenberg%E2%80%93Marquardt_algorithm

⁴ http://en.wikipedia.org/wiki/Atmospheric_Infrared_Sounder

⁵ <http://www.mmm.ucar.edu/mm5/mm5v3/data/gdas.html>

conclusion with some discussions is followed by together with future research works.

II. PROPOSED METHOD

A. Infrared Sounder

Fig. 1 shows atmospheric transmittance and transmittance of carbon dioxide and water vapor. Using absorption bands of molecules, carbon dioxide and water vapor content in the atmosphere can be retrieved. There is strong relation between air-temperature and carbon dioxide content. Meanwhile, water vapor content is a closely related to relative humidity. Therefore, air-temperature and relative humidity can be estimated using absorption characteristics of carbon dioxide and water vapor contents in the atmosphere.

Using slope characteristics of the absorption characteristics, vertical profile of carbon dioxide and water vapor can also be estimated. Optical depths⁶ of molecules are different and depend on the wavelength at the slope of the absorption characteristics. First derivative of optical depth against altitude is called weighting function.⁷ Therefore, weighting function is a function of altitude.

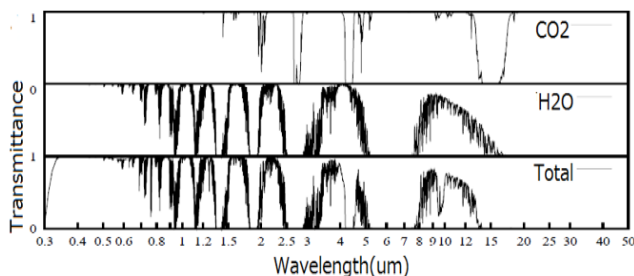


Fig. 1. Atmospheric transmittance and transmittance of carbon dioxide and relative humidity.

Fig. 2(a) shows an example of a transmittance as a function of altitude while (b) shows an example of a weighting function.

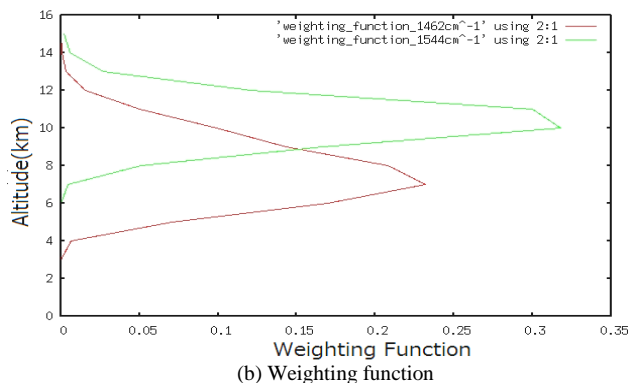
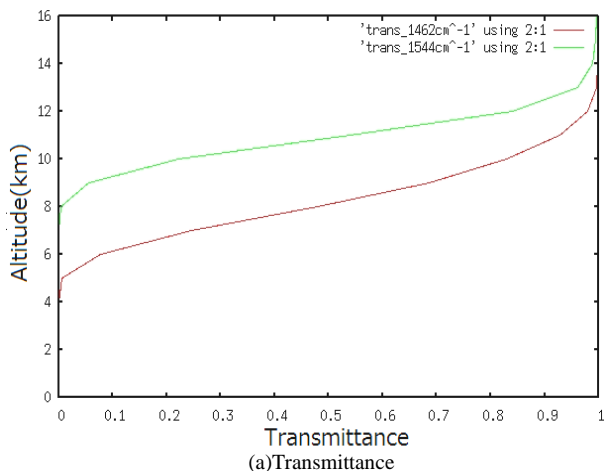


Fig. 2. Transmittance as a function of altitude and weighting function.

Mission instruments which allow vertical profile retrievals in the infrared wavelength region onboard satellite are called as infrared sounder. There are some infrared sounders such as TIROS-N / NOAA TOVS- High Resolution Infrared Sounder model 2 (HIRS/2)⁸. One of infrared sounders which are widely used with high quality of mission instruments is AIRS⁹: Atmospheric Infrared Sounder. Although AIRS data quality is high enough, retrieval accuracy is not good enough. Estimation accuracy of air-temperature and relative humidity at around tropopause is not good enough.

B. Conventional Retrieval Method

The previously proposed method for air-temperature and relative humidity retrievals is intended to improve retrieval accuracy. Typically, the conventional retrieval method is minimizing formulated error covariance using Newton-Raphson method. It, however, gives one of local minima sometime. Although the simulated annealing method gives global optimum solution, it takes huge computational resources.

The proposed method is based on non-linear least square method of Levenberg Marquardt: LM. LM method minimizes the difference between AIRS data derived radiance, R_{0i} and estimated radiance, R_i based on atmospheric model with the parameters of air-temperature and relative humidity.

$$S = \sum_{i=0}^n (R_i - R_{0i})^2 \quad (1)$$

Thus, the geophysical parameter of air temperature and water vapor (relative humidity) at the difference of radiance reaches to the minimum is estimated. Widely used, reliable and accurate enough atmospheric software code of MODTRAN¹⁰ which allows calculate at sensor radiance with a variety of parameters is used in the proposed method. Solution update equation of Newton-Raphson method is expressed by (2).

$$X_{n+1} = X_n - H^{-1}J(X_n) \quad (2)$$

⁶ http://en.wikipedia.org/wiki/Optical_depth

⁷ http://www.ecmwf.int/newsevents/training/rcourse_notes/DATA_ASSIMILATION/INVERSION_METHODS/inversion_methods2.html

⁸ www.ozone.noaa.gov/action/tovs.htm

⁹ http://aqua.nasa.gov/about/instrument_airs.php

¹⁰ <http://en.wikipedia.org/wiki/MODTRAN>

where H denotes Hessian matrix¹¹ which consists of the second order derivatives (residual square error, S which is represented by (1) by geophysical parameter, air-temperature and water vapor (relative humidity)). Also, J denotes Jacobean¹² which consists of the first derivative of vectors (S by geophysical parameters of air temperature and water vapor (relative humidity), x). On the other hand, solution update equation is expressed by (3).

$$X_{n+1} = X_n + (J^T J + I)^{-1} J^T (R_i - R_{0i}) \quad (3)$$

The first derivative is represented by (4) while the second order derivative is expressed by (5), respectively.

$$\frac{\partial S}{\partial x_i} = -2 \sum_{k=1}^n (R_k - R_{0k}) \frac{\partial R_{0k}}{\partial x_i} \quad (4)$$

$$\frac{\partial^2 S}{\partial x_i \partial x_j} = 2 \sum_{k=1}^n \left[\frac{\partial R_{0k}}{\partial x_i} \frac{\partial R_{0k}}{\partial x_j} - (R_k - R_{0k}) \frac{\partial^2 R_{0k}}{\partial x_i \partial x_j} \right] \quad (5)$$

On the other hand, the first and second order derivatives of R with x are expressed by (6) and (7), respectively.

$$\frac{\partial R_{0k}}{\partial x_i} \leftarrow \text{MODTRAN} \quad (6)$$

$$\frac{\partial^2 R_{0k}}{\partial x_i \partial x_j} \leftarrow \frac{\partial R_{0k}}{\partial x_i} * \frac{\partial R_{0k}}{\partial x_j} \quad (7)$$

Equation (6) can be derived from MODTRAN. Thus, the geophysical parameter, air temperature and water vapor relative humidity at the difference of radiance reaches to the minimum is estimated.

In the previously proposed method, these derivatives are calculated numerically. In order to determine the derivatives, 2% changes of relative humidity are considered for calculation of derivative R and S while 0.5K changes of air-temperature is also considered. Thus, the geophysical parameter of air temperature and water vapor (relative humidity) at when the difference of radiance reaches to the minimum is estimated.

C. Proposed Retrieval Method

Radiative transfer equation is expressed as follows:

$$Rv = (I_0)_v \tau_v(z_0) + \int_{z_0}^{\infty} Bv \{T(z)\} K_v(z) dz \quad (8)$$

where R, I, τ, z, B, T, K denotes the top of the atmosphere radiance (at sensor radiance), extinction component from the atmosphere, transparency of the atmosphere, altitude, brightness temperature of the ground surface, surface temperature, extinction function, respectively. This can be simplified and rewritten as follows:

$$R = BK \quad (9)$$

Then, R can be minimized as (10) through changing atmospheric transparency which is closely related to relative humidity content in the atmosphere which results in relative humidity profile retrieval.

$$R - R_0 = \frac{\partial R}{\partial q} (q - q_0) \quad (10)$$

In this process, R matrix of radiative brightness temperature is square matrix. Therefore, it can be solved relatively easily.

The estimated matrix can be expressed as follows:

$$\hat{x} = x_a + (A^T S_\varepsilon^{-1} A + S_a^{-1})^{-1} A^T S_\varepsilon^{-1} (R - R_a) \quad (11)$$

where x_a, A, S, R_a denotes designated matrix, Jacobean matrix which is expressed with (12), measuring error covariance matrix which is expressed with (13), and designated radiative brightness temperature.

$$A = \begin{pmatrix} \frac{\partial R_1}{\partial q_1} & \frac{\partial R_1}{\partial q_2} & \frac{\partial R_1}{\partial q_n} \\ \frac{\partial R_2}{\partial q_1} & \frac{\partial R_2}{\partial q_2} & \frac{\partial R_2}{\partial q_n} \\ \frac{\partial R_n}{\partial q_1} & \frac{\partial R_n}{\partial q_2} & \frac{\partial R_n}{\partial q_n} \end{pmatrix} \quad (12)$$

$$S_{ij} = \varepsilon (x_i - \hat{x}_i)(x_j - \hat{x}_j)^T \quad (13)$$

Then, optimum solution can be obtained by Steepest Descent: Method: SDM method as follows:

$$q_k = q_{k-1} + \alpha_k g_k \quad (14)$$

where α, g denotes step size and direction of the next solution, respectively. Process flow of the SDM method can be represented in Fig. 3.

Initial value is very important. If the initial value is close to the global optimum solution, SDM allows to reach the global optimum solution easily while the initial value is far from the global optimum solution and is close to one of a plenty of local minima, then SDM reaches a local minima easily. In the proposed method, the initial value is given by Lear Estimation Method: LEM. Generally, it is considered the LEM gives a solution which is situated at an appropriate solution space not far from the global optimum solution. Therefore, it is expected that the proposed method allows much appropriate solution rather than the SDM with arbitrary initial value not from the LEM in terms of residual error and computational resources.

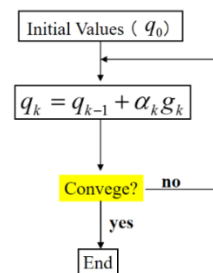


Fig. 3. Process flow of the SDM method.

¹¹ http://en.wikipedia.org/wiki/Hessian_matrix

¹² http://en.wikipedia.org/wiki/Jacobian_matrix_and_determinant

III. EXPERIMENT

A. Preliminary Experiment

AIRS data of the Mexican gulf at latitude of 20 degree North and longitude of 49-degree West which is acquired at 18:00 UTC on November 16, 2002 is used for the experiment. The intensive study area is shown in Fig. 4.

Also, Mid-Latitude Winter model of MODTRAN is used with the standard relative humidity, ozone and Carbone

dioxide together with the standard water vapor and air-temperature profiles because the Mexican gulf is situated in the Mid-Latitude region.

From the AIRS data, 9 channels of data are selected as shown in Table I. The top of the atmosphere brightness temperature (K) in the wave number range of 1460-1620 cm^{-1} is shown in Fig. 5.



Fig. 4. Intensive study area.

TABLE I. SELECTED AIRS CHANNELS TOGETHER WITH THEIR WAVE NUMBER

WAVE NUMBER	1478	1483	1508	1514	1519	1541	1544	1558	1585
AIRS CHANNEL	1684	1692	1731	1740	1748	1761	1765	1786	1824

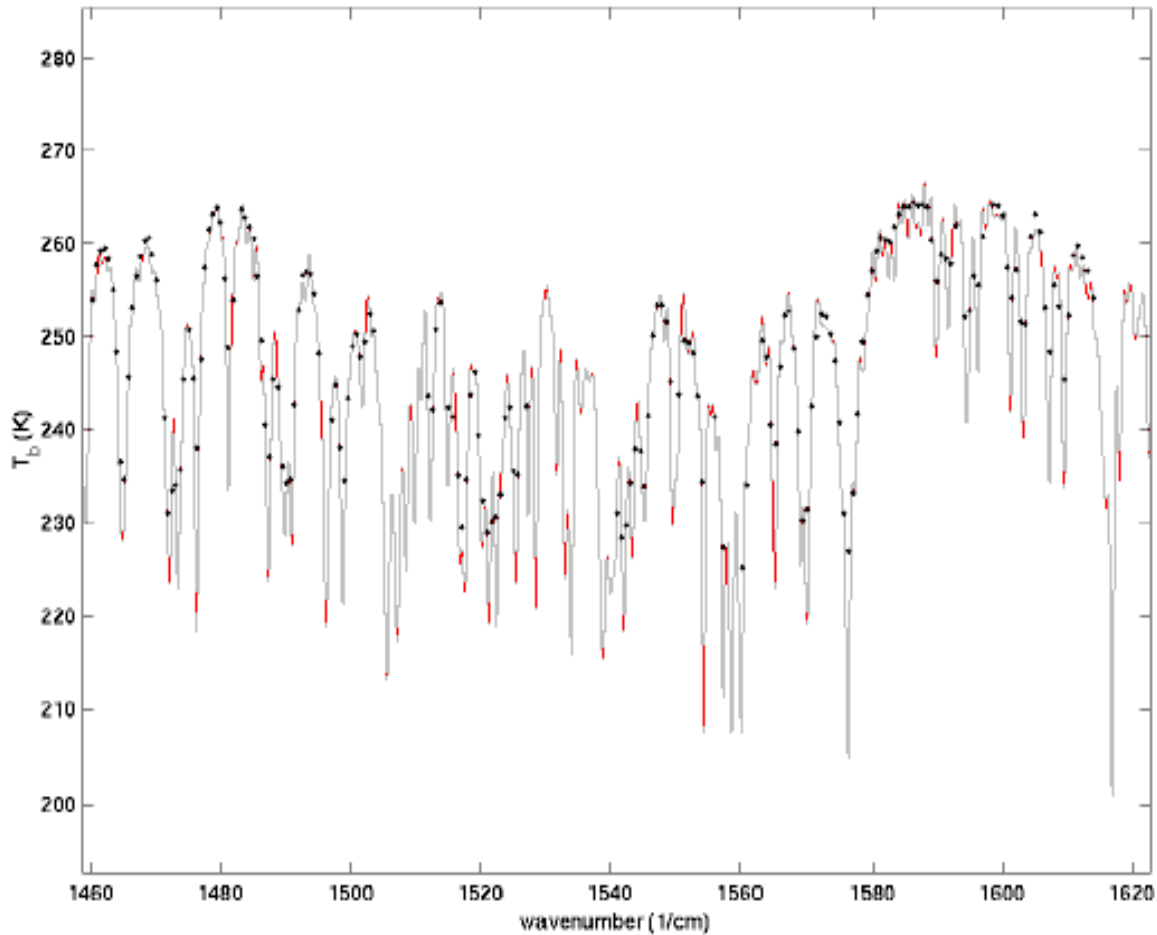


Fig. 5. The top of the atmosphere brightness temperature (K) in the wave number range of 1460-1620 cm^{-1} .

The peaks of the weighting functions imply that the most sensitive altitudes to the relative humidity profile.

By using these brightness temperature and weighting functions, relative humidity profile can be estimated based on the proposed method with the SDM method. The procedure of the proposed relative humidity profile estimation is shown in Fig. 6. Meanwhile, weighting functions of these channels are shown in Fig. 7.

The most sensitive wave number of AIRS channels are selected. Then the top of the atmosphere brightness temperature is calculated with MODTRAN. Starting with the initial point RH_0 which is estimated by using the method based on linear model of problem solving, the best matched brightness temperature is to be found by changing relative humidity profile.

First, wave number together with the corresponding weighting function is selected. Then, the top of the atmosphere brightness temperature is calculated using MODTRAN. The best relative humidity profile is estimated by minimizing the difference between observed brightness

temperature and the MODTRAN derived brightness temperature using the SDM method.

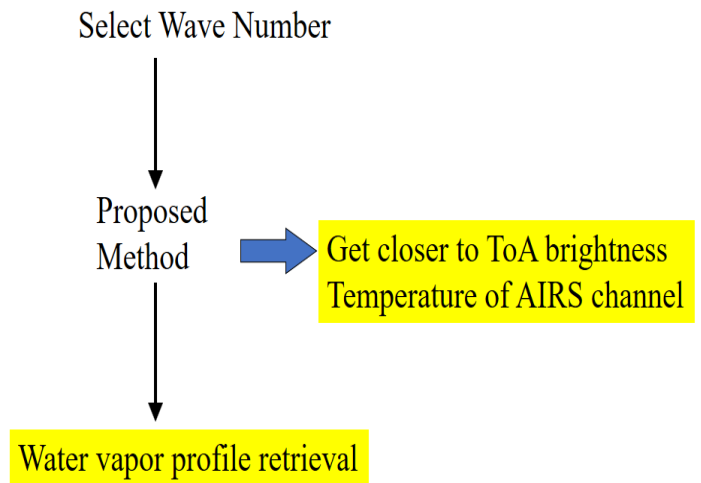


Fig. 6. Procedure of the proposed relative humidity profile estimation method.

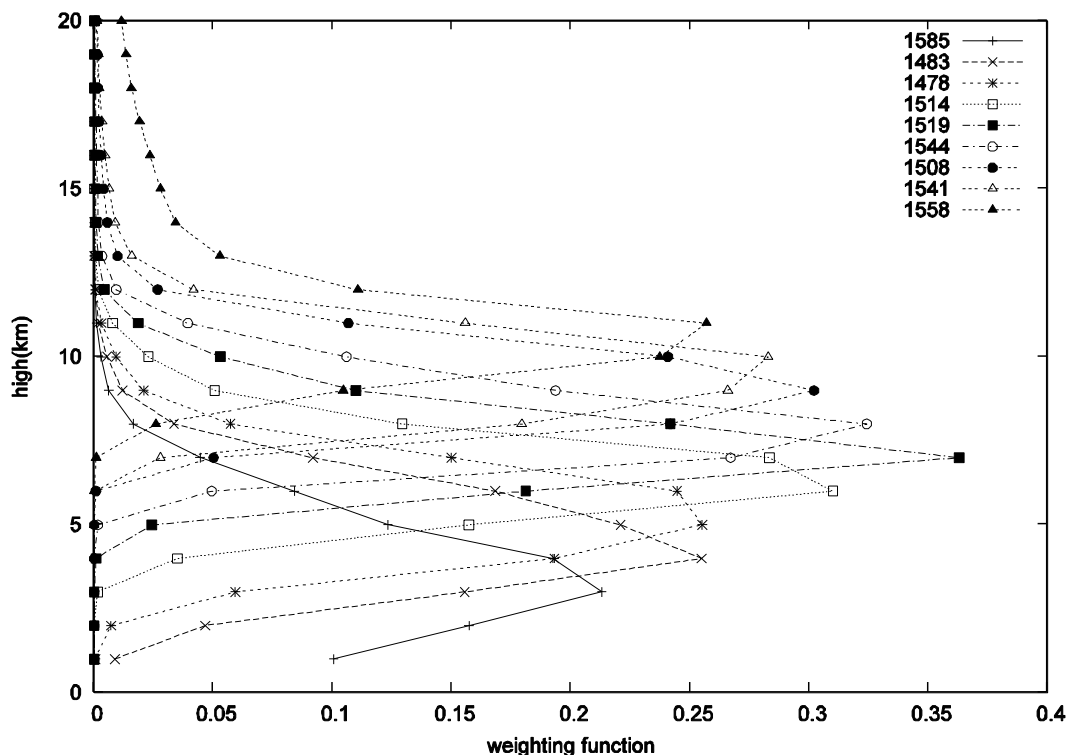


Fig. 7. Weighting functions of the selected 9 wave number channels

B. Retrieved Result

The retrieved result of relative humidity profile based on the proposed method is shown in Fig. 8. Meanwhile, the difference between calculated and acquired brightness temperature is shown in Fig. 9. There are three retrieved results in the Fig. 8, Mid-Latitude Winter, Mid-Latitude Summer and Tropic atmosphere models of MODTRAN. The AIRS data used is acquired in November 2002. Therefore, Mid-Latitude Winter would be better for relative humidity

profile estimation. It, however, is not true in some cases as shown in Fig. 8. Trend of the retrieved relative humidity profile with Mid-Latitude Winter coincide to the MODTRAN derived profile.

Meanwhile, the difference between calculated and estimated brightness temperature shows different trends depending on the wave number. It is not always the difference decreases in accordance with increasing of altitude.

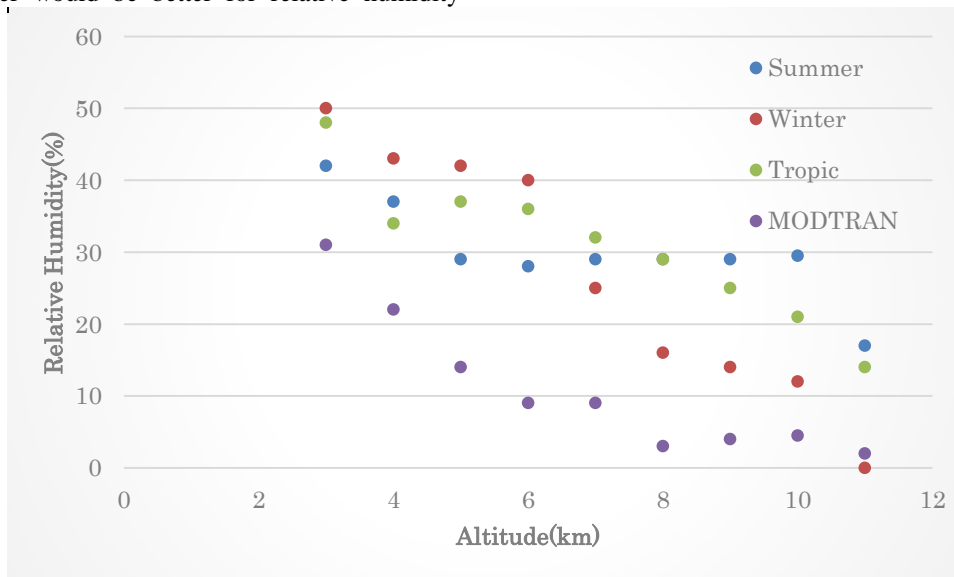


Fig. 8. Estimated relative humidity profile.

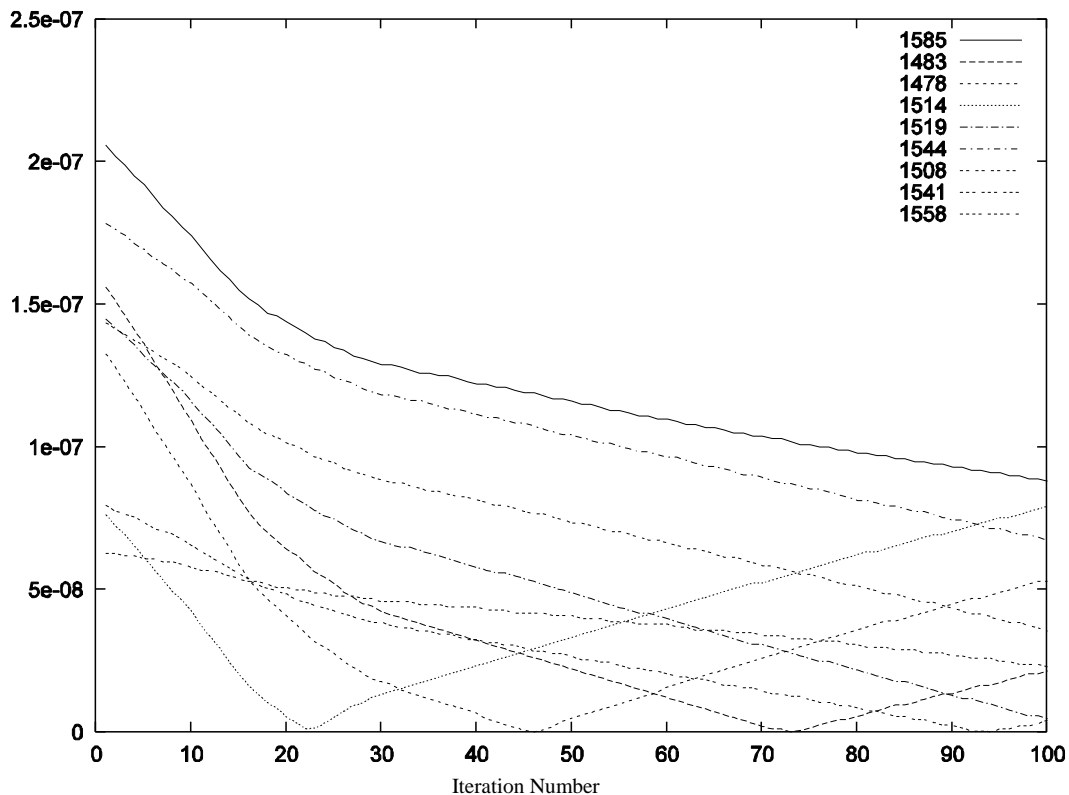


Fig. 9. Difference between calculated and acquired brightness temperature (K).

It is found that there is almost 15 (%) of relative humidity estimation error. Therefore, it can be said that the relative humidity is still tough issue for retrieval. It is also found that the estimation error does not depend on the designated atmospheric models, Mid-Latitude Summer/Winter, Tropic. Even if the assigned atmospheric model is not correct, the proposed SDM based method allows almost same estimated relative humidity. In other word, it is robust against atmospheric model.

VI. CONCLUSION

Non-linear optimization method based on Levenberg-Marquardt of non-linear least square method with numerically calculate deviation of Jacobean and Hessian is proposed for relative humidity and air-temperature profile estimations with infrared sounder data. Through comparisons with the conventional Newton-Raphson method, it is confirmed that the proposed method is superior to the conventional Newton-Raphson method in terms of retrieving accuracy. Also, the proposed method does not need Jacobean and Hessian because first and second order derivatives are calculated numerically with MODTRAN code based radiative transfer model.

It is found that there is almost 15 (%) of relative humidity estimation error. Therefore, it can be said that the relative humidity is still tough issue for retrieval. It is also found that the estimation error does not depend on the designated atmospheric models, Mid-Latitude Summer/Winter, Tropic. Even if the assigned atmospheric model is not correct, the proposed SDM based method allows almost same estimated

relative humidity. In other word, it is robust against atmospheric model.

Further investigation is required for improvement of relative humidity retrieval accuracy.

ACKNOWLEDGMENT

The author would like to thank Mr. Noriaki Yamada of Saga University for his effort to conduct the experiments.

REFERENCES

- [1] Kohei Arai, Lecture Note on Remote Sensing, Morikita-Shuppan publishing Co. Ltd, 2004.
- [2] NASA/JPL, "AIRS Overview". NASA/JPL. <http://airs.jpl.nasa.gov/overview/overview/>.
- [3] NASA "Aqua and the A-Train". NASA. http://www.nasa.gov/mission_pages/aqua/.
- [4] NASA/GSFC "NASA Goddard Earth Sciences Data and Information Services Center". NASA/GSFC. http://disc.gsfc.nasa.gov/AIRS/data_products.shtml.
- [5] NASA/JPL "How AIRS Works". NASA/JPL. http://airs.jpl.nasa.gov/technology/how_AIRS_works.
- [6] NASA/JPL "NASA/NOAA Announce Major Weather Forecasting Advancement". NASA/JPL. <http://jpl.nasa.gov/news/news.cfm?release=2005-137>.
- [7] NASA/JPL "New NASA AIRS Data to Aid Weather, Climate Research". NASA/JPL. <http://www.jpl.nasa.gov/news/features.cfm?feature=1424>.
- [8] Kohei Arai and Naohisa Nakamizo, Relative humidity and air-temperature profile estimation with AIRS data based on Levenberg - Marquardt, Abstract of the 50th COSPAR(Committee on Space Research/ICSU) Congress, A 3.1-0086-08,995, Montreal, Canada, July, 2008

- [9] Kohei Arai and XingMing Liang, sensitivity analysis for air temperature profile estimation method around the tropopause using simulated AQUA/AIRS data, *Advances in Space Research*, 43, 3, 845-851, 2009.

AUTHOR'S PROFILE

Kohei Arai, He received BS, MS and PhD degrees in 1972, 1974 and 1982, respectively. He was with The Institute for Industrial Science and Technology of the University of Tokyo from April 1974 to December 1978 also was with National Space Development Agency of Japan from January, 1979 to March, 1990. During from 1985 to 1987, he was with Canada Centre for Remote Sensing as a Post Doctoral Fellow of National Science and Engineering Research Council of Canada. He moved to Saga University as a

Professor in Department of Information Science on April 1990. He was a councilor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology during from 1998 to 2000. He was a councilor of Saga University for 2002 and 2003. He also was an executive councilor for the Remote Sensing Society of Japan for 2003 to 2005. He is an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Science Commission "A" of ICSU/COSPAR since 2008 then he is now award committee member of ICSU/COSPAR. He wrote 37 books and published 570 journal papers. He received 30 of awards including ICSU/COSPAR Vikram Sarabhai Medal in 2016, and Science award of Ministry of Education of Japan in 2015. He is now Editor-in-Chief of IJACSA and IJISA. <http://teagis.ip.is.saga-u.ac.jp/index.html>

Tuning of Customer Relationship Management (CRM) via Customer Experience Management (CEM) using Sentiment Analysis on Aspects Level

Hamed AL-Rubaiee

Department of Computer Science and Technology,
University of Bedfordshire
Bedfordshire, United Kingdom

Renxi Qiu

Department of Computer Science and Technology
University of Bedfordshire
Bedfordshire, United Kingdom

Khalid Alomar

Department of Information Systems
King Abdulaziz University
Jeddah, Kingdom of Saudi Arabia

Dayou Li

Department of Computer Science and Technology
University of Bedfordshire
Bedfordshire, United Kingdom

Abstract—This study proposes a framework that combines a supervised machine learning and a semantic orientation approach to tune Customer Relationship Management (CRM) via Customer Experience Management (CEM). The framework extracts data from social media first and then integrates CRM and CEM by tuning and optimising CRM to reflect the needs and expectations of users on social media. In other words, in order to reduce the gap between the users' predicted opinions in CRM and their opinions on social media, the existing data from CEM will be applied to determine the similar behavioural patterns of customers towards similar outcomes within CRM. CRM data and extracted data from social media will be consolidated by the unsupervised data mining method (association). The framework will lead to a quantitative approach to uncover relationships between the extracted data from social media and the CRM data. The results show that changing some aspects of the e-learning criteria that were required by students in their social media posts can help to enhance the classification accuracy in the learning management system (LMS) data and to understand more students' studying statuses. Furthermore, the results show matching between students' opinions in CRM and CEM, especially in the negative and neutral classes.

Keywords—Opinion mining; customer relationship management; customer experience management; sentiment analysis; Twitter

I. INTRODUCTION

Today, social media has become a method for maintaining strong relationships between users and companies.¹ With social media, companies are no longer in control of the relationship. Instead, customers are now driving the conversation [1]. To meet users' expectations and understand their opinions via social media platforms, companies are keen to adjust their CRM based on the feedback from social media platforms. In other words, CRM can be fine-tuned based on the differences between its own prediction of users' opinions and the actual

feedback from CEM. The differences between the opinions predicted by CRM and the actual opinions collected from social media are important to improve CRM further. The inconsistency in the opinions comes from two sources of error in CRM, the criteria or weights employed in CRM are not accurate, and the structure of CRM is not optimised.

In the former case, CRM and CEM cover similar aspects in a given domain; their difference is solely due to the inaccuracy of the criteria or weights employed in the CRM system. In the latter case, CRM could have missed a few key aspects that are important for CEM. The difference in the latter case is inherent in the structure of the CRM system.

Facing the problem, tuning CRM via CEM needs to be organised in two levels:

- Tuning on the aspect level: In this case, CRM and CEM have been constructed from similar aspects; sentiment analysis could be applied to adjust the criteria or weights in the CRM via the difference between CRM and CEM. This part of the problem will be investigated in this paper. A framework for adjusting CRM via CEM will be proposed and validated by a real-word example from the educational sector.
- Tuning on the sentence level: In this case, CRM and CEM have been constructed separately with different aspects. Opinions can only be collected on the sentence level for certain subjects in this domain.

This paper organised as follows: In Section 2, a related work for integrating CRM and CEM is introduced. In Section 3, a framework for integrating CRM and CEM is introduced. Section 4 contains the CRM tuning based on CEM at the aspect level—a case study from King Abdul-Aziz University. Section 5 presents sentiment analysis for CEM along with experimental results and evaluations. In Section 6, classification using CRM along with results and evaluations is

¹ <https://rapidminer.com>

shown. In Section 7, tuning of CRM via CEM is presented along with experimental results and evaluations. Finally, a summary of this paper is presented in Section 8.

II. RELATED WORK

In the first instance, a sentiment analysis system can be developed to determine consumer attitudes on products/services from review data. For example, when Dehkharghani and Yilmaz [2] conducted a review using a logistical classifier, they found that the average accuracy was 66.6% [3]. Bross and Ehrig [4] found that an aspect-based review to detect individual opinions and expressions about specific aspects of a product had a high accuracy [3]. Indhuja and Reghu [5] used an approach with novel fuzzy functions and achieved 85.58% accuracy. Wang [6] found that his model using a combined sentiment LDA and topic LDA was more effective than just topic sentiment analysis. Zhang and Varadarajan's [7] models, which incorporated features to predict utility scores of product reviews, achieved high performance, indicating that this is an effective approach [3].

The second example looks at the influence of Twitter; in terms of any relationship between sentiment analysis and the stock market, Bollen, Mao, and Zeng [8] reported that Twitter moods were used to predict the Dow Jones stock market index [9]. In their approach, public moods were measured using two tools: Google Profile of Mood (GPOM) and Opinion Finder [9] and a predictive model based on Self-Organised Fuzzy Neural Networks (SOFNN). The study found that the accuracy of the standard stock prediction model was significantly improved when mood dimensions and the value of the Dow Jones Industrial Average (DJIA) were included [9]. In addition, a study by Martin, Bruno, and Murisasco [10] showed a correlation between public opinion expressed on Twitter and the French stock market, using a neural network to find association patterns. They found that adding the sentiment feature on tweets two days before stock market closing values could improve accuracy [9].

Simsek and Ozdemir [11] also found a relationship between Turkish Twitter posts and the stock market index. By using the most common Turkish words representing happiness and unhappiness collected from tweets over a period of a month and a half, they worked out the frequencies of these two classes of words [9]. They found that the terms happy and trouble were commonly used in the emotional word database, and when the Twitter post contained stock market-related words, the average emotional value of the tweets changed from happiness to unhappiness by approximately 45% [9]. In a further study, Khatri, Singhal, and Johri [12] tried to train the neural network with words such as happy, hope, sad, and disappointing as input to predict the Bombay Exchange index. Their study indicated that an artificial neural network provided optimum results when set up with one hidden layer containing nine neurons [9]. Gao et al. [13] tried using a sentiment classification approach for Chinese stock news and found that pre-processing and a relevant sentiment dictionary affected the classification. Zhang and Skiena [14] showed that news sentiment can have an influence on market trading algorithms; they also found that news sentiment had a much sooner impact on stock markets than sentiment expressed on Twitter, which

could take up to three days [9]. All of these studies, however, point to a strong correlation between sentiments expressed on Twitter and global stock market indices.

In the third example, student feedback can highlight any issues students may have with the services provided by their colleges or universities. An example of this is when students cannot understand a lecturer or do not avail themselves of specific online services. Students have a habit of regularly using social media to express their opinions and describe activities they are involved in [15]. Therefore, universities utilise social media as a way of improving their teaching processes and generally to find out more about student experiences [15]. This is especially useful when finding out about online distance-education students, who do not give feedback face to face [15]. For instance, Tian et al. [16] developed an e-learner questionnaire and compared emotion words to measure the intensity of sentiment in each category. This approach enabled a positive result in dealing with challenges faced in the analysis of texts, such as those in Chinese, which are characterised by the richness of emotions. Wang, Zuo, and Diao [17] worked with the essential function of sentiment feedback in education over the Internet [15]. Wang [18] set up a Student Feedback Mining System (SFMS) to carry out an in-depth analysis of qualitative student feedback, which allowed insight into teaching practices, thus significantly improving student learning [15]. Donovan, Mader, and Shinsky [19] found that online student comments were much more detailed and informative than traditional paper-based feedback but are more time-consuming to analyse [15]. However, sentiment classification is important, as it gathers attitudes and opinions of users by mining and analysing personal information [18].

In general, a social CRM conversation between customers and enterprise agents over social media is called social CRM. Social CRM can influence the customer community to solve a customer's problem and turn a bad opinion into a good one [20]. Furthermore, social media has been used extensively by enterprises in the recent past to get insights about what users think about their products or services; this is typically achieved in a "listening" mode, i.e., a large amount of data from multiple social media sites is analysed in offline mode to extract aggregate-level business insights [21], [22]. Usually, the relational model database is the backbone for most companies or organizations; it stores the data in a structured order with rows and columns. However, a huge percentage of data in organizations or companies are located in unstructured data such as text data. This shows the need to analyse the unstructured data for companies' benefit [23].

Many researchers have shared their experiences with this subject. Ajmera et al. [24] built a social CRM that enabled firms to engage with customers by presenting analytical methods to identify actionable posts and analysing them. They presented novel features such as user intent and severity of issues in a customer complaint to determine a post's priority. In addition, Yaakub and Zhang [25] proposed a multi-dimensional model for opinion mining to integrate customers' characteristics and their related opinions about products. They used POS tagging to pre-process their data, trying to capture three parts from each document: nouns, which describe the

name of the product, and adverbs and adjectives, which describe the sentiment toward that product. However, it was hard to evaluate their model's conclusions about product attributes based on customers' opinions. In other words, it is hard to cover all details in CRM using only sentiment analysis. The baseline models for the above study were developed by trawling the reviews before putting them into a review database [26].

III. FRAMEWORK OF TUNING CRM VIA CEM

This study proposes a framework that combines a supervised machine learning and a semantic orientation approach to tune CRM via CEM. The framework extracts data from social media first and then integrates CRM and CEM by tuning and optimising CRM to reflect the needs and expectations of users on social media. In other words, in order to reduce the gap between the users' predicted opinions in CRM and their opinions on social media, the existing data from CEM will be applied to determine the similar behavioural patterns of customers towards similar outcomes within CRM. CRM data and extracted data from social media will be consolidated by the unsupervised data mining method (association). The framework will lead to a quantitative approach to uncover relationships between the extracted data from social media and the CRM data.

Fig. 1 illustrates the proposed framework for integrating CRM with CEM. The framework consists of three processes: (1) sentiment analysis for CEM, (2) classification using CRM, and (3) data tuning of CRM via CEM. In terms of data modelling, the main components of the three processes are very similar. The difference between process one and process two is in the input: one takes unstructured social media input, while the other takes a structured customer database. In process three, CRM data can be labelled automatically by CEM or vice versa. In this work, we will focus on the fine-tuning of CRM through CEM.

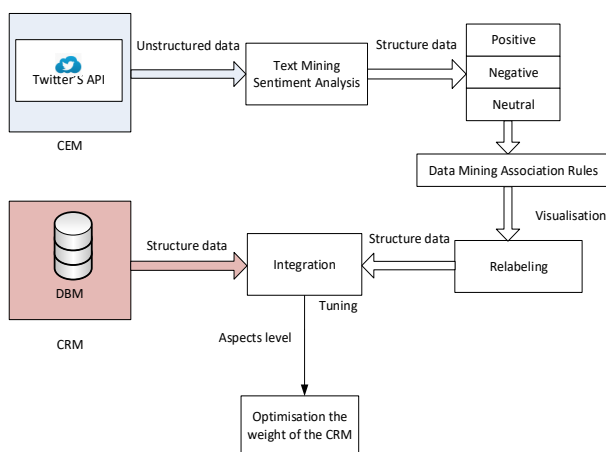


Fig. 1. The process of the proposed framework.

A. Sentiment Analysis for CEM

This is a process that extracts key features for CEM using social media. The extracted features will be represented by a semantic schema. The schema can be applied directly to the social media input. The process:

- Collects tweets based on a set of keywords that describe the case study using Twitter's streaming API.
- Performs text-processing techniques based on the proposed ontology model to reduce the amount of noise.
- Extracts key features with supervised learning. By exploring relevant data mining approaches such as sentiment analysis and NLP on building models, the component classifies tweets according to their sentiment polarity into one of the classes of positive, negative, and neutral.

B. Classification using CRM

This is a process that extracts key features for CRM using an existing customer database. The extracted features are also represented by a semantic schema, which can be applied directly to the consumer database.

- Collects structured data for the CRM using the database's API.
- Performs pre-processing techniques based on the existing CRM model to reduce the amount of noise.
- Extracts key features with supervised learning, which is again similar to the CEM process.

C. Tuning of CRM via CEM

A model has been developed that cross-validates features extracted from both CEM and CRM. In this case, if CRM and CEM are not consistent, CRM's semantic schema can be updated by CEM's output directly. In this process:

- Statistical algorithms will be applied to discover patterns and correlations in features extracted from both CEM and CRM.
- The confidence of the discovery will be examined automatically by comparing validation between the outputs. False positives, negatives, and neutrals will be identified during the process.
- CRM's semantic schema will be revised iteratively. If CRM and CEM are constructed from similar aspects, the tuning will be focused on the criteria or weights in CRM. Otherwise, CRM and CEM will need to be interpreted on the sentence level, so the structure of the CRM will be optimised based on its difference from CEM at this level.

IV. CASE STUDY FOR TUNING CRM VIA CEM AT THE ASPECT LEVEL: OPTIMISATION OF CRM OF KING ABDUL-AZIZ UNIVERSITY

The Deanship of e-Learning and Distance Education at King Abdul-Aziz University's Twitter account was chosen as a platform for opinion mining. The investigation was conducted on students' sentiments on the aspect level, with students being asked specific questions about e-learning criteria. The aim of this case study was to validate the proposed framework by adjusting the criteria or weights employed in the targeted CRM system.

A. Preparing the Sentiment Analysis for CEM using the Criteria from CRM

In order to prepare the sentiment and to clean the text to hand it off to the classifier, links, URLs, and hashtags were removed from the sentiments. In the following case study, the experimental collection process was based on certain hashtags that represent important topics, using Twitter hashtags as the domain.

B. CEM based on Sentiment on the Aspect Level

In this experiment, a hashtag was created to ask users for opinions on and reactions toward distance education criteria in order to identify students' positive, negative, and neutral opinions. The aspect level of sentiment analysis was followed for the students' opinions on the Twitter platform. The collection process was based on certain hashtags that represent the criteria, such as "#تطوير آليات التعليم عن بُعد" "# Enhance Distance education criteria". The experiment aimed to illustrate the relationship between the sentiments conveyed in Arabic tweets and the students' learning experiences at universities.

C. Experiment on Sentiment Analysis for CEM

During this stage, a hashtag was created to ask users for opinions on and reactions toward distance education criteria in order to identify students' positive, negative, and neutral opinions.

D. Experimental Parameter Settings

The classifier settings followed the default values for most of the parameters to get more accurate results. Ten-fold cross-validation was used several times to find the best values of these parameters. In addition, both SVM and NB had the same operator parameter settings used.

E. Data Collection and Description

To collect opinions that were comprehensive for the time period on the targeted objectives, tweets from different students on the Deanship of e-Learning and Distance Education page were obtained using Twitter's official developer's API. The data were collected over four months between 14 April and 13 August 2017. The data distribution depended on the methods Twitter's API utilised and the number of tweets posted on the distance education Twitter account page. Downloaded tweets were marked manually by employees of the distance education deanship as positive, negative, or neutral. In addition, the labelled tweets were stored in a database for experiments.

Distance education students' CRM data were obtained from the Learning Management System (LMS) Blackboard at King Abdul-Aziz University, which consists of the most important fields that can describe the students' activities. In this experiment, the collection process was based on students' data from one course chosen randomly from the student schedule of the second semester in 2017. Distance education students were labelled as "neutral" who had total marks in a course between 60 and 75 and attended 9 to 10 out of 12 lectures. For students who had total marks in a course less than 60 and attended fewer than 9 out of 12 lectures, the label was "negative". If a student had total marks above 75 and attended more than 10 out of 12 lectures, then he or she was considered an active

student and labelled as "positive". Table I shows the number of tweets that were involved in the classification experiment. Details of 143 students' were collected from the Blackboard CRM. All of them were anonymised. After that the data were labelled depending on distance education criteria, and then the data were stored. Rapidminer was used to load the dataset in normalized form with no duplicate records and no null attribute values.

After the data collection stage, the data model layout was created. This included a tweets description table (Tweet ID, Tweet Original, Tweet Filtered, Tweet Time, Tweet User, Tweet Label) and selected attributes from Blackboard databases. Again, this information was anonymised to protect privacy. Each student record contained the individual student's positive, negative, and neutral opinions, as well as the original tweet and the tweet labelled as an opinion from the CEM data for each individual student. Each student's record thus had two labels. The CRM label depended on the CRM criteria and the CEM label, which came from the sentiment analysis model.

Ten thousand six hundred and three students viewed the hashtag and agreed to start the experiment. Only 567 students registered their Twitter accounts and allowed the university to follow their tweets. One of the obstacles to letting students continue the survey was that the authentication from Twitter is in English. However, 242 of the 567 students have records in Blackboard, since only distance education students must use Blackboard regularly, whereas regular students and external students in some colleges still do not use Blackboard. Of these 242 students, 143 completed the survey.

The similar proposed model for sentiment analysis of Saudi Arabic (Standard and Arabian Gulf dialects) tweets was applied to extract feedback features from King Abdul-Aziz University data. The main idea was to examine the aspect level in sentiment analysis. In addition, the neutral class is important in the Arabic sentiment analyses, so our experiment was carried out with the neutral class as well as the negative and positive classes. The following tables show the type of tuning to provide the best accuracy in sentiment analysis for Arabic text in relation to King Abdul-Aziz University. For this experiment, 143 students' tweets responding to four questions were utilised. Table I shows the total number of tweets for the 143 students utilized for this experiment.

TABLE. I. NUMBER OF TWEETS UTILIZED FOR DATA MODEL COMPONENT ONE

Question	Positive	Negative	Neutral	Total
Q1	73	11	59	143
Q2	54	80	9	143
Q3	87	37	19	143
Q4	39	42	62	143
Total	338	369	261	572

F. Data Pre-Processing

A selection of 28 words and idioms in Arabic from the emotion corpus such as growth, good, excellent, problem, and inappropriate (see Table II) was then formed into the following three classes: positive, negative, and neutral. The most common Arabic words in the standard Saudi dialect represent positive and negative classes.

TABLE. II. COMMON ARABIC WORDS AND STANDARD SAUDI DIALECT REPRESENTING THE CLASSES

Arabic Positive Sentiment	English Translation	Arabic Negative Sentiment	English Translation
فتح تخصصات جديده	open new courses	مشكلة	problem
مناسب	suitable	تطوير بلاك يورد	develop Blackboard
جيد	good	لا يوجد تجاوب	no response
متكامل	complementary	غير مناسب	unsuitable
ممتاز	excellent	غير ملائم	inappropriate
سريع جدا	very fast	غير كافي	inadequate,
نظام جيدة	good platform	حرمان	absence exceeds limits
موافق	agree	زيادة النسبة	Increase percentage

Two distance education employees who have experience in e-learning labelled the data manually. Positive tweets were given the label "1", while negative tweets were given the label "-1". Neutral tweets were given the label "0", and irrelevant tweets were deleted from the database. A survey with four questions for e-learning criteria 2 was created and made available on Twitter via hashtags (see Table III). Three questions asked a specific question on the criteria, while the last one was an open question asking students about their opinions in general about e-learning. Table 4 shows the four questions asked of participants.

After data labelling was completed, the data were stored in our system in normalized form, as mentioned in the framework, with no hashtags, no duplicate tweets, no retweets, no URLs, and no special characters. After loading the dataset, the data were pre-processed with Rapidminer. The first step was to replace some Arabic words taking different shapes and icons.

TABLE. III. E-LEARNING CRITERIA

Number	E-learning criteria
1	Classwork: includes interaction with the instructor through the given activities available in the learning management system (30 marks of the grand total). Final examination: 70 marks of the grand total.
2	If student absence exceeded three lectures (equivalent to 25% of the synchronous online lectures) offered throughout the semester, then the student is prevented from taking the final exam of the distance learning course.
3	The quarterly work is divided as follows: assignments 4 (3 marks per assignment), activities 2 (6 marks per activity), forums (discussion board) 3 (2 marks per forum).
4	Any comments about the distance learning mechanism.

TABLE. IV. A SURVEY WITH FOUR QUESTIONS FOR E-LEARNING CRITERIA

Question No.	Question
Q1	What do you think about the following criteria: The student evaluates from 100 degrees: the quarterly work (30 degrees of the total) and the final test (70 degrees of the total).
Q2	What do you think about the following criteria: The student is denied entry to the final exam of the distance learning course if his absence exceeds 3 (25%) of the lectures.
Q3	What do you think about the following criteria: The quarterly work is divided as follows: four duties for each subject and the calculation of a single assignment three degrees, two periodic tests and each test six degrees, participation in the forums (discussion board) three posts, and each class two degrees (2).
Q4	Are there any comments you would like to share with us about the distance learning mechanism?

Then, the same pre-processing steps were performed: tokenization, removal of stop words, light stemming, filtering tokens by length, and application of the N-gram feature. Next, the 'Process Documents from Data' operator generated a word vector from the dataset after pre-processing and represented the text data as a matrix to show the frequency of occurrence of each term; then, relevant data mining approaches, i.e., NB and SVM, were explored for building models to classify tweets according to their sentiment polarity into positive, negative, and neutral. Finally, an evaluation was carried out using precision and recall methods.

G. Experiment Results and Evaluations

The results were divided into four groups for each question to show the sentiment analysis classification accuracy, precision, and recall for the NB and SVM classifiers with and without the N-gram feature.

- Sentiment analysis classification for Question 1

Table V shows Q1 classification accuracy, precision, and recall for the NB and SVM classifiers without the N-gram feature: Crosse-validation=10, sampling type=stratified sampling, prune=none. In addition, Table VI shows the class accuracy, precision, and recall for the NB and SVM classifiers with the N-gram feature, which is set to two: Crosse-validation=10, sampling type=stratified sampling, prune=none.

In conclusion, the experiment shows that NB performance was better when we used the N-gram feature with both schemas (TF-IDF and BTO). On the other hand, there was a slight performance increase when SVM used the same feature. However, the best accuracy was achieved by SVM with the TF-IDF schema when the N-gram feature was not involved.

²http://elearning.kau.edu.sa/Content.aspx?Site_ID=214&lng=EN&cid=24144

TABLE. V. ACCURACY, PRECISION, AND RECALL FOR ALL CLASSES USING SVM AND NB CLASSIFIERS WITHOUT N-GRAM FOR QUESTION 1

Classes	Classifier Name	Weighting schemes	Accuracy	Class recall	Class precision	Classification error
All Class	NB	BTO	39.24	43.95	48.18	60.76
		TF-IDF	55.19	54.59	54.54	44.81
	SVM	BTO	77.57	55.06	52.73	22.43
		TF-IDF	83.19	59.65	56.74	16.81

TABLE. VI. ACCURACY, PRECISION, AND RECALL FOR ALL CLASSES USING SVM AND NB CLASSIFIERS WITH N-GRAM FOR QUESTION 1

Classes	Classifier Name	Weighting schemes	Accuracy	Class recall	Class precision	Classification error
All Class	NB	BTO	55.24	53.36	54.42	44.76
		TF-IDF	67.86	61.65	62.66	32.14
	SVM	BTO	78.95	61.73	60.42	21.05
		TF-IDF	81.05	63.59	61.97	18.95

TABLE. VII. ACCURACY, PRECISION, AND RECALL FOR ALL CLASSES USING FOR SVM AND NB CLASSIFIERS WITHOUT N-GRAM FOR QUESTION 2

Classes	Classifier Name	Weighting schemes	Accuracy	Class recall	Class precision	Classification error
All Class	NB	BTO	77.62	66.44	63.18	22.38
		TF-IDF	76.24	65.47	62.26	23.76
	SVM	BTO	78.29	56.11	52.81	21.71
		TF-IDF	72.76	49.67	50.48	27.24

TABLE. VIII. ACCURACY, PRECISION, AND RECALL FOR ALL CLASSES USING FOR SVM AND NB CLASSIFIERS WITH N-GRAM FOR QUESTION 2

Classes	Classifier Name	Weighting schemes	Accuracy	Class recall	Class precision	Classification error
All Class	NB	BTO	78.33	66.86	63.54	21.67
		TF-IDF	77.67	66.31	63.06	22.33
	SVM	BTO	74.90	54.53	51.81	25.10
		TF-IDF	70.00	46.94	49.80	30.00

TABLE. IX. ACCURACY, PRECISION AND RECALL FOR ALL CLASSES USING FOR SVM AND NB CLASSIFIERS WITHOUT N-GRAM FOR QUESTION 3

Classes	Classifier Name	Weighting schemes	Accuracy	class recall	Class precision	Classification error
All Class	NB	BTO	72.05	64.58	66.34	27.95
		TF-IDF	76.95	67.22	67.20	23.05
	SVM	BTO	80.33	61.85	62.91	19.67
		TF-IDF	82.52	65.19	68.74	17.48

TABLE. X. ACCURACY, PRECISION AND RECALL FOR ALL CLASSES USING FOR SVM AND NB CLASSIFIERS WITH N-GRAM FOR QUESTION 3

classes	Classifier Name	Weighting schemes	Accuracy	class recall	Class precision	Classification error
All Class	NB	BTO	80.48	70.46	67.66	19.52
		TF-IDF	81.86	71.76	71.15	18.14
	SVM	BTO	79.67	62.69	63.60	20.33
		TF-IDF	81.05	61.39	61.32	18.95

TABLE. XI. ACCURACY, PRECISION, AND RECALL FOR ALL CLASSES USING FOR SVM AND NB CLASSIFIERS WITHOUT N-GRAM FOR QUESTION 4

Classes	Classifier Name	Weighting schemes	Accuracy	Class recall	Class precision	Classification error
All Class	NB	BTO	60.00	61.51	59.15	40.00
		TF-IDF	69.10	71.34	70.53	30.90
	SVM	BTO	58.00	53.27	65.65	42.00
		TF-IDF	68.43	67.48	74.74	31.57

TABLE. XII. ACCURACY, PRECISION, AND RECALL FOR ALL CLASSES USING FOR SVM AND NB CLASSIFIERS WITH N-GRAM FOR QUESTION 4

Classes	Classifier Name	Weighting schemes	Accuracy	Class recall	Class precision	Classification error
All Class	NB	BTO	59.29	61.03	58.54	40.71
		TF-IDF	69.81	71.90	71.59	30.19
	SVM	BTO	49.81	42.33	48.96	50.19
		TF-IDF	62.86	62.89	72.66	37.14

- Sentiment Analysis Classification for Question 2

Table VII shows Q2 classification accuracy, precision, and recall for the NB and SVM classifiers without the N-gram feature: Crosse-validation=10, sampling type=stratified sampling, prune=none. In addition, Table VIII shows the class accuracy, precision, and recall for the NB and SVM classifiers with the N-gram feature, which is set to two: Crosse-validation=10, sampling type=stratified sampling, prune=none.

In conclusion, the experiment shows that NB performance was better when we used the N-gram feature with both schemas (TF-IDF and BTO). On the other hand, there was a performance decrease when SVM used the same feature. However, the best accuracy was achieved by SVM with the BTO schema when the N-gram feature was not involved.

- Sentiment Analysis Classification for Question 3

Table 9 shows Q3 classification accuracy, precision, and recall for the NB and SVM classifiers without the N-gram feature: Crosse-validation=10, sampling type=stratified sampling, prune=none. In addition,

X shows the class accuracy, precision, and recall for the NB and SVM classifiers with the N-gram feature, which is set to two: Crosse-validation=10, sampling type=stratified sampling, prune=none.

In conclusion, the experiment shows that NB performance was better when we used the N-gram feature with both schemas (TF-IDF and BTO). On the other hand, there was a performance drop when SVM used the N-gram feature with both schemas (TF-IDF and BTO).

- Sentiment Analysis Classification for Question 4

Table XI shows Q6 classification accuracy, precision, and recall for the NB and SVM classifiers without the N-gram feature: Crosse-validation=10, sampling type=stratified sampling, prune=none. In addition, Table XII shows the class accuracy, precision, and recall for the NB and SVM classifiers with the N-Gram feature, which is set to two: Crosse-validation=10, sampling type=stratified sampling, prune=none.

In conclusion, the experiment shows that NB performance was better when we used the N-gram feature with both schemas (TF-IDF and BTO). On the other hand, there was a drop when SVM used the N-gram feature with both schemas (TF-IDF and BTO).

V. CLASSIFICATION USING CRM

Component two presented the design and implementation of students' CRM data (or classification through different algorithms, such as SVM and NB. In this part of the experiment, the collection process was based on students' Blackboard attributes assigned according to their points of view on different classes in Blackboard data. It was labelled according to the university's e-learning criteria.

A. Details of using the Classifiers with Rapidminer for CRM Classification

The classifier settings followed the default values for most of the parameters to get more accurate results. Ten-fold cross-

validation was used several times to find the best values of these parameters.

B. Support Vector Machine (SVM) Operator Parameter Settings

The SVM operator generates the SVM classification model. This model can be used for classification and provides good results for many learning tasks. In addition, it supports various kernel types, including dot, radial, polynomial, and neural³.

- SVM type: C-SVM, which is for classification tasks.
- Linear: A linear classifier works based on the value of a linear combination.
- C is the penalty parameter of the error term and was set to its default value, which is zero.
- Cache size is an expert parameter. It specifies the cache size in megabytes and was set to default value, which is 80.
- Epsilon: This parameter specifies the tolerance of the termination criterion and was set to the default value, which is 0.001.

C. Naïve Bayes (NB) Operator Parameter Settings

The NB operator generates a NB classification model. It is a probabilistic classifier based on applying Bayes' theorem with powerful independence assumptions. The NB classifier assumes that the presence or absence of a particular feature of a class is unrelated to any other feature.

- Laplace correction:

This parameter indicates whether Laplace correction should be used to prevent high influence of zero probabilities. Assume that our training set is so large that adding one to each count will make a small difference in the estimated likelihoods.

VI. EXPERIMENT OF THE CLASSIFICATION USING CRM

For this experiment 567 student registered their twitter account and allow the university to follow their tweets. However, only 234 students are distance education students and have full record in Blackboard, where other students do not use Blackboard. Table XIII shows the total number of tweets for the 242 students utilised for this experiment.

TABLE. XIII. NUMBER OF TWEETS UTILIZED FOR DATA MODEL COMPONENT TWO

Question	Positive	Negative	Neutral	Total
Q1	84	65	93	242
Q2	90	122	30	242
Q3	84	89	69	242
Q4	80	93	69	242
Total	338	369	261	968

A. Data Description for Question 1, 2, 3, and 4

Table XIV shows the criteria utilised to label record as positive, negative, or neutral for question 1.

³ <https://rapidminer.com>

TABLE. XIV. QUESTION 1 CRITERIA

	Labelled	Criteria
Student Mark at Blackboard (total mark is 30)	Positive	>=20
	Neutral	15-20
	Negative	<15
Student Final Exam Mark in ODUS from (total mark is 70)	Negative	<35
	Neutral	35-45
	Positive	>45
Student Final Mark in ODUS from (total mark is 100)	Negative	<50
	Neutral	50-75
	Positive	>75
TOT_MARK = BB_F_MARK_NO + ODS_MARK		

Table XV shows the criteria utilised to label record as positive, negative, or neutral for question 1.

TABLE. XV. QUESTION 2 CRITERIA

	Labelled	Criteria
Number of recorded attendance	Negative	>5 days
	Neutral	4-5 days
	Positive	<4 days
Number of online attendance	Negative	<9 lectures
	Neutral	9-10 lectures
	Positive	>10 lectures
ATTEND_MARK = REC_MARK_COUNT + ONLINE_MARK		

Table XVI shows the criteria utilised to label record as positive, negative, or neutral for question 1.

TABLE. XVI. QUESTION 3 CRITERIA

	Labelled	Criteria
Number of times student participate in the discussion, quizzes and assignments	Negative	>5
	Neutral	5-6
	Positive	>6
FORUM_TEST_MARK_TOTAL = FORUM_TEST_MARK_FORUM + FORUM_TEST_MARK_TEST		

Table XVII shows the criteria utilised to label record as positive, negative, or neutral for question 1.

TABLE. XVII. QUESTION 4 CRITERIA

	Labelled	Criteria
Total Mark	Negative	<3
	Neutral	<5
	Positive	>6
TOTAL_EVA_CHAR = (TOT_MARK + ATTEND_MARK + FORUM_TEST_MARK_RESULT) / 3		

VII. EXPERIMENTAL RESULTS AND EVALUATIONS

The results were divided into groups for each question to show the Blackboard students' data classification accuracy, precision, and recall for the NB and SVM classifiers.

A. CRM Classification for Question 1, 2, 3, and 4

The following table shows the type of tuning to provide the best accuracy in King Abdul-Aziz University's CRM classification. Table XVIII shows the classification accuracy, precision, and recall for the NB and SVM classifiers: Crosse-validation=10, sampling type=stratified sampling. The best accuracy was achieved by NB due to the advantages of NB, such as its simplicity, ease of implementation, and combination of efficiency with acceptable accuracy.

TABLE. XVIII. ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS

Classes	Classifier name	Accuracy	Class recall	Class precision	Classification error
Positive, negative, and neutral	NB	93.43	93.95	94.12	6.57%
	SVM	87.25	86.55	89.74	12.75%

The following table shows the type of tuning to provide the best accuracy in King Abdul-Aziz University's CRM classification.

Classes	Classifier name	Accuracy	Class recall	Class precision	Classification error
Positive, negative, and neutral	NB	95.05	96.69	96.48	4.95%
	SVM	81.80	62.82	55.83	18.20%

XIX shows the classification accuracy, precision, and recall for the NB and SVM classifiers: Crosse-validation=10, sampling type=stratified sampling. The best accuracy was achieved by NB due to the advantages of NB, such as its simplicity, ease of implementation, and combination of efficiency with acceptable accuracy.

TABLE. XIX. ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS

Classes	Classifier name	Accuracy	Class recall	Class precision	Classification error
Positive, negative, and neutral	NB	95.05	96.69	96.48	4.95%
	SVM	81.80	62.82	55.83	18.20%

The following table shows the type of tuning to provide the best accuracy in King Abdul-Aziz University's CRM classification.

Classes	Classifier name	Accuracy	Class recall	Class precision	Classification error
Positive, negative, and neutral	NB	98.33	72.44	72.31	1.67%
	SVM	96.72	65.44	64.52	3.28%

XX shows the classification accuracy, precision, and recall for the NB and SVM classifiers: Crosse-validation=10, sampling type=stratified sampling. The best accuracy was achieved by NB due to the advantages of NB, such as its simplicity, ease of implementation, and combination of efficiency with acceptable accuracy.

TABLE. XX. ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS

Classes	Classifier name	Accuracy	Class recall	Class precision	Classification error
Positive, negative, and neutral	NB	98.33	72.44	72.31	1.67%
	SVM	96.72	65.44	64.52	3.28%

The following table shows the type of tuning to provide the best accuracy in King Abdul-Aziz University's CRM classification. Table XXI shows the classification accuracy, precision, and recall for the NB and SVM classifiers: Crosse-

validation=10, sampling type=stratified sampling. The best accuracy was achieved by NB due to the advantages of NB, such as its simplicity, ease of implementation, and combination of efficiency with acceptable accuracy.

TABLE. XXI. ACCURACY, PRECISION AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS.

Classes	Classifier name	Accuracy	Class recall	Class precision	Classification error
Positive, negative, and neutral	NB	85.13	84.61	85.88	14.87%
	SVM	83.05	81.47	83.12	16.95%

VIII. TUNING OF CRM VIA CEM

This process aims to find a way to support CRM (or Blackboard in this case study). For instance, what exactly do students want or think about experiences? This is especially relevant for students who mostly depend on the Web, such as online distance education students, also proving that the value of social media information can bring better understanding of a student’s study situation. Therefore, similar tweets are grouped such that tweets within the same group bear similarity to each other, while tweets in different groups are dissimilar from each other. This will help to understand students’ behaviours and find out the most common problems. Moreover, this will give the university the ability to learn about and validate students’ data with more support from social media to develop new e-learning criteria to match the inputs from social media.

This case study investigates an alternative solution for supporting CRM by social media inputs on the aspect sentiment level and tuning CRM weights for some aspects of the e-learning criteria that students need, according to their posts on social media. Optimisation of the CRM weights was applied to update some aspects’ values in CRM. The results show closely CRM’s student labels match CEM’s, especially in the negative and neutral classes. Furthermore, they show that optimising CRM’s weights can enhance classification accuracy in the Blackboard data and help to understand more students’ studying statuses.

A. Experiment in Tuning of CRM via CEM

To be included in this experiment, students should have records in CRM and have completed the survey. Out of the 143 with records on Blackboard, only 79 completed the survey. Table XXII shows the difference between CRM and CEM in this case study before the tuning.

TABLE. XXII. NUMBER OF TWEETS UTILIZED FOR DATA MODEL COMPONENT THREE

Question	CEM/CRM	Positive	Negative	Neutral
Q1	CEM	46	5	28
	CRM	31	23	25
Q2	CEM	29	46	4
	CRM	37	24	18
Q3	CEM	28	26	25
	CRM	36	13	30
Q4	CEM	18	26	35

	CRM	31	20	28
Total CEM		121	103	92
Total CRM		135	80	101

B. Experiment Results and Evaluations

The results were divided into four groups for each question to show the comparison between the CRM and CEM labelling classes.

a) Integration results for Question 1

Fig. 2 shows a comparison between CRM labelling and CEM labelling. Thirty tweets were labelled similarly by CRM and CEM. On the other hand, 48 tweets were labelled dissimilarly by CEM and CRM. In other words, 62% of collected tweets were labelled differently by CRM and CEM, and only 38% were similarly labelled by CEM and CRM.

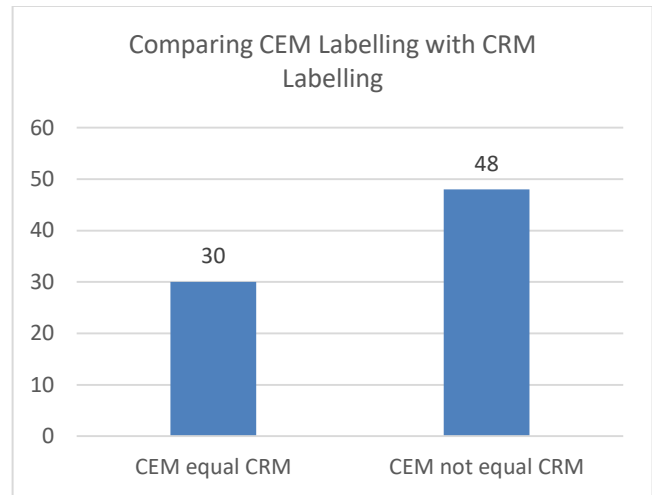


Fig. 2. Comparing CEM labelling with CRM labelling for Question 1.

From Fig. 3, there are 46 positive tweets in CEM, compared to 31 tweets in CRM. In contrast, there are only five negative tweets in CEM, compared to 23 negative tweets in CRM. For the neutral class, there are 28 tweets in CEM and 25 in CRM.

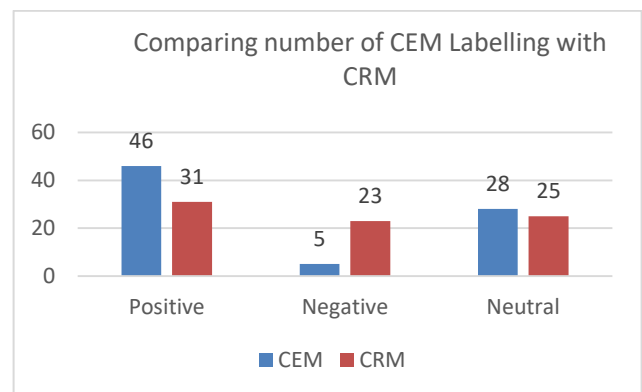


Fig. 3. Comparing CEM labelling class with CRM labelling class for Question 1.

b) Integration results for Question 2

Fig. 4 shows a comparison between CRM labelling and CEM labelling. Twenty-three tweets were labelled similarly between CRM and CEM. On the other hand, 56 tweets were labelled dissimilarly between CEM and CRM. In other words, 71% of collected tweets were labelled differently between CRM and CEM, and only 29% had the same labelling between CEM and CRM.

Fig. 5 shows the number of tweets with similar labelling between the CRM positive class and to the CEM (positive, negative, and neutral) labelling for Question 2. In CEM, 29 tweets were labelled as positive. By contrast, CRM labelled 11 of these tweets as negative, 10 as positive, and eight as neutral.

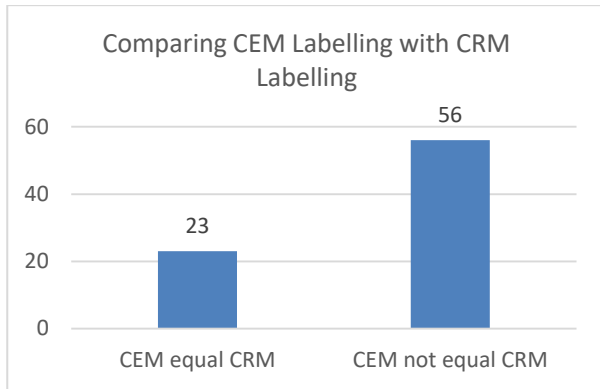


Fig. 4. Comparing CEM labelling with CRM labelling for Question 2.

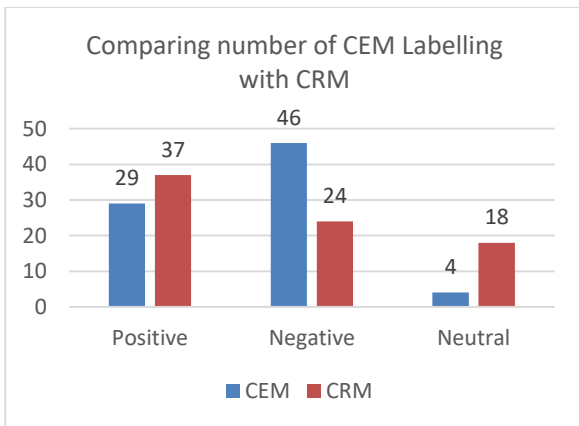


Fig. 5. Comparing CEM labelling with CRM labelling for Question 2.

c) Integration results for Question 3

Fig. 6 shows a comparison between CRM labelling with CEM labelling. Twenty tweets were labelled similarly between CRM and CEM. On the other hand, 59 tweets were labelled dissimilarly between CEM and CRM. In other words, 75% of collected tweets were labelled differently between CRM and CEM, and only 25% had similar labelling.

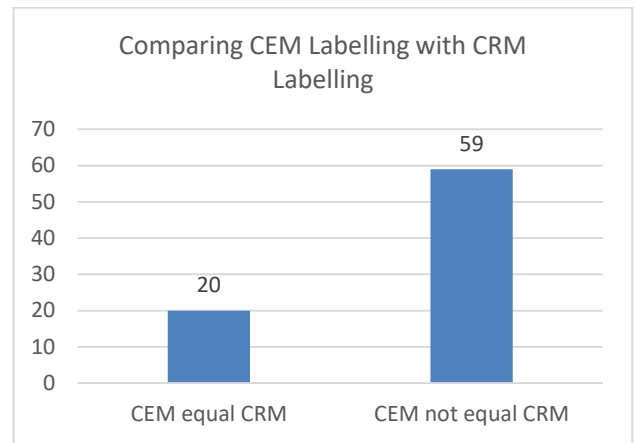


Fig. 6. Comparing CEM labelling with CRM labelling for Question 3

From Fig. 7, there were 28 positive tweets in CEM, compared to 36 tweets in CRM. In contrast, CEM had only 13 negative tweets, compared to 26 in CRM. The neutral class contained 25 tweets in CEM and 30 in CRM, a difference of only five tweets.

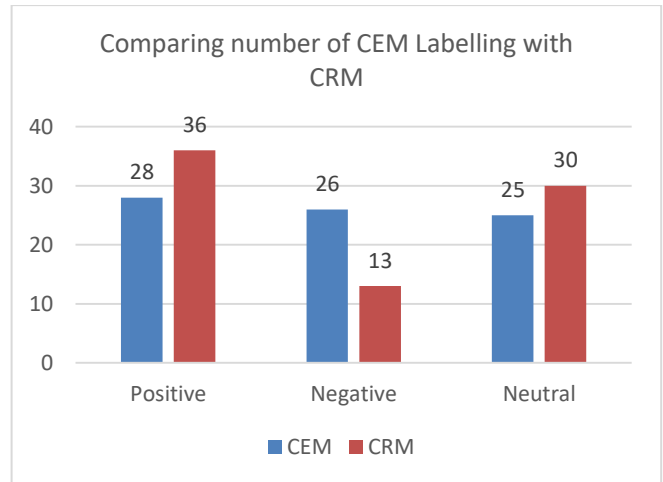


Fig. 7. Comparing CEM labelling with CRM labelling for Question 3.

d) Integration Results for Question 4

Fig. 8 shows a comparison between CRM labelling and CEM labelling. Twenty-four tweets were labelled similarly between CRM and CEM. On the other hand, 55 tweets were labelled dissimilarly between CEM and CRM. In other words, 70% of the collected tweets were labelled differently between CRM and CEM, and only 30% had the same labelling.

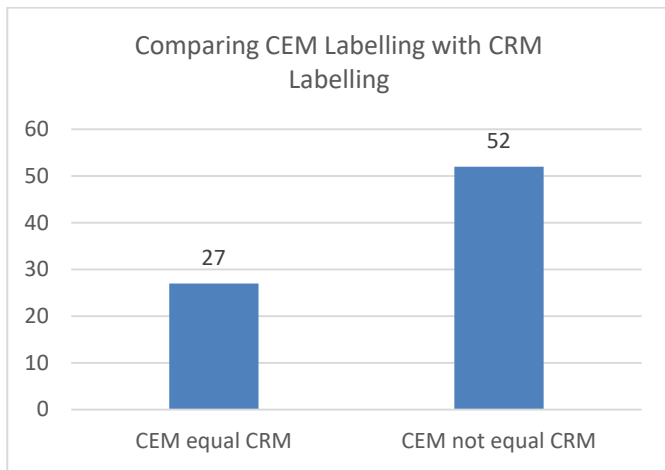


Fig. 8. Comparing CEM labelling with CRM labelling for Question 4.

From Fig. 9, there were 18 positive tweets in CEM, compared to 31 tweets in CRM. In contrast, there were 26 negative tweets in CEM, compared to 20 negative tweets in CRM. For the neutral class, there were 35 tweets in CEM, compared to 28 in CRM.

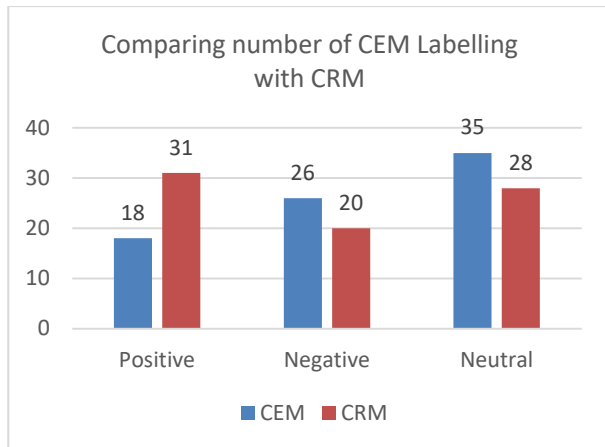


Fig. 9. Comparing CEM labelling with CRM labelling for Question 4.

IX. TUNING THE CRM WEIGHTS

The aim of this experiment was to tune CRM weights for some aspects of e-learning criteria that students need, according to their social media posts. To achieve this aim, first, CRM weights were updated by changing the values of some attributes that reflected the aspect-level opinions in CEM. After that, the criteria were also updated according to the input of the opinions in CEM. Then, the validation was carried out for the new criteria. Last, classification was carried out in order to compare the original criteria with the updated ones. The weight of CRM changed based on the input from social media. For this study, the values of the criteria were changed for aspects of Q1. Changes occurred for students with Blackboard marks between 30 and 40 and student final exam marks in ODUS from 70 to 60. After updating the weight of CRM, another classification was carried out in order to evaluate the accuracy, precision, and recall of all classes for the SVM and NB classifiers. After that, a comparison was carried to find out the accuracy, precision, and recall of all Classes for the SVM and

NB classifiers before and after CRM weight tuning. Table XXVI shows that the accuracy before the tuning is higher than after the tuning since the matching between the CRM criteria weights and students' opinions on social media are nearly the same for Q1. This indicates that there was no need to update the criteria weights for the aspects in Q1. Table XXIII shows the validation of the classification accuracy, precision, and recall for the NB and SVM classifiers for Question 1: Crosse-validation=10, sampling type=stratified sampling.

TABLE. XXIII. EVALUATION FOR ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS FOR QUESTION 1

Tuning	Classifier name	Accuracy	Class recall	Class precision	Classification error
Before	NB	93.43	93.95	94.12	6.57%
	SVM	87.25	86.55	89.74	12.75%
After	NB	92.58	92.38	92.59	7.42%
	SVM	85.18	83.37	84.70	14.82%

Values of the criteria were changed for aspects in Q2: Changes occurred in student attendance in CRM from seven days to 14 days. Table XXIV shows the comparison between the evaluation of the classification before and after the tuning. The results show that the accuracy after the tuning is higher than before for the aspects in Q2.

TABLE. XXIV. EVALUATION FOR ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS FOR QUESTION 2

Tuning	Classifier name	Accuracy	Class recall	Class precision	Classification error
Before	NB	95.05	96.69	96.48	4.95%
	SVM	81.80	62.82	55.83	18.20%
After	NB	95.85	96.09	95.76	4.15%
	SVM	84.70	65.96	57.98	15.30%

Values of the criteria were changed for aspects in Q3: Changes occurred in evaluating the total number of posts, including discussions and tests. Table XXV shows the comparison between the evaluation of the classification before and after the tuning. The results show that the accuracy after the tuning is almost the same.

TABLE. XXV. EVALUATION FOR ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS FOR QUESTION 3

Tuning	Classifier name	Accuracy	Class recall	Class precision	Classification error
Before	NB	98.33	72.44	72.31	1.67%
	SVM	96.72	65.44	64.52	3.28%
After	NB	98.35	72.44	72.37	1.65%
	SVM	95.92	64.96	63.92	4.08%

Q4 combines the aspects in Q1, Q2, and Q3. Table XXVI shows the comparison between the evaluation of the classification before and after the tuning. The results show that the accuracy after the tuning is higher than before. This indicates that changing some aspects of the criteria can help to enhance the classification accuracy in the CRM data.

TABLE. XXVI. EVALUATION FOR ACCURACY, PRECISION, AND RECALL WITH ALL CLASSES FOR SVM AND NB CLASSIFIERS FOR QUESTION 4

Tuning	Classifier name	Accuracy	Class recall	Class precision	Classification error
Before	NB	85.13	84.61	85.88	14.87%
	SVM	83.05	81.47	83.12	16.95%
After	NB	92.58	92.38	92.59	7.42%
	SVM	85.18	83.37	84.70	14.82%

To sum up King Abdul-Aziz University’s validation CRM experiments, the best accuracy was achieved by NB due to the advantages of NB such as its simplicity, ease of implementation, and combination of efficiency and acceptable accuracy.

Fig. 10 shows a comparison between the CEM labelling and the distance education criteria with suggestion one for Question 1. The number of positive tweets increases, as does the number of negative tweets, while the number of neutral tweets sharply decreases.

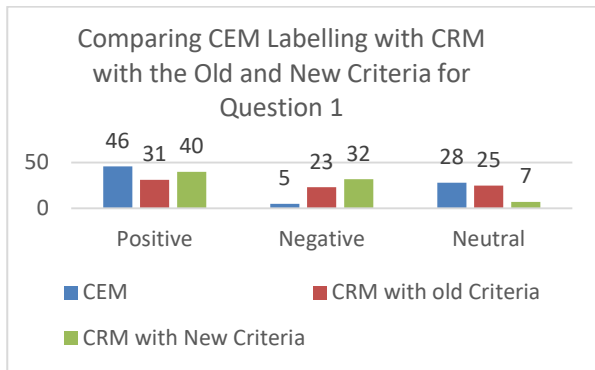


Fig. 10. Comparing CEM labelling with CRM with the old and new criteria for Question 1.

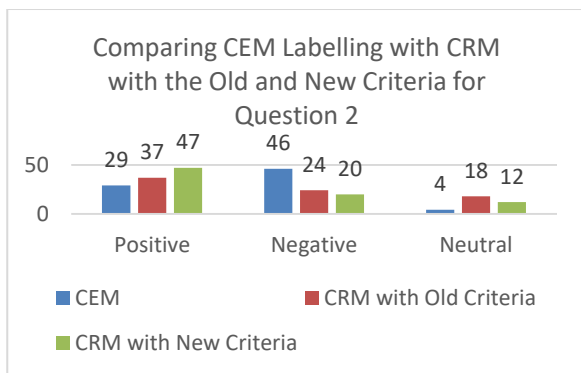


Fig. 11. Comparing CEM labelling with CRM with the old and new criteria for Question 2.

Fig. 11 shows a comparison between the CEM labelling and the distance education criteria with suggestion one for Question 1. The number of positive tweets increases, and the number of negative tweets decreases, as does the number of neutral tweets.

Fig. 12 shows a comparison between the CEM labelling and the distance education criteria with suggestion one for Question 1. The number of positive tweets approximately doubles, and the number of negative tweets approximately halves. In addition, the number of neutral tweets drops to zero.

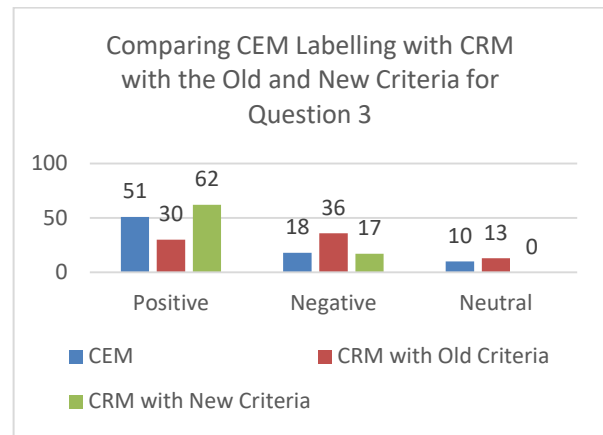


Fig. 12. Comparing CEM labelling with CRM with the old and new criteria for Question 3.

Fig. 13 shows a comparison between the CEM labelling and the distance education criteria with suggestion one for Question 1. The number of positive and neutral tweets increases, and the number of negative tweets decreases.

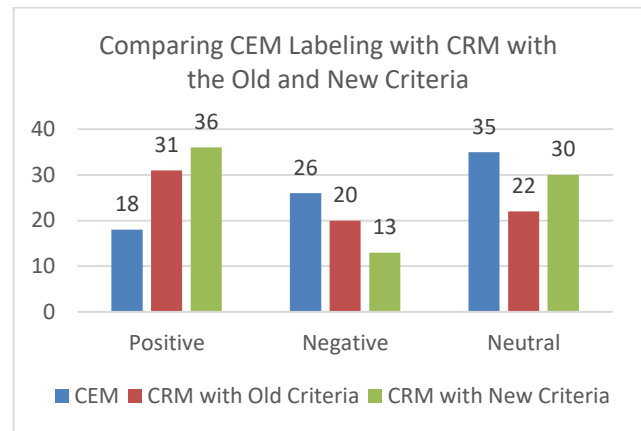


Fig. 13. Comparing CEM labelling with CRM with the old and new criteria for Question 4.

In conclusion, the result of this experiment shows that social media can support CRM with more details. The main aim of this study was achieved by showing the gap between the criteria and students’ needs. The validation result demonstrates the gap between the students’ perspectives and the criteria. This might help universities to adapt the distance education criteria in a way that helps the students to deal with the university. The validation results confirm that tuning the criteria will help students. However, the main reasons for the difference of the classification results with the given examples above is domain diversity, as well as the way that data were pre-processed, the size of the dataset, and of course the approaches that were used in each article. Moreover, In terms of the accuracy of the last survey, which studied recent work in Arabic sentiment analysis, SVM was applied successfully in several sentiment analysis tasks [27], [28].

X. CONCLUSION

This study carried out to investigate whether social media could help experts to understand users’ perspectives and could

support academics' knowledge about their students. Therefore, a framework was proposed for integrating CRM with CEM on the aspect level. The framework consists of three components: sentiment analysis for CEM, classification using CRM Blackboard data, and tuning of CRM via CEM, which integrates both results to study the level of matching between both resources' information, namely, social media and Blackboard data. In other words, there is good consistency between the CRM structure and students' opinions on the aspect level. This takes the study further by investigating King Abdul-Aziz University in e-learning and distance criteria, which brings the similarity between CEM and CRM opinions closer. Moreover, the final stage of this experiment shows an interesting result after changing the e-learning criteria according to the necessary requirements of input on social media requested in students' feedback comments through the Twitter platform. The results show that changing some aspects of the e-learning criteria that were required by students in their social media posts can help to enhance the classification accuracy in the Blackboard data and to understand more students' studying statuses. Furthermore, the results show matching between students' opinions in CRM and CEM, especially in the negative and neutral classes.

REFERENCES

- [1] Tripathi, G. and S. Naganna, Feature selection and classification approach for sentiment analysis. *Machine Learning and Applications: An International Journal*, 2015. 2(2): p. 1-16.
- [2] Dehkharghani, R. and C. Yilmaz. Automatically identifying a software product's quality attributes through sentiment analysis of tweets. in *Natural Language Analysis in Software Engineering (NaturaLiSE)*, 2013 1st International Workshop on. 2013. IEEE.
- [3] Al-Rubaiee, H., R. Qiu, and D. Li. Identifying Mubasher software products through sentiment analysis of Arabic tweets. in *Industrial Informatics and Computer Systems (IIICS)*, 2016 International Conference on. 2016. IEEE.
- [4] Bross, J. and H. Ehrig. Generating a context-aware sentiment lexicon for aspect-based product review mining. in *Web Intelligence and Intelligent Agent Technology (WI-IAT)*, 2010 IEEE/WIC/ACM International Conference on. 2010. IEEE.
- [5] Indhuja, K. and R.P. Reghu. Fuzzy logic based sentiment analysis of product review documents. in *Computational Systems and Communications (ICCSC)*, 2014 First International Conference on. 2014. IEEE.
- [6] Wang, W. Sentiment analysis of online product reviews with Semi-supervised topic sentiment mixture model. in *Fuzzy Systems and Knowledge Discovery (FSKD)*, 2010 Seventh International Conference on. 2010. IEEE.
- [7] Zhang, Z. and B. Varadarajan. Utility scoring of product reviews. in *Proceedings of the 15th ACM international conference on Information and knowledge management*. 2006. ACM.
- [8] Bollen, J., H. Mao, and X. Zeng. Twitter mood predicts the stock market. *Journal of Computational Science*, 2011. 2(1): p. 1-8.
- [9] Hamed, A.-R., R. Qiu, and D. Li. Analysis of the relationship between Saudi twitter posts and the Saudi stock market. in *Intelligent Computing and Information Systems (ICICIS)*, 2015 IEEE Seventh International Conference on. 2015. IEEE.
- [10] Martin, V., E. Bruno, and E. Murisasco, Predicting the French Stock Market Using Social Media Analysis. *International Journal of Virtual Communities and Social Networking (IJVCSN)*, 2015. 7(2): p. 70-84.
- [11] Şimşek, M.U. and S. Özdemir. Analysis of the relation between Turkish twitter messages and stock market index. in *Application of Information and Communication Technologies (AICT)*, 2012 6th International Conference on. 2012. IEEE.
- [12] Khatri, S.K., H. Singhal, and P. Johri. Sentiment analysis to predict Bombay stock exchange using artificial neural network. in *Reliability, Infocom Technologies and Optimization (ICRITO)(Trends and Future Directions)*, 2014 3rd International Conference on. 2014. IEEE.
- [13] Gao, Y., et al. Sentiment classification for stock news. in *Pervasive Computing and Applications (ICPCA)*, 2010 5th International Conference on. 2010. IEEE.
- [14] Zhang, W. and S. Skiena. Trading Strategies to Exploit Blog and News Sentiment. in *Icwsn*. 2010.
- [15] Hamed AL-Rubaiee, R.Q., Khalid Alomar and Dayou Li, Sentiment Analysis of Arabic Tweets in e-Learning. *Journal of Computer Science*. 12(11): p. 553-563.
- [16] Tian, F., et al. Can e-Learner's emotion be recognized from interactive Chinese texts? in *Computer Supported Cooperative Work in Design, 2009. CSCWD 2009*. 13th International Conference on. 2009. IEEE.
- [17] Wang, P., M. Zuo, and L. Diao. Design of the monitoring and analysis system of the Internet Education Public Sentiment. in *Future Computer and Communication (ICFCC)*, 2010 2nd International Conference on. 2010. IEEE.
- [18] Fu, G. and X. Wang. Chinese sentence-level sentiment classification based on fuzzy sets. in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*. 2010. Association for Computational Linguistics.
- [19] Donovan, J., C.E. Mader, and J. Shinsky. Constructive Student Feedback: Online vs. Traditional Course Evaluations. *Journal of Interactive Online Learning*, 2010. 9(3): p. 283-296.
- [20] Mohan, S., E. Choi, and D. Min. Conceptual modeling of enterprise application system using social networking and Web 2.0 "social CRM system". in *Convergence and Hybrid Information Technology, 2008. ICHIT'08*. International Conference on. 2008. IEEE.
- [21] Biere, M., *Business intelligence for the enterprise*. 2003: Prentice Hall Professional.
- [22] Spangler, S., et al., COBRA—Mining web for corporate brand and reputation analysis. *Web Intelligence and Agent Systems: An International Journal*, 2009. 7(3): p. 243-254.
- [23] Sukumaran, S. and A. Sureka, Integrating structured and unstructured data using text tagging and annotation. *Business Intelligence Journal*, 2006. 11(2): p. 8.
- [24] Ajmera, J., et al. A CRM system for social media: challenges and experiences. in *Proceedings of the 22nd international conference on World Wide Web*. 2013. ACM.
- [25] Yaakub, M.R., Y. Li, and J. Zhang, Integration of sentiment analysis into customer relational model: the importance of feature ontology and synonym. *Procedia Technology*, 2013. 11: p. 495-501.
- [26] Hu, M. and B. Liu, Mining and summarizing customer reviews, in *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*. 2004, ACM: Seattle, WA, USA. p. 168-177.
- [27] Boudad, N., et al., Sentiment analysis in Arabic: A review of the literature. *Ain Shams Engineering Journal*, 2017.
- [28] Kaseb, G.S. and M.F. Ahmed, Arabic Sentiment Analysis approaches: An analytical survey.

An Accurate Multi-Biometric Personal Identification Model using Histogram of Oriented Gradients (HOG)

Mostafa A. Ahmad^{1,2}, Ahmed H. Ismail^{1,3}, Nadir Omer¹

¹College of Computing & Information Technology, University of Bisha, Bisha, Saudi Arabia

²Faculty of Computers & Information, Menofia University, Shibeen El-Kom, Egypt

³Faculty of Science, Menofia University, Shibeen El-Kom, Egypt

Abstract—Biometrics is the detection and description of individuals' physiological and behavioral features. Many different systems require reliable personal identification schemes to either prove or find out the identity of an individual demanding their services. Multi-biometrics are required inside the current context of large worldwide biometric databases and to provide new developing security demands. There are some distinctive and measurable features used to distinguish individuals known as Biometric Identifiers. Multi-biometric systems tend to integrate multiple identifiers to increase recognition accuracy. Face and digital signature identifiers are still a challenge in many applications, especially in security systems. The fundamental objective of this paper is to integrate both identifiers in an accurate personal identification model. In this paper, a reliable multi-biometric model based on Histogram of Oriented Gradients (HOG) features of a face and digital signature and is able to identify individuals accurately is proposed. The methodology is to adopt many parameters such as weights of HOG features in merging process, the HOG parameters itself, and the distance method in matching process to gain higher accuracy. The proposed model achieves perfect results in personal identification using HOG features of digital signature and face together. The results show that the HOG feature descriptor significantly performs target matching at an average of 100% accuracy ratio for face recognition together with the digital signature. It outperforms existing feature sets with an accuracy of 84.25% for face only and 97.42% for digital signature only.

Keywords—Biometric identifiers; personal identification; multi-biometric systems; face recognition; digital signature; Histogram of Oriented Gradients (HOG)

I. INTRODUCTION

Biometrics refer to personality check of people as indicated to their physical or behavioral qualities [1]. Biometric technology is a technology that allows people to have a digitally authentication using different parts of their physical bodies [2]. Numerous physical body parts and individual highlights have been utilizing for biometric Systems: faces, digital signatures, and DNA. Person verification in view of biometric highlights has pulled in more consideration in planning security systems [3]. The spread of biometrics is enacted by two parts: technical specialized and the necessity for security. Biometric is critical to keep data secured in our day by day life [4], [5]. The requirement for unique sensors to acquire biometrics was for quite some time thinking about

disadvantages, particularly if multi-biometrics was considered [5]. Nonetheless, no single biometrical feature can meet all the execution necessities in reasonable systems. Face recognition has as of late gotten critical consideration. It assumes a vital part in numerous application regions, for example, human-machine communication, verification, and surveillance. Digital signature are recognized by current standards and legislation as a term that use a key pair of user for sign and verify a document using biometric systems [2], [6]. Today there are many advantages for digital signature such as offers more security than any electronic signature, independent verification cannot be alter by unauthorized parties and long-term retention and access. Face recognition and digital signature have been a long-standing issue in PC vision [7]. As of late, Histograms of Oriented Gradients (HOGs) have turned out to be an effective descriptor as feature extraction for object recognition in general and face recognition and digital signature in particular. Face and digital signature have been a long-standing problem in many applications, especially in computer vision. The main contribution of this paper is to integrate both identifiers in an accurate personal identification model. The proposed multi-biometric model is based on HOG features of a face and digital signature and this model able to identify individuals accurately.

This paper is organized as follows. The related work is describes in Section 2, as well as our topic and techniques. Section 3 introduces the HOG descriptor. In Section 4, we describe the methodology of our multi-biometric system with its different modules. Section 5 explains the description of the dataset used in the testing process as well as the results and experimental analysis. Finally, the main conclusion is drawn in Section 6 with a hint for the future work we have engaged in follow-up this work.

II. RELATED WORK

Biometrics deals with innovations used of gauge human physical or behavioral characteristics to distinguish and perceive people [8]. For the of biometric there are two types features: physiological (e.g. iris, face, unique mark) and behavioral (e.g. voice and digital signature) [9]. The mix of biometric systems, otherwise called "biometric fusion", can be ordered into unimodal biometric in the event that it depends on a single biometric characteristic and multimodal biometric in the event that it utilizes a few biometric qualities for individual verification [8]. A few systems and structures identified with

the mix of biometric systems, both unimodal and multimodal is examine and grouped by a given scientific classification. Face recognition is the most broadly recognized people most of the time by recognizing the face of the individuals and advancement in computing skill over the past few decades [10]. There are three stages namely face detection, feature extraction and face recognition in face recognition system [11]. Techniques for face identification and recognition frameworks can be influence by posture, nearness or nonattendance of auxiliary segments, outward appearance, impediment, picture introduction and imaging conditions. It is difficult to implement a strong face recognition framework, which work in all condition. In computer vision, Face recognition has been a long-standing problem. Recently face recognition system attracted significant attention due to the accessibility of inexpensive digital cameras and computers, and its different applications in biometrics and surveillance [12]. In any case, the wide-run varieties of a human face, because of stance, brightening, and demeanor, result in an exceedingly complex appropriation and fall apart the acknowledgment execution. What's more, the issue of machine recognition of human faces keeps on pulling in scientists from orders, for example, pattern recognition and digital signature [10]. First, the Face recognition system detects the presence of a face in an image. If is found, the system's role is to trace the position of one or more faces in the image.

To make robust use for face recognition, O. Deniz [12] in his study investigated a powerful approach based on HOG features. The use of HOG for face recognition and he used the HOG to extract features from overlapping cells because it is important for this case. Also, applied four databases in the study and obtaining significant result based on FERER database. Also in his paper Alberto [6] proposed HOG-EBGM for face recognition. He used HOG descriptor with three databases and FERET is one them and he obtained better performance by change the properties of HOG to get maximum accuracy of the face graphs acquired compared to classical Gabor-EBGM ones. In their research Bin li [2] introduce LSHOG for face recognition to extract features based on reduce the dimension of the features compared with HOG idea, which distributed an image into several cells, and computed a histogram of gradient orientations over each cell. Unlike traditional HOG their proposed LSHOG tell a histogram over gradient orientations atop the complete photo at each pixel location. Experimental outcomes confirm the feasibility and efficiency of LSHOG and their face recognition method.

Now a days , Digital signature is a popular term that uses a key pair of user for sign and authentication of a document. Using biometric technology professionals can create their digital signature. Biometric technology is a technology that allows people to have an digitally authentication using different parts of their bodies [13]. Mustafa [14] proposed an offline signature verification system based on a signature's local histogram features using classifiers combination of HOG and histogram of local binary patterns (LBP) features. The combination of all classifiers (global and user-dependent classifiers trained with each feature type), achieves a 15.41% equal error rate in skilled forgery test, in the GPDS-160 signature database without using any skilled forgeries in

training. The signature might change over some undefined period and are impacted by physical and enthusiastic states of a subject. The signature may change over some vague time span and are affected by physical and excited conditions of a subject. Further, capable falsifiers may have the ability to imitate signatures that trap the structure Due to outer assembling imperatives in detecting innovations and also innate confinements inside each biometric, no single biometric technique to date can warrant a 100% verification exactness and utilization independent from anyone else [5], [14], [15]. These frameworks are additionally ready to meet the strict execution prerequisites forced by different applications. In [10], 1-median filtering as a spoofing-resistant summed up contrasting option to the entirety administer focusing on the issue of fractional multi-biometric spoofing where m out of n biometric sources to be joined are attacked. Section 3 introduces the HOG.

III. HISTOGRAM ORIENTED GRADIENT (HOG)

The histogram of situated angles was proposed for the utilization of person on pedestrian detection [7]. HOG is a feature extraction strategy that figures the situated gradients of a picture utilizing angle finders. Due to its victories, it has been utilized as a part of numerous PC vision frameworks [1]. For example, it has been utilize for face and on-street vehicle identification applications. It has been connect to face recognizable proof and in addition feeling and gesture recognition. Applicable descriptors assume a critical part in face parameterization. Wavelet, contour lets, and Gabor wavelets have been generally utilize for face recognition. Different parameters like example arranged edge greatness POEM utilizing nearby twofold example LBP and histogram of situated angle HOG have been as of late connect to human location and face recognition. The utilization of introduction histograms has numerous forerunners. Freeman and Roth utilized introduction histograms for hand signal recognition. Histogram of Oriented Gradients (HOG) is highlight descriptors that were first presented by Dalal and Triggs in their CVPR paper [1] to distinguish people on foot in pictures. The examination was then extended to identify human in recordings and creatures and questions in static pictures [16]. In this work, HOG descriptors are extracted to recognize the area of target appearance in face and digital signature images [17]. Fig. 1 shows an image that divided into equal size cells of size 8x8 pixels. Moreover, each cell is initialized with a 9-bin histogram range from 0 to 180 degrees or 0 to 360 degrees. The magnitude and orientation of each pixel are calculated using (1) and (2), where G_x and G_y are the horizontal and vertical gradient, respectively.

$$\text{Magnitude, } |G(x, y)| = \sqrt{G_x(x, y)^2 + G_y(x, y)^2} \quad (1)$$

$$\text{Orientation, } \tan(\theta(x, y)) = \frac{G_y(x, y)}{G_x(x, y)} \quad (2)$$

Consequent to getting the introduction and result, every pixel will vote to the 9-bin histogram similarly to its orientation [9]. The quantity of voting will be chosen by its relative size. Therefore, more grounded sizes will largely affect the histogram. By arranging the extent and direction of every cell into a histogram, we are decreasing the gradient mechanisms down to a vector of only 9 values, which are the sum of sizes

of each bin. In other words, the gradient histogram quantifies the components of every cell in the picture. It is additionally vital to realize that HOG does not keep the data about the gradient or edge positions, but the dissemination of neighborhood power inclination or edge headings.



Fig. 1. Normalized face and its spatial cell by HOG descriptor.

Considering that slope is generally affected by brightening changes, standardization is expected to deal with this issue. Instead of normalizing every histogram separately, the cells are first collect into pieces and standardization in view of the considerable number of histograms in the block. In Dalal and Triggs, any block is work of 2x2 cells as shown in Fig. 1, whereby each blocks cover by half. The histograms of the four cells inside a block are linked into a vector with 36 parts (4 histograms x 9 bins for each histogram) and after that gap this vector by its greatness to normalize it. The most common block normalization used is L2-normalization, as denoted in (3).

$$f = \frac{V}{\sqrt{\|V\|_2^2 + e^2}} \quad (3)$$

To decrease the calculation time of removing HOG descriptors by moving the window overall GPR picture, the concentration of recognition is limited to the areas that contain potential target reflection. We additionally control the measure of the picture so hyperbolas can be recognizing at various scales. This is vital as the span of hyperbolas shifts in like manner to a few viewpoints like the measurement of the objective, nature of the medium and the setting of the face and signature system itself.

In our work, each spatial cell is square of 8x8 pixels. This size is selected based on the distance between eyes of the normalized faces, which in our work is 32 pixels and also in our work we used different values to get our result such as the number of bins is choose with different values 9, 12 and 15 also the block size is 1, 2, 4 and 8. Finally, the cell size also they have different values like 4, 8, 12 and 16.

IV. METHODOLOGY OF OUR MULTI-BIOMETRIC SYSTEM

Different Multi-biometric systems share a public general flow as shown in Fig. 2, which is described the four main mechanisms:

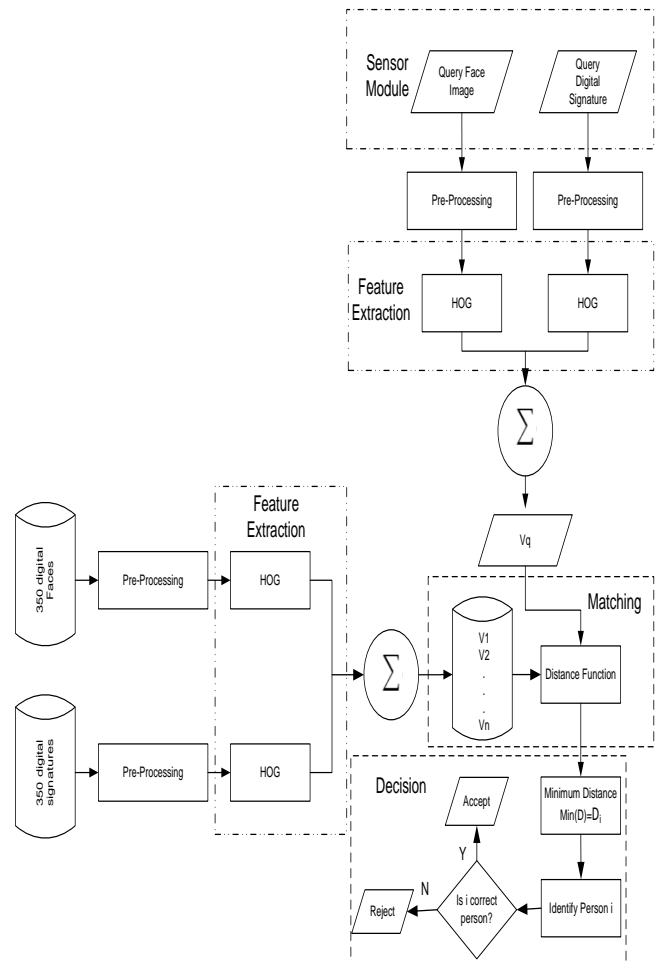


Fig. 2. Block Diagram of suggested multi-biometric system.

- Sensors:** The first segment of a multi-biometric system is obtaining of the biometrics information of a person from biometric sensors hardware. For face recognition and digital signature, the sensor is regularly a camera, the sensor is commonly a scanner, for face information, the sensor is a mouthpiece. The nature of the procurement module has a significant effect on the execution of the system that is delicate to the ecological conditions (i.e. changes in the brilliance of a picture), nature of sensor (i.e. dpi of the picture), human factor (i.e. posture varieties).
- Feature extraction module:** The gain information is pre-prepared to expel commotion or different anomalies present and afterward subjected to the component extraction process with a specific end goal to extricate biometrical values that in a perfect world must depict extraordinarily an individual, so biometric information gathered from one individual, under various circumstances, are “comparable”, while those gathered from various people are “desperate”. For instance, the position and introduction of particulars focus in a fingerprint picture are utilized as a part of a fingerprint framework. The highlights removed amid enlistment are put away in a format, which is a potent little and

simple to process. Keeping in mind the end goal to enhance interoperability among various biometric frameworks there exist recommendations of the standard configuration of layouts, i.e. for fingerprint they are constructed just with respect to particulars focuses.

- **Matching module:** In this module, which is not used during enrollment, the feature values from an unknown individual are compared to those in the stored template by generating a matching score which shows the level of similarity between a set of biometrics data. The score should be good enough for features from the same individuals and unacceptable for those from different ones. In a signature system, for example, return the number of matching minutiae points between the query and the template can as a matching score. Usually matching is a difficult pattern-recognition problem due to large intra-class variations (caused by bad acquisition, noise, varying environmental conditions, alterations, etc.) and large inter-class resemblance (i.e. differentiating identical twins is always difficult in face recognition).
- **Decision component:** In this module, the user's identity is established (identification) or a claimed identity is accepted/rejected based on the matching score. Usually, the final decision is taken by comparing the matching score to a fixed level, which is selected according to consideration of the degree of security required by the application.

The methodology of this paper is to study previous systems in personal identification using multi-biometrics then select the most suitable biometric identifiers such HOG to use multiple identifiers for person identification, finally Measure the accuracy of the developed algorithm according to standard and real data sets. Verify the ability of the developed algorithm to work in real-time.

V. EXPERIMENTAL RESULTS AND ANALYSIS

In this section, we present the database of faces and a digital signature, which our experiments are, taking place. After that, we discuss the primary results of our multi-biometric identification system.

A. Datasets

Training and testing of a biometric system dataset plays a main role in achieving better recognition performance. All the experimentations dataset carried out in this paper using two set of images. The first one we got sample dataset for faces from the Face Recognition Technology (FERET) program. The FERET database is a standard database of face imagery which was necessary to evaluate the FERET program, both to supply standard imagery to the algorithm developers and to supply a sufficient number of images to allow testing of these algorithms [9]. FERT dataset contains 3365 full frontal facial images of nearly 1000 subjects. FERT dataset images are organized into a gallery set (fa) and four probe sets (fb, fc, dup1, dup2). The second dataset we added some from local images on the way to build our face database from individuals in Bisha University. Our dataset may contain color or

grayscale images with different resolutions and file formats. Our experiment is considered 50 persons with 7 face images for each person, a total of 350 faces. The seven samples for each person contain full frontal face views, head rotation, different emotions, and even different clothes or background as shown in Fig. 3. The offline signature database consists of 160 individuals' signatures: each individual has 7 genuine signatures, and 7 forgeries of his signature. The 7 genuine samples of each finger were collected in a single writing session. All signatures have binary bitmap picture format, with 300 dpi resolution, as shown in Fig. 4.



Fig. 3. A sample of faces dataset.

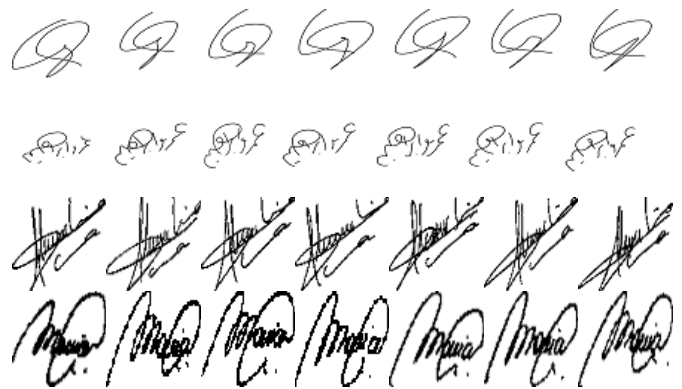


Fig. 4. Sample of digital signature dataset.

B. Experimental Results

We start our experiment by examining the HOG feature extractor and descriptor on faces and digital signature separately. Preliminary results for faces only were not encouraging, while it ranges between 80% and 84% of identification accuracy; the digital signature was ranged between 96% and 98%. Since we propose to use multi-biometric in order to increase the identification accuracy, so we combine both of face and digital signature together to produce a single HOG feature vector f as a weighted sum as shown in (4):

$$f = \alpha f_1 + (1-\alpha) f_2 \quad (4)$$

Where, $0 \leq \alpha \leq 1$, f_1 is the HOG feature vector for the face, f_2 is the HOG feature vector for the digital signature.

To measure the matching between feature vectors of faces or digital signatures, we use some of the distance functions. Some most popular distance functions are Manhattan, Euclidean, Angle-Based, and modified Manhattan. To show the concept of different functions, let x and y be two feature vectors of length n , then we can calculate the following distances between these feature vectors as in the following (5)-(8):

Manhattan distance:

$$d(x, y) = L_{p=1}(x, y) = \sum_{i=1}^n |x_i - y_i| \quad (5)$$

Euclidean distance:

$$d(x, y) = L_{p=2}(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (6)$$

Angle – based distance:

$$d(x, y) = -\cos(x, y)$$

Where,

$$\cos(x, y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2}} \quad (7)$$

Modified Manhattan distance:

$$d(x, y) = \frac{\sum_{i=1}^n |x_i - y_i|}{\sum_{i=1}^n |x_i| \sum_{i=1}^n |y_i|} \quad (8)$$

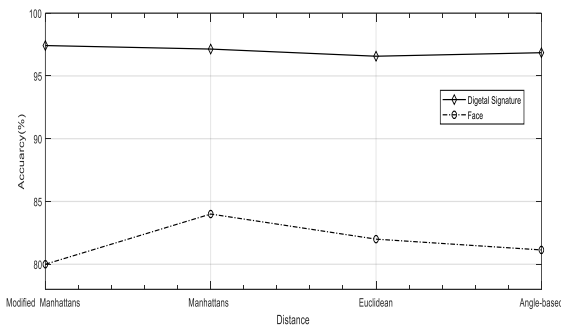


Fig. 5. Performance comparison of HOG among the different distance functions for face and digital signature features in accuracy ratio of identifying.

In Fig. 5, the digital signature identification has outperformed the face in the accuracy ratio all over the distance functions. This is because of the different challenges in face feature detection and description. The accuracy ratio of identification is maximized by the Manhattan distance function in both the face and digital signature. It produces 96.99% in the digital signature, and 81.78% in the face on average. In order to achieve better accuracy in individual identification, and benefit from both face and digital signature features, we proposed to merge the two feature vectors of the face and digital signature using different weights according to (4). The results of comparing the accuracy ratios among as shown in Table I, different distance functions at different α values (from 0.01 to 0.1). It is shown that the Manhattan distance function is outperformed the other functions for all the values of parameter α . The accuracy ratio is at its maximum value (99.43%) at $\alpha = 0.06$ for the Manhattan distance function, which means less weight for the face features in comparing with the digital signature features. The nearest distance function to Manhattan matching accuracy is the modified Manhattan function.

TABLE I. THE ACCURACY RATIOS IN PERCENTAGE FOR DIFFERENT MATCHING DISTANCES AND VARIOUS α WITH A NUMBER OF HOG BINS EQUAL TO 9

Distance \ α	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09	0.1
Modified Manhattan	97.43	97.71	98.57	98.86	99.14	98.86	98.86	98.29	98.29	97.71
Manhattan	98	98.86	98.86	98.86	99.14	99.43	99.14	99.14	98.86	98.57
Euclidean	96.57	96.57	96.57	96.57	96.57	96.57	96.57	96.57	96.57	96.57
Angle-based	96.57	96.57	96.57	96.86	96.86	96.86	96.86	96.86	96.86	97.14

Some parameters in HOG are affecting its performance, such as the number of bins, cell size, and block size. To reach the optimal value of each parameter, we try to change one parameter while fixing some others. Starting with the bins number, Fig. 6 measures the performance of HOG at a different number of bins with different α values for Modified Manhattan distance function, a cell size of 8, and block size equal to 2. It shows that the best accuracy ratio is acquired at the number of bins equals to 9 and 15 at α equals to 0.05, and 0.04, respectively. It is a little disturbance to give the same accuracy ratio at different bin values at different α .

Fig. 7 and 8 introduces the effect of changing the number of bins on the HOG performance for the Manhattan distance function at different α values, cell size = 8, and block size = 2. The results show that the Manhattan distance function has superior performance at $\alpha = 0.06$ and number of bins equal to 12. We get an accuracy ratio of 100% at these parameters.

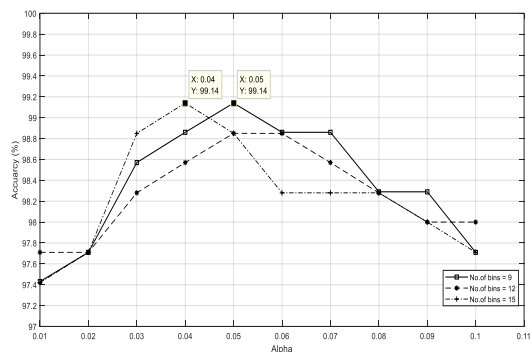


Fig. 6. Performance comparison of HOG at a different number of bins with different α values for Modified Manhattan distance function.

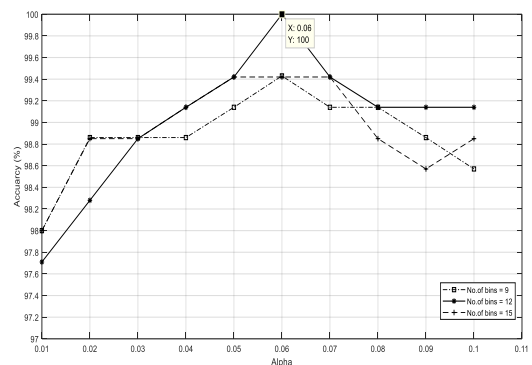


Fig. 7. Performance comparison of HOG at a different number of bins with different α values for Manhattan distance function.

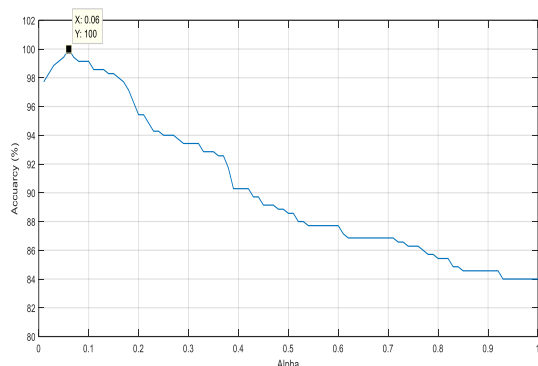


Fig. 8. Accuracy values of HOG at different values of α , and with number of bins equal to 12 for Manhattan distance function.

From this point up, we will use the Manhattan distance function as the best function matches feature vectors, with the number of bins equal to 12 and $\alpha = 0.06$. Now, the cell and block size is checked with the same number of bins and α . Tables II and III show the effect of change the cell size and block size in the HOG performance, respectively. It is shown that the optimal number of cell size is equal to 8, while the block size is 2. The accuracy of individual identifying is still 100% at $\alpha = 0.06$ and a number of bins is equal to 12.

TABLE II. THE ACCURACY RATIOS IN PERCENTAGE FOR DIFFERENT CELL SIZES AND VARIOUS α WITH A NUMBER OF HOG BINS EQUAL TO 12

α / Cell Size	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09	0.1
4	96	96.57	96.57	96.85	97.71	98	97.42	98	98	98
8	97.71	98.28	98.85	99.14	99.42	100	99.42	99.14	99.14	99.14
12	98.57	99.14	99.14	99.14	98.85	98.57	98.57	98.28	98	97.71
16	98.57	99.14	98.85	99.42	99.14	99.14	99.14	98.57	97.71	97.71

TABLE III. THE ACCURACY RATIOS IN PERCENTAGE FOR DIFFERENT BLOCK SIZES AND VARIOUS α WITH A NUMBER OF HOG BINS EQUAL TO 12 AND CELL SIZE EQUAL TO 8

α / Block Size	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09	0.1
1	96	96	95.42	96	96	96	96	96	96	96.28
2	97.71	98.28	98.85	99.14	99.42	100	99.42	99.14	99.14	99.14
4	98	99.42	99.42	99.42	99.42	99.71	99.71	99.71	98.85	98.57
8	81.42	81.42	81.42	81.42	81.42	81.42	81.42	81.42	81.42	81.42

Finally, in comparing the performance of face or digital signature features only with the merging of features between both of them as proposed, Fig. 9 shows the outperformance of multi-feature proposed method.

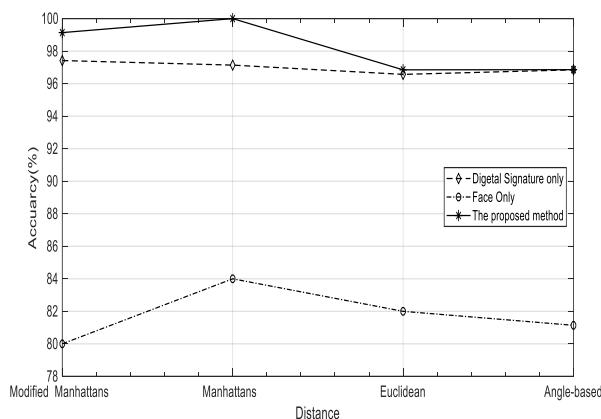


Fig. 9. Performance comparison of HOG at a different distance methods for face only, digital signature only, and the proposed method, at number of bins equal to 12, and with $\alpha = 0.06$.

VI. CONCLUSION

This study presents a new hyperactive system that depends on HOG descriptor as features extraction for face recognition and digital signature together. Multi-biometric personal identification model using Histogram of Oriented Gradients (HOG) as feature extraction for face recognition and digital signature was present in this paper. The contributions are threefold: Firstly, to provide robustness to facial and signature feature detection, we propose to uniform sample the HOG as features. The result presented that our method performs better result for the multi biometrics system based on face recognition and digital signature instead of using them as an individual. Secondly, the matching result for our methods shows better result compared to other face recognition and digital signature only. This better performance is explained by the properties of HOG descriptor that is more robust for the hybrid. Finally, the results show that the HOG feature descriptor significantly performs target matching at an average of 100% accuracy ratio for face recognition together with the digital signature. It outperforms existing feature sets with an accuracy of 84.25 % for face only and 97.42% for digital signature only. In near future, we hope to apply the deep-learning approaches for feature extraction instead to HOG. This way, we hope to gain a big range for Alpha (α) selection with 100% accuracy.

ACKNOWLEDGMENT

This study has been supported by research support program, University of Bisha, Kingdom of Saudi Arabia, grand number (UB-06-1438).

REFERENCES

- [1] Dalal, N. and B. Triggs. *Histograms of oriented gradients for human detection*. in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. 2005: IEEE.
- [2] Rahmawati, E., et al. *Digital signature on file using biometric fingerprint with fingerprint sensor on smartphone*. in *Engineering Technology and Applications (IES-ETA), 2017 International Electronics Symposium on*. 2017: IEEE.
- [3] Benaliouche, H. and M. Touahria, *Comparative study of multimodal biometric recognition by fusion of iris and fingerprint*. *The Scientific World Journal*, 2014. **2014**.
- [4] Jain, A.K., A. Ross, and S. Prabhakar, *An introduction to biometric recognition*. *IEEE Transactions on circuits and systems for video technology*, 2004. **14**(1): p. 4-20.
- [5] Kataria, A.N., et al. *A survey of automated biometric authentication techniques*. in *Engineering (NUICONE), 2013 Nirma University International Conference on*. 2013: IEEE.
- [6] Hernandez-Ardieta, J.L., et al., *A taxonomy and survey of attacks on digital signatures*. *Computers & security*, 2013. **34**: p. 67-112.
- [7] Karaaba, M., et al. *Robust face recognition by computing distances from multiple histograms of oriented gradients*. in *Computational Intelligence, 2015 IEEE Symposium Series on*. 2015: IEEE.
- [8] Lumini, A. and L. Nanni, *Overview of the combination of biometric matchers*. *Information Fusion*, 2017. **33**: p. 71-85.
- [9] Janbandhu, P.K. and M.Y. Siyal, *Novel biometric digital signatures for Internet-based applications*. *Information Management & Computer Security*, 2001. **9**(5): p. 205-212.
- [10] Chelali, F.Z. and A. Djeradi. *Zernike moments and histogram of oriented gradient descriptors for face recognition from video sequence*. in *Complex Systems (WCCS), 2014 Second World Conference on*. 2014: IEEE.
- [11] hihaoui, M., et al., *A survey of 2D face recognition techniques*. *Computers*, 2016. **5**(4): p. 21.
- [12] Déniz, O., et al., *Face recognition using histograms of oriented gradients*. *Pattern Recognition Letters*, 2011. **32**(12): p. 1598-1603.
- [13] Tripathi, S.K. and B. Gupta, *An Extension to Modified Harn Digital Signature Scheme with the Feature of Message Recovery*, in *Networking Communication and Data Knowledge Engineering*. 2018, Springer. p. 183-193.
- [14] Omara, I., et al. *Discriminative Local Feature Fusion for Ear Recognition Problem*. in *Proceedings of the 2018 8th International Conference on Bioscience, Biochemistry and Bioinformatics*. 2018: ACM.
- [15] onnor, P. and A. Ross, *Biometric recognition by gait: A survey of modalities and features*. *Computer Vision and Image Understanding*, 2018. **167**: p. 1-27.
- [16] Lee, K. and M. Mokji. *Automatic target detection in GPR images using Histogram of Oriented Gradients (HOG)*. in *Electronic Design (ICED), 2014 2nd International Conference on*. 2014: IEEE.
- [17] Mathew, S. and G. Saranya. *Advanced biometric home security system using digital signature and DNA cryptography*. in *Innovations in Green Energy and Healthcare Technologies (IGEHT), 2017 International Conference on*. 2017: IEEE.

Detection of Mass Panic using Internet of Things and Machine Learning

Gehan Yahya Alsalat, Mohammad El-Ramly, Aly Aly Fahmy
Faculty of Computer and Information
Computer Science Department, Cairo University
Cairo, Egypt

Karim Said, M.D.
Faculty of Medicine,
Cardiovascular Department
Cairo University
Cairo, Egypt

Abstract—The increase of emergency situations that cause mass panic in mass gatherings, such as terrorist attacks, random shooting, stampede, and fires, sheds light on the fact that advancements in technology should contribute in timely detecting and reporting serious crowd abnormal behaviour. The new paradigm of the ‘Internet of Things’ (IoT) can contribute to that. In this study, a method for real-time detection of abnormal crowd behaviour in mass gatherings is proposed. This system is based on advanced wireless connections, wearable sensors and machine learning technologies. It is a new crowdsourcing approach that considers humans themselves as the surveillance devices that exist everywhere. A sufficient number of the event’s attendees are supposed to wear an electronic wristband which contains a heart rate sensor, motion sensors and an assisted-GPS, and has a wireless connection. It detects the abnormal behaviour by detecting heart rate increase and abnormal motion. Due to the unavailability of public bio-dataset on mass panic, dataset of this study was collected from 89 subjects wearing the above-mentioned wristband and generating 1054 data samples. Two types of data collected were: firstly, the data of normal daily activities and secondly, the data of abnormal activities resembling the behaviour of escape panic. Moreover, another abnormal dataset was synthetically generated to simulate panic with limited motion. In our proposed approach, two-phases of data analysis are done. Phase-I is a deep machine learning model that was used to analyze the sensors’ collected readings of the wristband and detect if the person has indeed panicked in order to send alerting signals. While phase-II data analysis takes place in the monitoring server that receives the alerting signals to conclude if it is a mass panic incident or a false positive case. Our experiments demonstrate that the proposed system can offer a reliable, accurate, and fast solution for panic detection. This experiment uses the Hajj pilgrimage as a case study.

Keywords—Internet of Things; IoT; Mobile Crowd Sensing (MCS); wearables; mass panic; mass gatherings; accelerometer; Optical Heart Rate (HR) sensor; abnormal crowd behaviour; deep learning; Recurrent Neural Network (RNN); Long Short Term Memory (LSTM); Gated Recurrent Unit (GRU); time series

I. INTRODUCTION

The increase of emergency situations cause mass panic, such as terrorist attacks, random shooting, stampede, natural catastrophes, and fires requires fast detection and swift action to save lives. Advances in technology can greatly contribute to the timely detection and reporting of serious crowd abnormal behaviour. Early detection of an incident will give the

emergency authorities valuable time to deal with the situation and prevent it from getting worse by implementing immediate and possibly automated actions.

Mass gathering is an event involving the gathering of a large number of people, at least 1000 but can rise up to millions, at a specific location for a defined period of time and for a specific purpose, it can be either organized or unplanned [1]. Mass gatherings are exposed to unpleasant incidents due to the large number of people and limited space and exit routes. Examples of mass gatherings are major sporting, religious, and cultural events (e.g. Olympic Games, religious pilgrimages, etc.).

A. Mobile Crowd Sensing (MCS)

Mobile Crowd Sensing (MCS) is a new sensing paradigm based on collecting real-time data from two participatory sources: sensing and social media platforms. Therefore, it allows ordinary users to contribute by sharing real-time data with data sensed or collected from their mobile devices/wearables. MCS collects data from users’ devices to analyze them and identify spatiotemporal patterns [2].

The revolution of cost-effective hardware, emergent computing and communication trends such as the Internet of Things (IoT), big data, machine learning, wearable sensing and cloud computing have enabled the existence of mobile crowd sensing applications and made our environment smarter.

B. Internet of Things

Internet of Things (IoT), also known as machine-to-machine (M2M) is a paradigm in which smart sensors and devices collect data and interact with one another without human intervention. Gartner expects that “8.4 billion connected things will be in use worldwide in 2017, up 31 percent from 2016, and will reach 20.4 billion by 2020” [3]. Owing to the rapid growth and advancement in the field of IoT wearable sensors, it became possible to monitor physiological signals continuously, accurately and in a real-time manner [4].

C. Problem Statement and Motivation

Despite the best efforts of authorities to secure mass gatherings, unfortunate incidents still occur in such occasions and cause loss of lives. The causalities of any disaster may not be prevented. However, the early detection in a timely manner

may enable swift actions and reduce the losses. The minutes or even the seconds can contribute to saving people lives.

The main motivation for our work is to overcome the limitations of the existing approaches of detecting mass panic such as video surveillance, audio surveillance, and human surveillance. The details of their shortcomings are explained in Section III. So our proposed system will leverage the popularity of smart wearables and the power of crowd-sourcing, participatory sensing to enable a new wave of effective techniques to detect and report incidents in a timely manner.

D. The Contribution of this Paper

- This paper is – to our best knowledge – the first proposal that employs human physiological bio-data (heart rate and accelerometer) in a mobile crowd sensing (MCS) and IoT application; to contribute to real-time mass panic detection in a mass gatherings.
- We believe that, our dataset collection and data preprocessing is an important contribution that can be a good seed for similar studies. Due to ethical and impractical limitations restricting collecting a real mass panic dataset, it is inevitable to produce artificial datasets reflecting the typical panic behaviour. We plan to disclose our dataset for benchmarking and to make it available for researchers.
- Data analysis phase-I in the wristband: Detecting the panic of an individual was done by analyzing his/her heart rate with respective motion. This problem is a time series problem. We used promising sequential deep machine learning models; Recurrent Neural Networks (RNN), particularly Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU) as time series classification detection algorithm and we compared their performance with non-sequential models, Vanilla Neural Network (VNN) and Support Vector Machine (SVM).
- Data analysis phase-II: Currently, we propose a simple data analysis module that works on the monitoring server. It can monitor and cluster the incoming signals and differentiates between real mass panic and false positive. This proposed module has data visualization capabilities which can provide event's officials with a real-time view of one or more critical incidents at the same screen. This cannot be achieved using conventional surveillance approaches. Therefore, officials can quickly grasp the threats and take proper fast actions.

E. Paper Scope

It is important to emphasize the following: (1) the scope of this paper is mainly proposing a new detection method based on MCS and IoT, collecting and analyzing experimental data, and running a proof of concept experiment to validate this proposed approach. (2) We did not cover or experiment the details of various possible types of communication networks that may be used to send/receive the sensor data. (3) In this study, we focused on incidents that happen in rapid onset

disasters and occur swiftly in a mass gathering within seconds or minutes and often without any warning. (e.g., terrorist attacks, random shooting, stampede, and fires). (4) Data analysis phase-II prototype is under development and it is planned in future work.

The remainder of this paper is organized as follows. Section II describes the characteristics of mass panic. Section III describes the related work of detecting mass panic. In Section IV, we explain our proposed system architecture and modules. Section V is the case study that includes the experiments and results. Section VI is the discussion section. We conclude the paper in Section VII, and finally, recommend future work.

II. CHARACTERISTICS OF MASS PANIC

Mass Panic is “type of anomaly in a human crowd, which appears when a group of people start to move faster than the usual speed. Such situations can arise due to a fearsome activity near a crowd such as stampede, fire, fight, robbery, riot, etc.” [5]. Here we describe the characteristics of mass panic. Studying such characteristics helps developing detection approaches that exploit technologies which can discover these characteristics. These characteristics can be categorized or explored from two perspectives; first, from physical behaviour of the mass people and second, from physiological factors.

A. The Physical Characteristics of Mass Panic

The physical behaviour characteristics are reflecting the fight and flight response and they are highlighted by Helbing *et al.* and by Bracha [6], [7]. They include (1) People move faster than normal. (2) Individuals start pushing and interactions among people become physical in nature. (3) The physical interactions in the jammed crowd cause dangerous pressures up to 4,450 Nm⁻¹ (refs which can bend steel barriers or push down brick walls). (4) People show a tendency towards mass behaviour, that is, to do what other people do. (5) Many studies [8]-[10] has shown that during disasters, humans evolved to fight, flight, or freeze; the phrase “fight or flight” normally describe the key behaviours that are triggered by fear and occur as a defensive reaction to threat. The temporal relations are summed up by Bracha [7] in terms of “freeze-flight-fight-fright-faint” and he has ordered these sympathetic responses as they occur. A freeze response is believed to be the first response to a threat by some individuals. freeze response was mentioned by Leach [9] as the “stop, look, and listen” to assess a dangerous situation and sometimes it is called “play dead.”. Fight response is when acting impulsively and combat to prevent something undesirable. Flight response is escaping very fast toward a safe place or to run away from serious danger.

B. The Psychological Characteristics of Mass Panic

Here we discuss the characteristic of escape panic from psychological factors. It is mainly about several biological changes in our body that are mediated by the autonomic nervous system. During severe fear, there are significant changes in our physiological measures. For example, the heart rate increases, our breathing becomes rapid, our muscles tense[11]. The sympathetic nervous of the body stimulates the

adrenal glands, which in turn trigger the release of adrenaline and nor-adrenaline hormones. This is what causes the body to respond to stress or fear and as a trigger for the fight-or-flight response which eventually increases the heart rate, breathing rate and movement[12].

III. RELATED WORK

A. General Approaches of Mass Panic Detection

Firstly, we reviewed the possible existing and emergent techniques for mass panic detection in mass gatherings. We classified them into two categories; first, surveillance based techniques; Second, mobile crowd sensing-based techniques.

1) Surveillance based Techniques

Someone or something is monitoring the crowd in the event and reporting any abnormal behaviour to event's authorities. Examples of these approaches are video surveillance, audio surveillance, and human surveillance.

These three approaches have common limitations such as; they are mostly suitable for limited spaces and they cannot cover the whole space for either cost or privacy concerns. Furthermore, each one of these approaches has limitations in the detection method itself. In the video surveillance: (1) If the cameras are smart, they may fail to detect the abnormality in a high dense crowd due to ambiguities and severe occlusions. While if the cameras are traditional; they require human operators to monitor the surveillance cameras continuously over a long period of time, human are subjected to fatigue and loss their attention. (2) The cameras are almost useless during night-time, due to sudden illumination changes. (3) The high computational requirements of the detection algorithms to work in real life [13], [14]. In Audio surveillance, some events may have sounds of very high volume. This makes detection of abnormal sounds difficult [15]. In Human surveillance, the public safety officers are humans that are subject to the same dangers that cause mass panic. They themselves might be victims of the attack or incident that took place or they might panic out of fear for their lives and forget/fail to report the incident.

2) Mobile Crowd Sensing based Techniques

This technique depends on participatory crowd sensing; when normal people who are attending an event contribute to detecting and reporting incidents using their own mobile phones in a timely manner. The Advantages of this approach is that it is cheap and it takes an advantage of attendees' resources (e.g. Mobile phones) and availability. Examples of this approaches are: (1) Contacting emergency authorities by calling 911 [16], [17]. (2) Using social media to report the incident (e.g. posting in Twitter)[18-20]. The limitations of these approaches are: (1) They may not provide real-time detection but near real-time, usually within minutes of the incident, this depends on how fast people will report the incident. (2) Due to the awe and shock of the incident, people might be too busy trying to escape and save their lives. Also, it is possible that they lose their mobile phones during the hide or escape.

It is important to emphasize that, our proposed system is categorized in this approach. The key features of our proposed

system over the previously mentioned approaches are the followings: (1) it provides real-time and fast detection because our technique does not require deliberate action from the user at the time of the incident; the detection and the reporting are done automatically. (2) It covers the whole event premises and it has a good visualization capabilities, where a threat or more than one can be detected and visualized in one screen.

B. Biosensors based Techniques

Secondly, we reviewed the studies that use bio-sensors to detect escape panic. It is important to highlight that – to our best knowledge – there is no real deployment of a large-scale bio-crowd sensing for mass panic detection. But there are a scarce number of proposals or researchers that propose a method of detecting abnormal behaviours of individuals or a small number of people in small buildings.

The U.S. Patent (No. 8,477,035) [21], is a proposed security system triggered by heart rate detection; it is proposed for banks and other enterprises which use security systems in their facilities, to detect robbery or other types of threatening activities. The method includes a wearable heart rate sensor that monitors the heart rate of an individual (security guard or employee); determining if the detected heart rate is abnormal; generating a signal which is wirelessly transmitted to the central monitoring station. The detection is based on threshold-based technique.

The U.S. Patent (No. 9,858,794) [22] is a proposed system for checking the state of workers who are performing tasks in a hazardous environment. The proposed sensing device includes a plurality of biometric sensors to report the corresponding hazardous state of a person. The detection is based on threshold-based technique.

Using accelerometer data, a system [23] with the name Emergency Rescue Support System (ERESS) proposed a disaster detection algorithm. It uses plural sensors such as acceleration, angular velocity, and geomagnetism using data from sensors mounted on ERESS terminals such as smartphones or tablets. ERESS detects the disaster from the behaviour analysis of people by the sensors using machine learning based on group learning using support vector domain description (SVDD). ERESS operates under mobile ad-hoc networks (MANET) between the hand-held terminals.

Using electro-dermal activity and acceleration (skin conductance) the work [24] proposed an algorithm which detects people's extraordinary condition using Skin Potential Level (SPL) sensor (SPN-01) which can measure Skin Potential Level (SPL) with a 9 axis motion sensor to detect an extraordinary condition in case of panic. An electro-dermal activity is a parameter which shows man's mental activity condition and they are electrical signals that appear on an outer skin and sweat glands.

The key differences between our proposed system and the previously mentioned studies [21]-[24] are the followings: (1) They use threshold-based approaches which were designed to detect abnormality of individuals and were not designed for the mass gatherings where knowing the abnormal threshold of a huge number of attendees in advance is challenging. (2) They are designed for small buildings compared to our

proposed system which is intended to cover a very large space using long-range wireless networks. (3) Our proposed system is fusing heart rate sensor with accelerometer sensor. (4) We employed state-of-the-art deep learning models, namely, (Recurrent Neural Network RNNs) and compare them with conventional machine learning approaches.

IV. PROPOSED SYSTEM BASED ON IOT

Our proposed system consists of three modules as follows:

A. Data Acquisition and Data Analysis Phase-I Module

A sufficient number of event attendees are supposed to wear the electronic waterproof wristband. This wristband contains HR and motion sensors, Also it will be embedded by an A-GPS (Assisted-Global Positioning System) for finding the location of a person (by considering the longitude and latitude). The HR and motion sensors are continuously monitoring the heart rate and movement of a person. Once a panic is detected, an alerting signal will be transmitted to the central monitoring server containing the following: User Unique Identifier (UID), location coordinates, heart rate reading, motion data, and timestamp. The monitoring server will distinguish whether the coming signals are a false positive or a real panic incident and then will report the incident severity.

Our assumption is that the wristband will sense the user's heart rate and motion passively, (not connecting to the monitoring server for 24/7) and it will start a connection and send alerting signals once the panic is detected. We have built a machine learning detection model that classifies sensors data as normal or abnormal. Fig. 1 is illustrating the flow chart of the data analysis that occurs in the wristband device. The system modules and the components of the wristband are illustrated in Fig. 2.

B. Communication Module

Every wristband will be embedded with communication modem for wireless connection. The proposed system can use the existing cellular network infrastructure for the wireless communication-between the wristband and the monitoring server- for the following reasons: (1) Service providers have already covered the event area with such networks (GSM/GPRS, 2G, 3.5G or 4G); no need for extra cost of deploying other types of wireless sensors network (WSN). (2) Cellular networks are long-range wirelesses which are able to cover almost all event area.

C. Data Analysis Phase-II and Visualisation Module

In this module, data analysis phase-II is handled by the monitoring server that receives the signals coming from the attendees' wristbands. According to our assumption that was mentioned previously, this monitoring server will receive only the alerting signals. If there is a massive number of alerting signals from the same location, a cautionary warning will be issued and a pop-up colorful pinpoint will appear in the event map at that particular location, with a counter that shows the number of alerting signals per second and this number may increase or decrease depending on the development of the situation.

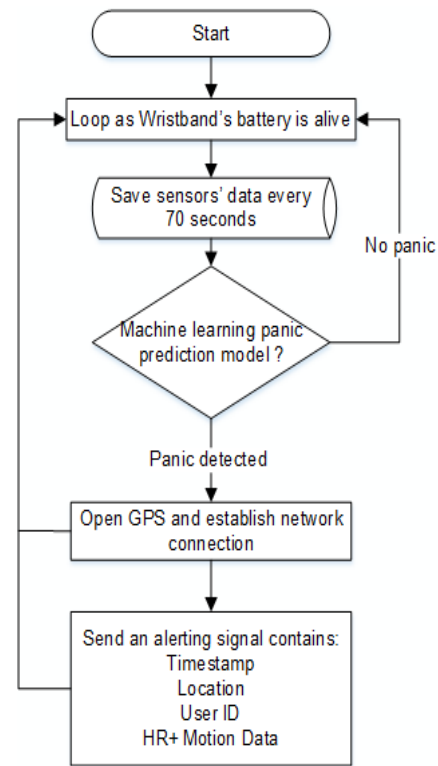


Fig. 1. Flow chart of the data analysis phase-I module.

The color of the pinpoint indicates the severity of the problem; the green color means no threat, it may be a “false positive” alarm, while the yellow color means a caution that should be considered, as there is a probability of having a serious problem. However, when the color goes red, this means a very “critical problem” is going on. Every color represents a range of numbers depending on the number of the event's attendees and the location depth and width. These analytics and data visualization capabilities provide event's officials with a real-time view of critical incidents on one screen and therefore they can take proper action and fast response (see Fig. 3).

V. CASE STUDY: PANIC DETECTION DURING HAJJ PILGRIMAGE

Hajj is the fifth pillar in Islam, It is an annual religious duty which lasts for six days, and it must be performed at least once in a Muslim's lifetime, only for those who are physically and financially capable to afford it [25]. It is one of the most congested mass gatherings in the world. Millions of pilgrims are performing the Hajj rituals in Saudi Arabia and they often crowd into packed spaces at densities of about six or seven persons per square meter [26]. Such crowd densities present a huge challenge for both public safety authorities and any computer-aided system to predict and detect crowd disasters. Therefore we chose Hajj pilgrimage as our case study. In the following, we described: (1) Hardware details. (2) Data collection and preprocessing. (3) Machine learning models. (4) Result and performance metrics.

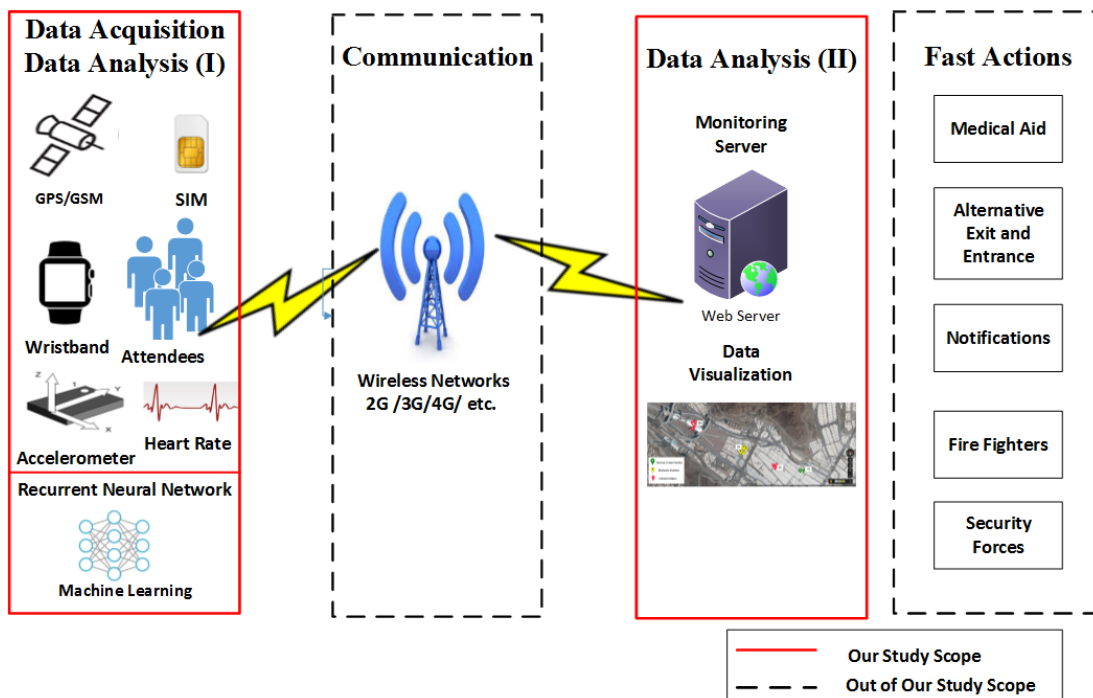


Fig. 2. Proposed system components and modules.

A. Hardware Details

TomTom Spark Cardio wristband [27] was used in this study as our bio-data collection device (see Fig. 4). It is a waterproof wristband that has a built-in optical heart rate sensor that measures the heart rate beats per minute (bpm) using Photoplethysmography (PPG) optical technology. The heartbeat is measured by using a light that measures changes in the blood flow. This is done on the top of the wrist by shining light through the skin. The PPG sensor monitors changes in the light intensity via reflection from or transmission through the tissue and detecting the changing light reflections [28]. Furthermore, it also has an x-axis accelerometer and a Gyro-meter. Moreover, it records the heart rate value with respective acceleratory data in one second interval time. The data is extracted as a csv file.

B. Data Collection

Unfortunately, there exists no publicly available dataset on mass panic as a bio-data. Producing such real dataset is not possible in the real world since it often requires exposing real people to the actual, possibly dangerous environment. These pose ethical and safety concerns. Thus we proposed in this paper, a way to produce an artificial/pretended data on

abnormal activities reflecting on typical behaviour of escape panic. Two types of data were collected, firstly the data of normal daily activities, and secondly the data of abnormal activities (subjects pretended a typical behaviour of escape panic) and another abnormal data was synthetically generated.

C. Subjects and Activities Details

The data collection and subjects' selection were supervised by the fourth author. We consulted him on the physiological factors that cause heart rate fluctuation and also about the data collection activity types, to ensure that we cover all possible normal and abnormal activities that reflect the typical normal and abnormal activities during Hajj rituals. This ensures that the machine learning will learn all possible scenarios to reduce false negative and false positive rates. We conducted our experiments on 89 individuals that are fit, healthy and medication free, aged between 20 and 48 years. This age group represents a good portion of human-beings. Therefore, in our proposed system we targeted this age group (fit and medication free) to be as human sensors for the following reasons: (1) People in this age group are proactive and have the ability to respond to danger and move faster. (2) They are medication free, so there are no other factors that increase HR rather than the real panic situation.



Fig. 3. Data visualization, data analysis phase-II module.



Fig. 4. Tom tom spark cardio wristband.

All participants were informed about the experimental setup and they agreed to contribute to this study. Subjects were asked to perform the following activities: (1) To pretend the escape panic by doing abnormal actions (e.g. running suddenly and rapidly); their heart rate and respective movement data are recorded in one second interval time. It is important to emphasize that, from a physiological perspective, this pretended dataset is valid and represents the escape panic because in the real escape panic the value of HR and movement will be almost the same or mostly much higher (due to the release of more adrenaline and nor-adrenaline hormones). (2) To perform normal daily activities (without heavy physical activities such as playing sport). (3) We found that climbing stairs or a high hill will potentially increasing the heart rate while it is not a panic case. So we collected such data to enable the machine learning model to learn the correlation between the heart rate and the movement of climbing stairs as a normal activity. (4) On the other hand, we manually synthesized 112 abnormal activity samples that reflect the freezing action; when frightened people are freezing in their place and do not escape; but for sure their heart rate is pumping very fast. So we produced these synthetic datasets and we took in our consideration the following factors; First, we set the HR values to be

extraordinary from normal, at the same time do not exceed $(220 - \text{age of a person})$ [29]. Second, the respective movement data was very minimal because during the freezing action, people had limited motion.

Labelling the normal and abnormal cases was done via a handcrafted labelling process. The Full details of data segmentation and pre-processing are in the following sections.

D. Data Pre-processing

1) Feature Selection from Row Data

Synchronized raw data from the wristband sensor (HR-monitor and movement data) is merged into 1 data file per subject per session, available as CSV-files (.csv). Each of the data-files contains 5 features/columns, and 70-time steps/rows, the columns contain the following data:

- 0 Timestamp (s)
- 1 Age
- 2 Gender
- 3 Weight
- 4 Movement
- 5 Heart Rate (bpm)

2) Data Segmentation

To segment the data, we used a sliding window of fixed length; the length of the window is 70s, with 1-second interval. The number of instances (segments) obtained after using this sliding window configuration is 1054 data samples (532 are normal and 522 are abnormal using handcrafted labeling).

The dataset of this study is a time series which is a sequence of real-valued data points with timestamps. We denote a time series as $T = \{t_1, t_2, \dots, t_n\}$, where t_i is the HR and motion values at time stamp i and there are n timestamps for each time series ($n=70$ seconds). In our study, our assumption is that all training series have the same number of timestamps. We denote a labeled time series dataset as $D = \{(T_j, y)\}_{j=1}^N$ which contains N time series data ($N=1054$) and their associated labels. For each $j=1 \dots N$, T_j represents the time series and its label is y . Our classification problem is binary

where y is a binary value $y = \{0, 1\}$ the label 0 is an indication that it is normal (no panic) and the label 1 is an abnormal situation.

It is important to emphasize that the sequence length ($n = 70s$) was not arbitrarily chosen; according to our experiments, we came up with the following justifications: (1) We ran our experiment with ($n = 70s, 80s, 90s, 100s$). We realized that the accuracy was better with $n=70s$ and $n=80s$, but we have chosen 70s because our system's objective is to achieve swift detection and it is better to make it as small as possible. On the other hand, we excluded any values that are less than 70s, because it may be computationally expensive.

3) Normalization

Artificial neural networks are one of the machine learning models that need data to be normalized, so we performed min/max normalization approach.

E. Using Deep Learning Sequence Classification Models for Mass Panic Detection

Deep Learning is a subfield of machine learning. It uses a layered structure of artificial neural networks (ANN). In recent years deep learning has attracted tremendous research attention and proven outstanding performance in many applications compared to conventional machine learning algorithms [30], [31].

Recurrent Neural Networks (RNNs), such as Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU), have achieved a significant outstanding performance over traditional machine learning models. RNNs have the ability to capture long-term temporal dependencies and variable-length observations and use contextual information when mapping between input and output sequences. They are able to deal with and overcome the vanishing gradient problem. (LSTM) was introduced by Hochreiter and Schmidhuber (1997) [32]. While (GRU) was recently proposed by Cho et al. [33], it is like LSTM but with fewer parameters.

According to our data collection and experiments, it was very difficult to find an empirical panic threshold that works with every human, due to the physiological fact that, every human being has different heart rate beating zones, as heart rate varies from person to person. There are many interior and exterior factors that can cause heart rate fluctuation (e.g. emotional or physical exertion will speed up the pulse). This depends on several factors such as age, gender, weight, activity, and fitness [34]. So we decided to use a machine learning classification approach.

Machine learning models have attracted both academia and industry entities and has proven outstanding performance in several applications. They are able to learn, adapt, and perform good classification. Furthermore, the detection itself is considered as a sequence classification and labeling problem. In order to differentiate whether there is a panic or not, there is a need to analyze the data as a time series (learning the relationship between the increment of HR and the respective accelerated movement and trace the situation for a specific time. To this end, it was necessary to handle this detection using a sequence deep learning model like RNN.

In this work, supervised deep learning sequence classification models are used, LSTM and GRU (Type: many to one). To run experiments on RNNs models, we used dropout with a value of 0.2. Number of epochs = 50 iterations. The internal architecture of LSTM and GRU is one layer consisting of 32 nodes and a time step of the sequences ($n=70$). Fig. 5. illustrates the used machine learning models.

F. Using Non-Sequential Machine Learning Models

For the sake of comparison, we also used the non-sequential machine learning models one-class SVM with default parameters, and Vanilla Neural Network (VNN). Our dataset is a time series being of equal length (time steps=70). The classification can be done using RNNs as well as VNN and SVM, so we compared the performance of sequential and non-sequential models. Note that the data representation in the non-sequential models is by inserting the data one after another so it is a one dimension vector of size 350 (5 features * 70 time steps).

G. Impelementaion Environmnet

We performed our experiments in Ubuntu environment using Keras [35] (on top of TensorFlow [36]) using Python (version 2.7). In this study, the models training and classification are performed on Intel (R) Core (TM) i7-3630QM CPU @2.40MGz and GPU NVIDIA Quadro K2000M With 8GB of shared GPU memory, and 16 GB RAM.

H. Results and Evaluation Metrics

To evaluate this study, we split the dataset ($N=1054$) into a training set (70%) which was used to train the model on k-fold cross validation $k=5$ and we left out 30% of the dataset for testing (unseen data). Accuracy metric represents the percentage of correct classification. We got accuracy of (97.48) % using LSTM, (95.58)% using GRU, (94.32)% using SVM and (94.01)% using VNN. Also, we calculated the False Positive Rate (FPR) which is the percentage of normal activities that are detected falsely as abnormal activities. Also we calculated the False Negative Rate FNR, as the percentage of abnormal activities that are classified falsely as normal activities (see Table I).

TABLE I. RESULTS OF TEST EXPERIMENTS (30% UNSEEN DATA)

Model		False Positive Rate (FPR)	False Negative Rate (FNR)	Accuracy (30% Unseen Data)
Sequential	LSTM	1.83 %	3.27 %	97.48 %
	GRU	6.71 %	1.96 %	95.58%
Non-Sequential	SVM	4.27%	7.19 %	94.32%
	VNN	9.76 %	1.96 %	94.01%

A confusion matrix is summarizing the performance of the classification model on a set of test data, it gives a better idea of what types of errors the classifier is making. It provides a good visualization of the performance of a classifier. Table II-V are summarizing the performance of LSTM, GRU, SVM, and VNN respectively.

VI. DISCUSSION

Regarding our experiments, we found that using sequence labeling methods such as LSTM and GRU can capture the temporal and spatial relationship between the HR and movement readings, compared to non-sequential models like SVM and VNN. Our results show that LSTM performed the best, followed by GRU, while SVM and VNN achieved less accuracy. Note that a small false positive rate (FPR) is acceptable because there is still data analysis phase II that can handle the false positives and differentiate between the real panic and normal HR and motion increment. Furthermore small false negative rate (FNR) is also acceptable because this false detection happens at an individual level, not for all attendees.

It is important to emphasize that our proposed system suits any mass gathering that does not involve extreme movement or excitement. To apply our system to such excited gatherings, data collection should be involved in a way that emulates the normal and abnormal activities for that particular event.

Furthermore, we discuss the limitations of our proposed system in the followings:

1) *Cost of IoT Wearable Wristband:* While the cost of wearable devices is declining significantly, still there is a significant cost in buying a huge number of sensors. (Although distributed over a large number of people as part of the event's cost). Perhaps when wristbands become standardized and become a common gear for everyone like cellular phones, this limitation will be lifted and people will participate in such mobile crowd sensing using their own wristbands.

2) *Deployment in Portable Devices with Limited Resources:* Power consumption and battery drain are very critical challenges because embedded wristbands have low computational/memory/energy resources. At the present, there are substantial challenges in deploying and performing inference of deep learning models in wearables. However, recently Google has announced TensorFlow Lite and TensorFlow Mobile [37]. They are two lightweight solutions for deploying AI on mobile and embedded devices. Furthermore, these lightweight solutions are powered by small hardware chips such as NVIDIA's Jetson TX2 and USB-based Intel's Movidius. Although these powerful chips still cost significantly. In the recent future, their prices will fall off and they will be used at large-scale.

3) *User Acceptance and Willingness to Contribute:* Our proposed system requires the obligatory or voluntary participation of thousands or even millions of individuals. First, individuals need to be persuaded or obliged to buy and/or carry the suitable devices. For example, in Hajj the Saudi Ministry of Hajj has already begun considering a plan to make electronic bracelets mandatory for all pilgrims that contain their personal data, according to ARAB NEWS [38].

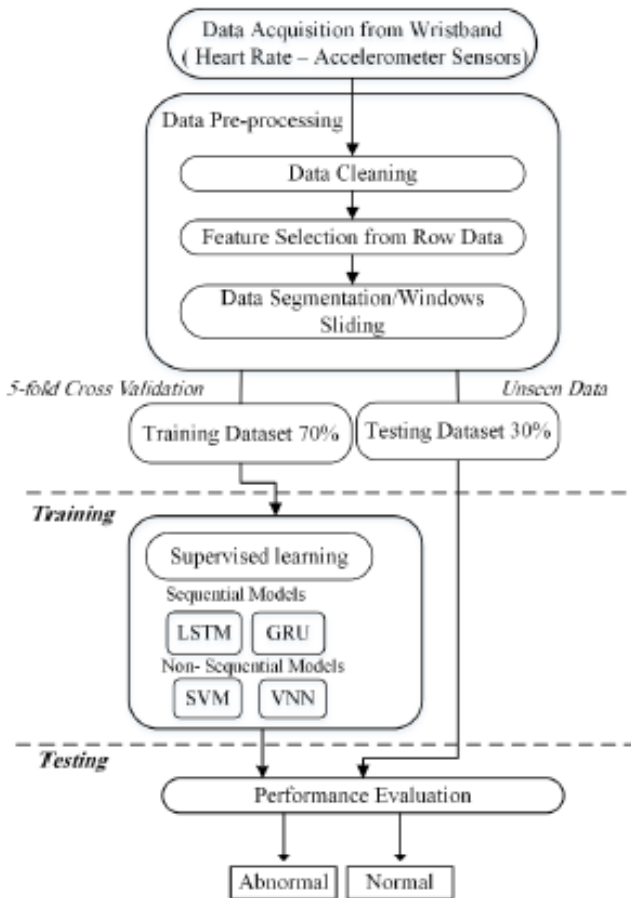


Fig. 5. Machine learning classification models.

TABLE. II. CONFUSION MATRIX OF TEST RESULTS (LSTM)

	Predicted	
	Abnormal	Normal
Actual		
Abnormal	148	5
Normal	3	161

TABLE. III. CONFUSION MATRIX OF TEST RESULTS (GRU)

	Predicted	
	Abnormal	Normal
Actual		
Abnormal	150	3
Normal	11	153

TABLE. IV. CONFUSION MATRIX OF TEST RESULTS (SVM)

	Predicted	
	Abnormal	Normal
Actual		
Abnormal	142	11
Normal	7	157

TABLE. V. CONFUSION MATRIX OF TEST RESULTS (VNN)

	Predicted	
	Abnormal	Normal
Actual		
Abnormal	150	3
Normal	16	148

VII. CONCLUSION AND FUTURE WORK

IoT and mobile crowd sensing applications help us to feel, see, and hear things that we never imagined before. This paper has proposed a new mobile crowdsourcing real-time detection method for abnormal crowd behaviour in mass gatherings. This system is based on advanced wireless and wearable sensors technologies and communication networks. It considers the humans themselves as the surveillance devices that exist everywhere. The proposed approach can detect mass panic swiftly in real-time manner by detecting the heart rate increase and abnormal motion.

In this study, we collected our own dataset from 89 subjects wearing an electronic wristband that has HR and motion sensors and we generated 1054 dataset samples. In this proposed approach, two-phases of data analysis were involved. Phase-I is a deep machine learning model that was used to analyze the sensors' collected readings of the wristband and detect if the person has indeed panicked so then sending alerting signals. While phase-II data analysis takes place in the monitoring server that receives alerting signals to conclude if it is a mass panic incident or a false positive case. The experimental results indicate that our developed deep learning sequential models LSTM and GRU got good accuracy of 97.48%, 95.58%, respectively. Compared to non-sequential models achieved only 94.32% for SVM and 94.01% for VNN.

In the future work, we will further develop data analysis phase-II module. Furthermore, using the same dataset, we will develop another detection algorithm that depends on anomaly detection approach; where we will train the RNN on normal data and predicting abnormality if there is a deviation of normality. Furthermore, we are planning to add another sensor, which measures skin conductance response (Electrodermal Activity EDA) that increases in sympathetic responses.

ACKNOWLEDGMENT

We would like to acknowledge participating subjects in the data collection phase for their willingness to contribute to this research and for the effort and time dedicated to data collection.

REFERENCES

- [1] L. Soomaroo and V. Murray, "Disasters at mass gatherings: lessons from history," *PLoS currents*, vol. 4, 2012.
- [2] B. Guo, Z. Yu, X. Zhou, and D. Zhang, "From participatory sensing to mobile crowd sensing," in *Pervasive Computing and Communications Workshops (PERCOM Workshops)*, 2014 IEEE International Conference on, 2014, pp. 593-598.
- [3] Egham. (2017, Accessed 29/11/ 2017). Gartner Says 8.4 Billion Connected "Things" Will Be in Use in 2017, Up 31 Percent From 2016. Available: <https://www.gartner.com/newsroom/id/3598917>
- [4] S. C. Mukhopadhyay, "Wearable sensors for human activity monitoring: A review," *IEEE sensors journal*, vol. 15, pp. 1321-1330, 2015.
- [5] A. Kumar, "Panic Detection in Human Crowds using Sparse Coding," 2012.
- [6] D. Helbing, I. Farkas, and T. Vicsek, "Simulating dynamical features of escape panic," *Nature*, vol. 407, pp. 487-490, 2000.
- [7] H. S. Bracha, "Freeze, flight, fight, fright, faint: Adaptationist perspectives on the acute stress response spectrum," *CNS spectrums*, vol. 9, pp. 679-685, 2004.
- [8] N. B. Schmidt, J. A. Richey, M. J. Zvolensky, and J. K. Maner, "Exploring human freeze responses to a threat stressor," *Journal of behavior therapy and experimental psychiatry*, vol. 39, pp. 292-304, 2008.
- [9] J. Leach, "Why people 'freeze' in an emergency: temporal and cognitive constraints on survival responses," *Aviation, space, and environmental medicine*, vol. 75, pp. 539-542, 2004.
- [10] A. R. Mawson, "Understanding mass panic and other collective responses to threat and disaster," *Psychiatry: Interpersonal and biological processes*, vol. 68, pp. 95-113, 2005.
- [11] A. Haag, S. Goronzy, P. Schaich, and J. Williams, "Emotion recognition using bio-sensors: First steps towards an automatic system," in *Tutorial and research workshop on affective dialogue systems*, Springer, 2004, pp. 36-48.
- [12] K. Kozłowska, P. Walker, L. McLean, and P. Carrive, "Fear and the defense cascade: clinical implications and management," *Harvard review of psychiatry*, vol. 23, p. 263, 2015.
- [13] T. Li, H. Chang, M. Wang, B. Ni, R. Hong, and S. Yan, "Crowded scene analysis: A survey," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 25, pp. 367-386, 2015.
- [14] B. Zhan, D. N. Monekosso, P. Remagnino, S. A. Velastin, and L.-Q. Xu, "Crowd analysis: a survey," *Machine Vision and Applications*, vol. 19, pp. 345-357, 2008.
- [15] M. Crocco, M. Cristani, A. Trucco, and V. Murino, "Audio surveillance: a systematic review," *ACM Computing Surveys (CSUR)*, vol. 48, p. 52, 2016.
- [16] R. Barnes and B. Rosen, "911 for the 21st century," *IEEE Spectrum*, vol. 51, pp. 58-64, 2014.
- [17] L. K. Moore, *Emergency communications: The future of 911*: Congressional Research Service, 2009.
- [18] T. Cheng and T. Wicks, "Event detection using Twitter: a spatio-temporal approach," *PloS one*, vol. 9, p. e97807, 2014.
- [19] T. Simon, A. Goldberg, and B. Adini, "Socializing in emergencies—A review of the use of social media in emergency situations," *International Journal of Information Management*, vol. 35, pp. 609-619, 2015.
- [20] A. A. Press, "CSIRO software monitors Twitter to guide emergency services to fires," in *theguardian*, ed, 2013.
- [21] S. O. Goldman, R. E. Krock, K. F. Rauscher, and J. P. Runyon, "Security system triggered by heart rate detection," ed: Google Patents, 2013.
- [22] D. W. McCleary, S. M. Rosato, J. O. Uchidiuno, W. Xiyang, and J. D. Weisz, "Detecting and notifying of various potential hazards," ed: Google Patents, 2018.
- [23] T. Wada, H. Higuchi, K. Komaki, H. Iwahashi, and K. Ohtsuki, "Disaster Detection Using SVDD Group Learning for Emergency Rescue Evacuation Support System," *Journal of Advanced Simulation in Science and Engineering*, vol. 3, pp. 79-96, 2016.
- [24] H. Iwahashi, H. Higuchi, K. Kogo, T. Kitamura, T. Wada, H. Okada, and K. Ohtsuki, "Extraordinary Judging Using Electrodermal Activity and Acceleration for Emergency Rescue Evacuation Support System," in *Parallel Processing Workshops (ICCPW)*, 2014 43rd International Conference on, 2014, pp. 355-360.
- [25] D. E. Long, *The Hajj Today: A Survey of the Contemporary Pilgrimage to Makkah*: SUNY Press, 1979.
- [26] D. Helbing, A. Johansson, and H. Z. Al-Abideen, "Dynamics of crowd disasters: An empirical study," *Physical review E*, vol. 75, p. 046109, 2007.
- [27] TOMTOM SPARK CARDIO WRISTBAND. Available: <https://www.tomtom.com>
- [28] T. Tamura, Y. Maeda, M. Sekine, and M. Yoshida, "Wearable photoplethysmographic sensors—past and present," *Electronics*, vol. 3, pp. 282-302, 2014.
- [29] H. Tanaka, K. D. Monahan, and D. R. Seals, "Age-predicted maximal heart rate revisited," *Journal of the American College of Cardiology*, vol. 37, pp. 153-156, 2001.
- [30] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks*, vol. 61, pp. 85-117, 2015.

- [31] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, p. 436, 2015.
- [32] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, pp. 1735-1780, 1997.
- [33] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [34] Heart.org. (July 2015, Accessed: 24/5/2017, 11:30 pm). American Heart Association, All About Heart Rate (Pulse) Available: <http://www.heart.org/>
- [35] F. Chollet. (2015, Accessed:13/6/2016). Keras. Available: <https://keras.io/>
- [36] Google. (2015, Accessed: 11/5/2016). Available: <https://www.tensorflow.org/>
- [37] TensorFlow. (2018, Accessed:15/2/2018). TensorFlow Lite versus TensorFlow Mobile. Available: <https://www.tensorflow.org/mobile/>
- [38] ARABNEWS. (2016, Accessed: 21/8/2016). E-bracelets a 'must' for pilgrims this Haj. Available: <http://www.arabnews.com/node/947221/saudi-arabia>

Motif Detection in Cellular Tumor p53 Antigen Protein Sequences by using Bioinformatics Big Data Analytical Techniques

Tariq Ali¹, Sana Yasin², Umar Draz³, Tayyaba Tariq⁵,
Sarah Javaid⁶
CS. Department
(CIIT) Sahiwal, Pakistan

M. Ayaz Arshad⁴
SCNC Research Centre
UoT, Tabuk,
Kingdom of Saudi Arabia

Abstract—Due to the rapid growth of data in the field of big data and bioinformatics, the analysis and management of the data is a very difficult task for the scientist and the researchers. Data exists in many formats like in the form of groups and clusters. The data that exist in the group form and have some repetition patterns called Motifs. A lot of tools and techniques are available in the literature to detect the motifs in different fields like neural networks, antigen/antibody protein, metabolic pathways, DNA/RNA sequences and Protein-Protein Interactions (PPI). In this paper, motif detection is done in tumor antigen protein, namely, cellular tumor antigen p53 (Guardian of the protein and genome) that regulate the cell cycle and suppress the tumor growth in the human body. As tumor is a death causing disease and creates a lot of other diseases in human beings like brain stroke, brain hemorrhage, etc. So there needs to investigate the relation of the tumor protein that prevents the human from not only brain tumor but also from a lot of other diseases that is created from it. To find out the gap between the motifs in the tumor antigen the GLAM2 is used that detects the distance between the motifs very efficiently. Same tumor antigen protein is evaluated at different tools like MEME, TOMTOM, Motif Finder and DREME to analyze the results critically. As tumor protein exists in multiple species, so comparison of homo tumor antigen protein is also done in different species to check the diversity level of this protein. Our purposed approach gives better results and less computational time than other approaches for different types of user characteristics.

Keywords—Bio-informatics; motif detection; guardian protein Tp53; DNA; tumor antigen; cancer; un-gapped motifs; MEME

I. INTRODUCTION

Big data became an active research from the last few years due to its immeasurable range of applications. Due to rapidly increasing trends and interest of research in this domain, there are many improvements have been done in this field that become famous among the research society due to manage the large amount of data that cannot be handled by the traditional

databases [1]. Instead of the momentous work on the motif discovery, motif detection in tumor proteins remains a difficult task for computer scientists and biologists. Lot of encouraging tools and algorithm is purposed in this field to make progress. Huge attempts have been done for the enlargement of the computational techniques for the identification of the sequence motifs in proteins. In the field of bioinformatics motif detection is an exigent problem due to the variety of protein motifs. In [2] author divides the biological motifs in three classes. Each class contains different type of motifs like the class 'A' contains the motifs that are in small size and appear at the functional sites of the biopolymer, cleavage and binding sites is the example of such type of motifs. The class 'B' contains large size motifs that are frequently crop up due to the divergent evolution and these motifs are highly associated with spherical structural domain. The recurring motifs fall in the class 'C' and these motifs are appearing due to the innovative recent replications. Due to diverts and complex nature of each class it's too much difficult to tackle all type of motif through single motif searching method (SMSM). There are multiple techniques to discover the over-represented motifs in the protein sequence to maximize the expectation in the sequence, but these techniques do not give the appropriate result. In this paper, graphical approach is used to detect the motifs in the tumor protein p53 that is the "*guardian of the proteins and genome*" because of its role in conserving stability by preventing genome mutation [3]. Among the field of protein-protein interactions the protein p53 has immune effect in the medical health sciences such as controls the oxidative stress, DNA damage; manage the functionality of ribosomal dysfunction and Hypoxia. In order to determine the rate of some metabolic and anabolic reactions when a cell protects from one step to the path for cancer and different diseases, therefore it is also called the tumor protein Tp53. The unique features and functionalities of the tumor protein p53 are shown in Fig. 1.

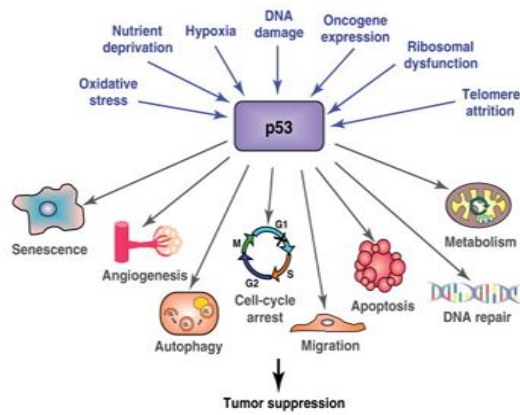


Fig. 1. Characteristics of Tumor protein p53.

The methods that were used in past for the motif detection was very slow and they search the motif of equal length. This motif does not work accurately in the class identification process in which motifs has different length and type. Motifs are located at innumerable distances in the protein sequence so to discover the gapped and un-gapped motifs is an essential task in the field of bioinformatics and big data. The importance of the gapped motifs is demonstrated also by the fact that many databases exist that contains motifs like PROSITE and ELM that contain gapped and different length motifs [4]. Newly purposed graph-based motifs detection technique efficiently searches the gapped and un-gapped motifs in the protein sequence. The tumor protein p53 data set that is selected for the motifs detection performs a momentous role in the body to control the cell cycle and apoptosis. The motifs that exist in the tumor protein p53 are minor persistent patterns that are accredited to have a conventional task. Defective p53 could be conceivably allowing the abnormal cells to promulgate that resulting in a tumor. As well as 55% of all human tumors comprise p53 mutants but in common cells of the human body, p53 protein level is small [5]. There are a lot of factors in the homo species that increase the p53 protein ratio like stress signal and DNA damage. There are multiple functions of p53: growth arrest, DNA repair and apoptosis (cell death) [6]. Tumor suppression becomes reduced in the human body if the p53 protein becomes damaged. A disease Li-Fraumeni syndrome occurs in the childhood if people inherit only one functional copy of the p53. The proposed research addresses the following issues:

- Motif detection in tumor protein
- Investigate the suitable parameters for the detection of motifs
- Reduce the tumor ratio in homo species by analyzing the proteins tumor that suppresses the tumor
- Identification of a suitable and match motif in the tumor protein
- Comparison of tumor protein motifs in varied species

In Fig. 2, the structure of different tumor proteins is represented with the root and leaf motifs hierarchy. The

alignment of different motifs of tumor protein at some different level is shown as a red label. This level further divides the motif roots where all the propagation of the motifs is present. In another way, the motif is basically the sub-part of motif root.

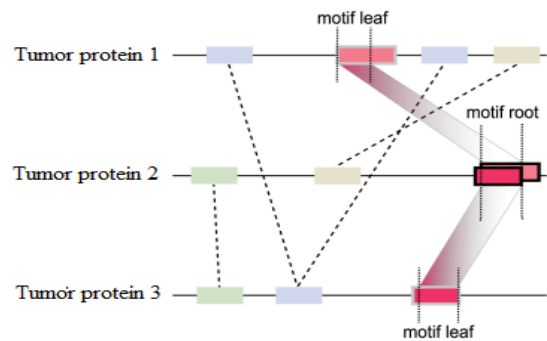


Fig. 2. Alignment of Motifs in the Tumor protein p53.

The major objective of this paper is to analyze the detection of tumor protein in the form of Motif Detection Algorithm (MDA) with the help of some suitable parameters. Furthermore, this research gives an idea of how it is possible to reduce the ratio of tumor proteins that suppress the tumor. After this, with the help of different tools, the two categories are assigning like residue motifs and site motifs that help to identify the matched motif in a different location. At the end, for the reliability of the results and its efficiency, the comparison is done between different species and identifies the number of matched motifs. The same protein dataset of p53 has been evaluated in different tools and Motif Finder technique, for the detection and evaluation of gapped motifs the GLAME2 is used and for the detection of un-gapped motifs, the MEME data bank is used. To perform the comparison analysis between the Tp53 the Motif Enrichment Analysis (MEA) is used that determines the position of best sites of the motifs against its possible probability of the detective motifs among the given data sets of the protein. Finally, after the comparison, the motif between different species the residue and site-based motifs are categorized.

Rest of this paper is organized as: section II is discussed the related work. Section III discussed the motif representation. Section IV describes the purposed methodology. Section V contains Motif Detection Algorithm. Section VI deal with the results and discussion. At the end, the conclusion is represented in Section VII.

II. RELATED WORK

A lot of research has been done at the motif detection in the field of big data and bioinformatics, but still more attention is required in this field. Recently, the research on efficient mining of previously unfamiliar, recurrently emerging patterns has received much attention in the field of medical health sciences [7]. With the advancement of technology and trend of social media; the amount of data is growing very rapidly. This data is spread across different places, in different formats, in large volumes ranging from Gigabytes to Terabytes, Petabytes, and even more. Today, the data is not only generated by humans, but enormous amounts

of data are being generated by machines and it surpasses human-generated data [8]. This size aspect of data is referred to as volume. In the Big Data world lot of work has been done but there still more work is pending to the amount of different data aspects like in Bioinformatics Big Data (BBD). Motif detection is some of the highly focused topics among the researchers in the big data research community. These motifs are useful for various time-series and data mining tasks. The relation between DNA and protein is a key motivating force. Binding of protein-sites and the specially targeted proteins are two important moves to understand the concept of biological activities [9]. A lot of techniques that gives high-throughput have recently purposed that try to enumerate the similarity between proteins motifs and protein [10]. In spite of the strong achievement, these techniques have some limitations and go down towards the strict classification of motifs. As a result, need further critical analysis of protein and protein sequences to dig out useful and modifiable information from a stack of strident and raw data. In [11] Motif Mark Algorithm (MMA) is purposed to find out the regular motifs in the protein sequence. This algorithm based on the graph theory and machine learning that finds binding sites of protein sequences. It also analyzes investigational data that is derived from universal protein against two of the most precise motif detection methods [12]. Gene mapping has been considered as one of the challenging tasks for researchers who belong to the field of bioinformatics and data analysis. Previously such tools are developed which are used for gene analysis and gene mapping for example MAPMAKER [10]. MAPMAKER has been applied to the construction of linkage maps in a number of organisms; including the human beings [13]. Mutations [14] in its structures can cause various diseases for this purpose simulation in proteins can reveal many new structures. One of the other techniques introduced for protein unfolding is Steering MD [14]. In [15] analysis is done on the whole genomes to find out the repetitive protein sequences called non-B motifs. These motifs are capable to predict the non-canonical structure of the protein and can autonomously report for deviation in mutation density. Graph visual motifs are also helpful for distinguish between applications protocol and to determine the known behavior of unlabeled traffic [16]. The most widely used tool for motif discovery is the MEME. MEME is a complete suite and performs a series of operations on the dataset thus discovering, analysis, finding enrichment and comparison with the existing motif databases [17]. Some other tools like DMINDA; Ensemble Genome Browser also performs a sequence of operations [18]. The combined effort of p53 and p63 in some Differential composition of DNA-binding sites may contribute to distinct functions of these protein homologs in some different species. To identify the legends and nuclei of the p54 proteins the SELEX (systematic evolution of legends by exponential enrichment) tool has been used [19]. To arrange the sequence in some protein input the long chain of the sequence has been used, for example; AGTGCGGCCGCTCAGGTTGACTTCCCCGCG.

In Western Bolt Analysis (WBA) [20]-[22] take the data sets of p53 and p63 for the proper cure of cancer in Health-Nutrient laboratory, they found that the p53 is most effective and the dominant parameters that play part and parcel role in the disease of Cancer. So, its need to investigate the p53 tumor

protein at different tools to find out the exact relationship with the different disease that plays our role in the sequence. In our proposed algorithm the p53 Motif Detection Algorithm (P53MDA) detects the gapped and un-gapped motifs.

III. MOTIF REPRESENTATION

Graphical based approach to find the gapped and un-gapped motifs in the sequence of the protein is in the form of regular expression is presented as:

$$R1 - p(n1; m1) - R2 - p(n2; m2) - \dots - Rr$$

$R1 =$ Base Class

$N1 =$ Least number of the base class

$M1 =$ Most number of the base class

Where 'R1' is an un-interrupted sequence with $l \leq i \leq r$ of amino acids that are called components, while $-p(n1; m1)$ represents a gap of length at least number 'n1' and at most 'l'. There are three major types of motifs depending on the size of the gap, the first one is contiguous motifs, these motifs have no gaps between them and 'n1' and 'l' values are zero 'n1' = 'm1' = 0 (for un-gapped motifs). The motifs that contain gaps called rigid motifs that fall into the second category of the motifs. The length of the rigid motifs is fixed i.e., 'n1 = m1' for all $l \leq i \leq r - l$. The third category of motifs is flexible gap motifs, these motifs contains different size gaps between the two motifs, i.e., 'n1 \leq m1; for all $l \leq i \leq r - l$.

IV. METHODOLOGY

The p53 tumor suppressor is implicated in cell cycle control, DNA repair, explicative sequence and programmed cell death. In-activation of the p53 contributes to the wide range of human tumors; including Glial neoplasm's. Due to its lot benefits and features its need to analyze the tumor protein more critically. In this paper, the proposed algorithm is introduced to detect the motifs of different lengths with gaps and without gaps through motif detection algorithm. The sequence analysis of the tumor protein p53 is performed by using MEME tool. To find the rigid motifs in the tumor protein sequence the GLAME2 tool is used that is very efficient for the detection of gapped motifs in the sequence. Tumor protein Sequence clusters are downloaded from UniProt database. The sequences of p53 are taken as a class of homo species. These sequences have their own significance in genomics. They were selected due to the unique feature of being the "guardian of the genome" and example of mutation caused in Homo sapiens. P53MD algorithm is used for finding the motif within the sequences. Discovered Motifs are shortest motifs that found out, than compared using TOMTOM and a resultant table is derived according to the number of matches along with the PROTEINS motifs are found using DMINDA tool box. The motif which is found in maximum number of sequences is compared with the existing database and results are derived. Simulation parameters for the proposed work are discussed in Table I. All the simulation is done with the help of proposed p53MDA and the alignment of this work is considering for both random and discreet. The reason is that for difference between the residue and site motif it is necessary to take the data set is evaluated for some random and discriminative fashion.

Up to our best knowledge this work is firstly done on the basis of both data types format. To select the discovery mode discriminative then the number of protein are aligned is some order, otherwise random order is apply. The novelty of this work is not only the detection of motif inside the protein but also provides the detail comparison among different available tools.

TABLE. I. GENERAL SIMULATION PARAMETERS OF P53 MDA AT MEME, DMINDA, DREME, TOMTOM & MOTIF FINDER

Simulation Parameters	Value
Discovery Mode	Random/Discriminative
Standard Custom Alphabets	Protein Sequence
Expected Motif Distributed	Zero/One Occurrence per sequence
No. of Motifs/Alignment	1000-5000
Strands Identification	ON
Presence of strands	+
Absence of strands	-
Length of Protein	5,000 for Each Protein/Sequence
Computational Timing	Depends on Input Sequence e.g 10 and 15 Seconds Each
Start/End Value	ON
Motif Score/Enrichment	ON
Best Motif Value	Mentioned
Location/Position of Motif	ON
X dimension of topography	1000
Y dimension of topography	4500

V. P53MDA FOR GAPPED AND UN-GAPPED MOTIFS

A. Pseudo code

Algorithm 1: p53 Motif Detection Algorithm

In this algorithm, Tumor protein sequence is taken as an input and novel gapped and un-gapped motif are detected that have fixed and variable length.

Input: ‘N’ numbers of Tumor protein sequence of TP53
Output: Gapped and Un-gapped motifs with variable length

1. Input tumor protein sequence
2. For s=1,.....,S-1 do
3. For i=1,.....,I-s do while {
4. For j=1,.....,J-s-i do while {
5. N1 = 0 // initialize the gap of length at least domain of the motifs variable (selected) from zero
6. M1 != 0 // initialize the gap of length at most number variable (unselected) from zero
7. for (s in 1: mid) // this loop runs from 1 to mid for selecting motifs from the first segment of input sequence
8. {
9. for (i in 1:2) // inner for loop is running from 1 to 2 times used for locations
10. {
11. for (j in 1:1/2:2) // inner for loop is running from 1 to half and half to 2 times used for locations
12. {

13. k=1
14. if (s [i, j] = =0 OR a [i, j=j+1] = =0)
15. {
16. gapped=gapped+1 // counts the gapped motifs from segment 1
17. gapped
18. }
19. else
20. {
21. Un-gapped =un-gapped+1 // counts the un-gapped motifs from segment 1
22. end
23. }
24. }
25. }

Fig. 3. High-level pseudo code for P53MDA.

VI. RESULT AND DISCUSSION

In this section, the graphical results have been displayed against the tumor proteins by using different tools with different parameters. Every homo species has nucleus, DNA, RNA, genes and lot of chromosomes. The genes that want to express motifs create their replication in the chromosomes that replication is called RNA and it generates proteins. In this paper, the tumor protein p53 is under consideration and it is analyzed critically. The tumor protein p53 dataset was first aligned and then analyzed on different tools of the MEME suite and Motif Finder. Every protein sequence contains a lot of residues. Fig. 3 shows all residues and their frequency.

A	Alanine	0.065	0.065
C	Cysteine	0.045	0.045
D	Aspartic acid	0.047	0.047
E	Glutamic acid	0.075	0.075
F	Phenylalanine	0.031	0.031
G	Glycine	0.059	0.059
H	Histidine	0.032	0.032
I	Isoleucine	0.029	0.030
K	Lysine	0.051	0.051
L	Leucine	0.089	0.089
M	Methionine	0.028	0.028
N	Asparagine	0.034	0.034
P	Proline	0.085	0.085
Q	Glutamine	0.037	0.038
R	Arginine	0.069	0.069
S	Serine	0.083	0.083
T	Threonine	0.051	0.051
V	Valine	0.056	0.056
W	Tryptophan	0.009	0.009
Y	Tyrosine	0.024	0.024

Fig. 4. Parameters for Motif detection.

B. Result through MEME: Motif Detection

Tumor protein sequence p53 is evaluated at MEME tool (Multiple Em for Motif Elicitation) that is a most famous tool for motif discovery and gained much attention in the field of bioinformatics. P53 based on probabilistic model and detects motifs by defining the probabilistic measures. p53 is the most important novel motif detection tool that takes the sequence as

an input and finds the innovative motifs that are repeated patterns in the sequence. MEME divides these variable length patterns into unique un-gapped motifs. MEME work on some parameters to finds the unique motifs with length 166 to 1134. Table II shows all those parameters that are used to detect motifs in MEME. Based on these parameters, the motifs are detected through MEME tool some of these un-gapped motifs are shown in Fig. 4, 5 and 6, respectively at different alignment. Through MEME tool box the first alignment is just to repeat the motifs in a given sequence. The second alignment finds the motifs in the whole sequence. In order to determine the reminder sequence from the given dataset, the third alignment is performed so that all the left and the right possibility of motifs have been determined. This alignment further helps to detect the best motif sites in the Motif Enrichment Analysis (MEA). Between the width of 6 and 50, the motifs, MEME is analyzed the total three basic motifs that frequently presents among the whole sequence. In this experiment, it has been noted that the detection of motifs is usually present in the middle of the p53 protein sequence.

TABLE. II. SIMULATION PARAMETERS OF MEME

Parameters	Description
Sequences	A set of 393 protein sequences of p53
Length	Between 166 and 1134
Background	Order-0 background generated by MEME according to given sequences
Distribution	One occurrence per sequence
Motif width	Between 6 wide and 50 wides
Site	EPLQVAHYREYWEYSIMCENKRTEQSVF EAYLIYVCMKHICTGEEELRVKES CSLQPSYSVFLFGYLDMCABQERMRTYI
Relative Entropy	31.5
Bayes Threshold	10.2515
p-value respectively	6.29e-12 9.71e-10 2.77e-9



Fig. 4 First Alignment of Motif detection: for the sake of replication Scenario in order to determine the Motif discovery among the whole sequence.



Fig. 5. Second alignment of Motif detection: for the sake of replication Scenario in order to determine the Motif discovery among the whole sequence.

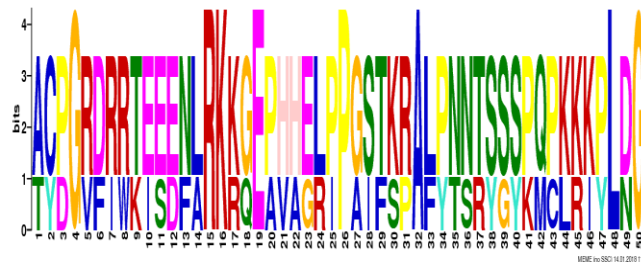


Fig. 6. Third alignment of Motif detection: for the sake of replication Scenario in order to determine the Motif discovery among the whole sequence.

VII. RESULT THROUGH MOTIF FINDER TOOL BOX: MOTIF DETECTION OF TUMOR PROTEIN

The tumor protein p53 is also analyzed by Motif Finder tool Box to check the multiple behaviors of this protein. The protein sequence is first aligned and then output motifs are detected. The resultant output of the Motif Finder is shown in Fig. 7.

Motif in the sequence

Query: UniRef100_UPI0009820F92
Pfam ID: Ank
Description: PF00023, Ankyrin repeat
Appearance:

Position	965..995
Alignment Query Database	GITALHNAVCAGHTEIVKFLVqFGVMVNAAD GnTPLHIAatngkkriik1LL.hGAdInald
i-Value	9.1e-05

Position	997..1024
Alignment Query Database	DGWTPLHCAASCNMVQCKFLVeSGAAV dGnTPLHIAatngkkriik1LL.hGAdI
i-Value	4.5e-06

Sequence:
MTGSVDNLASHLEEEKPEPNFEKNEILPSVNPFEYEEEDNIVYPDLLNYYNNLQNYVHP
NTNTNIQFQFTFGGGEDQDWYKSLDKLFKMEKIVPMRFHWEQILNPLGLYIRTKMV
YKVEQYRNEPVRRCNHMHAPTYHINQKLDPSVLPPYVHCVNHRGASVEVDNHLISLTQL
GSYEPGTQYEPMCFQFLCKNSCPSGMNRPELLEFTLEDSEGHVLGQKELSVRVCSPKR
DKLKEEKELKKTSEELIRQNSGNKMSSCDTHPYKVDISVAGKENFLSVLKYAHDVMAQQA
SRTGQYQAFKPYMDAIRKIP

Fig. 7. Motif detection through Motif finder.

A. Result through GLAME2: for the Sake of Gapped Motif Detection

To find out the gapped or rigid motifs in the tumor protein sequence p53 GLAME2 tool is used. The unique nature of this tool is that it finds the motifs of variable length. By GLAME2 the Best Motif is finding out in Fig. 8.

Best Motif Found:

NAME	START	SITES	END	STRAND	MARGINAL SCORE
sp P04637 P5	214	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	263	+	105.
sp P04637 -2	214	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	263	+	105.
sp P04637 -3	214	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	263	+	105.
sp P04637 -4	175	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	224	+	105.
sp P04637 -5	175	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	224	+	105.
sp P04637 -6	175	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	224	+	105.
sp P04637 -7	82	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	131	+	105.
sp P04637 -8	82	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	131	+	105.
sp P04637 -9	82	HSVVVYEPPEVGSDCITTIYNYVNCSSCGGNRRRIITITILEDSSGN	131	+	105.
sp Q12933 TR	395	Y.LN.G.DGTGRITLSL.FFVVKGN..DALLRWYFNQKVTLMWLQD	438	+	93.5
sp Q12933 -2	447	Y.LN.G.DGTGRITLSL.FFVVKGN..DALLRWYFNQKVTLMWLQD	490	+	93.5
sp Q12933 -3	394	Y.LN.G.DGTGRITLSL.FFVVKGN..DALLRWYFNQKVTLMWLQD	427	+	93.5
sp Q12933 -4	370	Y.LN.G.DGTGRITLSL.FFVVKGN..DALLRWYFNQKVTLMWLQD	413	+	93.5

Fig. 8. Best Motif detection through GLAME2.

B. Gapped Motif Detection at different Alignment

In order to determine the gapped motifs, the three different alignments have been performed in this Fig. 9. There is an inverse relationship between the alignment and its relative score. As the number of alignment is increased the relative score is decreased. The reason is that the motif enrichment is decreased when same data sets are evaluated again and again for the same data sets.

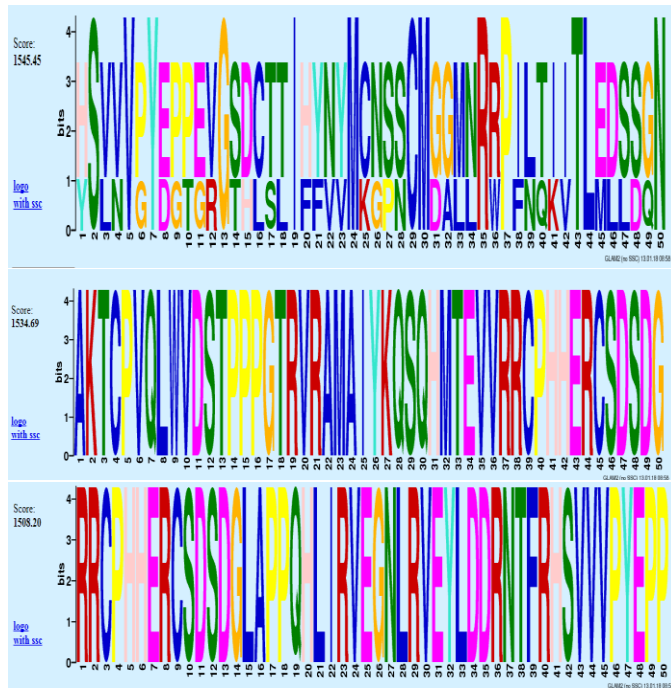


Fig. 9. Motif detection through GLAME2 with their relative score.

C. Result through DREME Tool Box: Un-gapped Motif Detection

In Fig. 10 and 11 un-gapped motifs are found by using DREME tool that discovers the short un-gapped motifs. The resultant output shows that there are only three motifs in the sequence that is un-gapped.

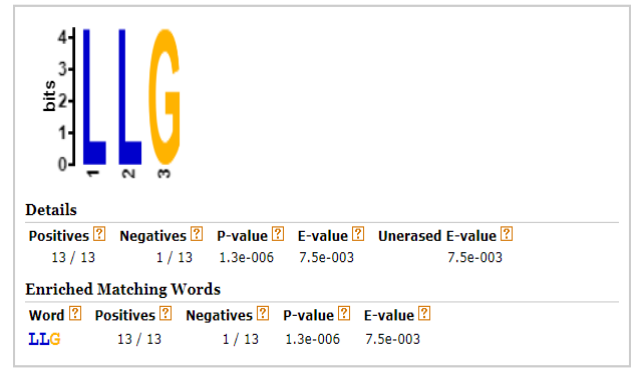


Fig. 10. Motif detection through DREME.

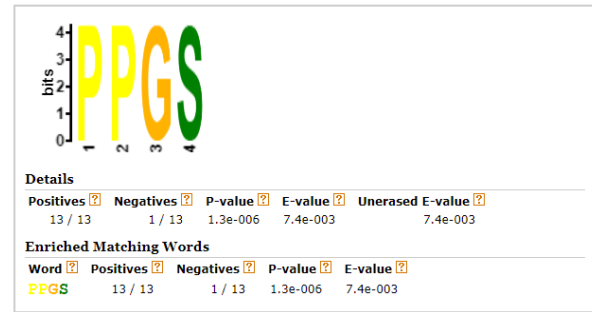


Fig. 11. Motif detection through DREME.

D. Motif Locations in the Tumor Protein Sequence

Fig. 12 shows the motifs location in the tumor protein sequence. To find out the motifs matching factor in the protein sequences, different color scheme is used that represent it clearly. The un-matched sequence appears in red color while the matched sequences across the tumor protein are displayed with the blue color. In order to highlight the location of best motif sequence, the green color is displayed.

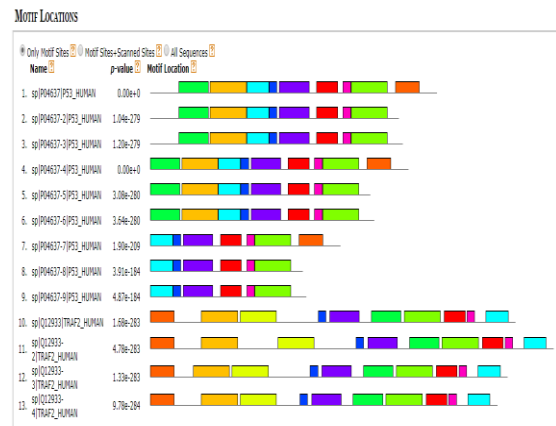


Fig. 12. Motif location in tumor protein at different colors.

E. Motif Enrichment Analysis in Tumor Protein p53

Motif Enrichment Analysis (MEA) of tumor protein is done by using the Centrimo tool that selects those motifs that have a maximum number of repetitions. It takes the previous motifs that were detected by the MEME, Motif finder, GLAME2, DREME as input and applies the enrichment on it.

The enrichment is done through some specific parameters like a number of motifs, motifs width and the match score of the motifs. Fig. 13 shows the Enriched Motif Graph on the basis of parameters that are described in Table III. The sequence is extracted for the probability value 1.0 against the position of best sites in the sequence. Three sub-datasets are used that are commonly evaluated among all the above tools. These three sub-datasets in the form of motifs is described as WTPFHCAASC, WWWWARLGD, and CHLAEVWCG.

TABLE III. MOTIF ENRICHMENT PARAMETERS FOR THE ENRICHMENT ANALYSIS OF MEME, DREME, GLAME2 AND MOTIF FINDER TOOL BOX

Parameters	Description
Motifs	set of 10 motifs
Width	10 and 15
Average width	15
Background	Order -0 backgrounds were generated
Match Score	Sites considered where match score ≥ 5
Region E-value	E -value ≤ 10

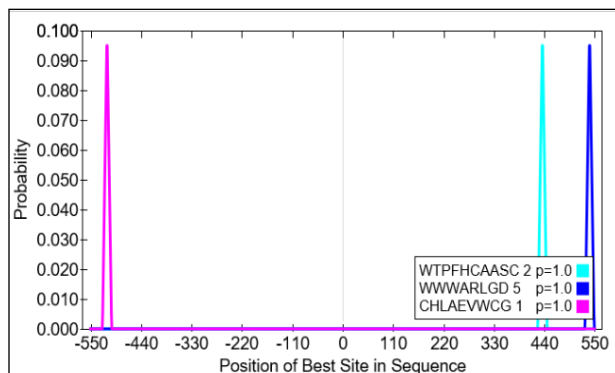


Fig. 13. Motifs Enrichment of Tumor Protein p53

F. Comparison of Tumor Protein p53 by TOMTOM TOOL BOX

The comparison is done of tumor protein p53 motifs with other protein motifs by using TOMTOM that is the best tool for motif comparison. The matching parameters for the TOMTOM comparison are given in Table IV. Fig. 14 shows the matched motifs of protein structure. Only one hit was found in the TOMTOM database that has been displayed below.

TABLE IV. MOTIF COMPARISON PARAMETERS OF MOTIF DETECTION AMONG OTHER SPECIES

Comparison Parameters	
p-value	2.04e-04
e-value	2.04e-04
q-value	2.04e-04
Overlap	7
Offset	0
Orientation	Normal

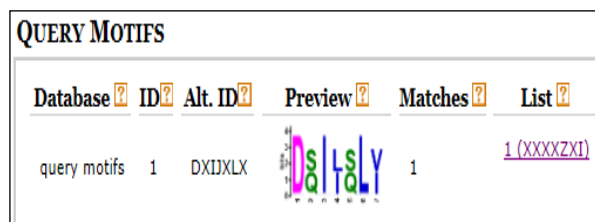


Fig. 14. Motifs comparison by TOMTOM

VIII. COMPARISON OF P53 WITH OTHER SPECIES

A. Comparison of Tumor Protein p53 Homo Species with Other Species with same Tumor Protein p53

Tumor protein p53 exists in multiple species. In Table IV the comparison of homo tumor protein p53 is done with other species that contain same tumor protein. The comparison is done on the previously detected motifs by the MEME, GLAME2, DREME and Motif Finder and number of matched motifs is found that is mentioned in Table V.

TABLE V. P53 MOTIF COMPARISON PARAMETERS WITH OTHER SPECIES

Knowledgebase	No of Matches
MOUSE	3
JASPER	4
Homo Sapiens	8
FLY (combined drosophila database)	2
CIS-BP Single species	1
Prokaryote DNA (CoLlecTF (bacterial	5
Ray 2013 all species (DNA Encoded)	3
YEAstract	3
Swiss Regulon e coli	9
DAP Motifs	0
Vertebrates	1
Malaria	9

Motifs exist in the homo species with different shapes formats and length at a different location. Due to the diversity of nature of motifs, the discovery of motifs is challenging task. Some motifs are site-based and some are residue-based. In Tables VI and VII, the motif detection is done by using multiple tools against the same tumor protein sequences. Firstly, the site based motif is detected by using MEME, GLAME2, HHMOTIF and SLIM Finder than residue based motifs is detected by using same tools and protein sequences. The results show that MEME tool is most appropriate to find the site and residue-based motifs in the tumor proteins

respectively. The two parameters are used for this propose like recall and precision. Recall states that the relevant motifs among the retrieved motifs and the precision states that the relevant motif that should be retrieved and the collection of precision and recall is F1-measure.

TABLE. VI. SITE BASED COMPARISON OF MOTIF

Site-based			
Tools	Recall	Precision	F1
MEME	0.236	0.564	0.333
GLAM2	0.249	0.099	0.142
HH-MOTIF	0.413	0.164	0.235
SLIM Finder	0.272	0.389	0.320

TABLE. VII. RESIDUE BASED COMPARISON OF MOTIF

Residue-based			
Tools	Recall	Precision	F1
MEME	0.210	0.420	0.280
GLAM2	0.219	0.061	0.095
HH-MOTIF	0.380	0.073	0.123
SLIM Finder	0.203	0.350	0.257

In Fig. 15 the accuracy comparison of different tools that are available for motifs detection is done and the performance of each tool is measured in the graphical format. Two parameters are under consideration, the first one is an F1 site that shows the site-based motifs in the sequence and the second one is the F1-residue that shows the residue-based motifs in the protein sequences. MEME tool is the only tool that finds the site-based and residue-based motifs with a high score and stood first in all of the other motif detection tools that are GLAM2, HHMOTIF, and SLIM Finder for the hundreds number of iterations.

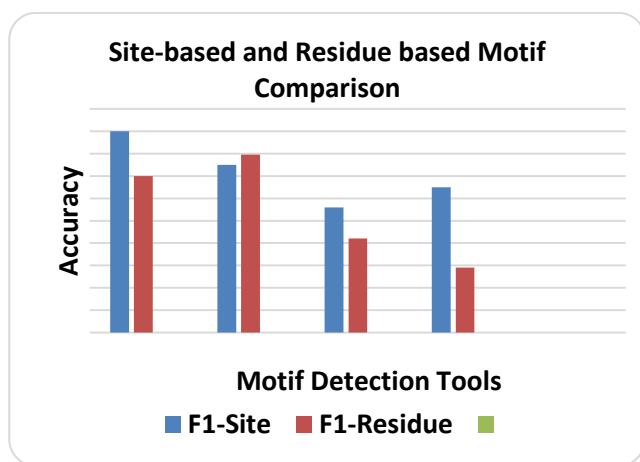


Fig. 15. Performance as measured By F1 of MEME and other tested de novo slim search methods.

IX. CONCLUSION

Big data became an active research from the last few years due to its immeasurable range of applications. A lot of research has been done at the motif detection in the field of big data and bioinformatics, but still more attention is required

in this field. Motif detection is some of the highly focused topics among the researchers in the big data research community. These Motifs are useful for various time-series and data mining tasks. The relation between DNA and protein is a key motivating force. Binding of protein-sites and the specially targeted proteins are two important moves to understand the concept of biological activities. The combined effort of p53 and p63 in some differential composition of DNA-binding sites may contribute to distinct functions of these proteins homologs in some different species. The p53 tumor suppressor is implicated in cell cycle control, DNA repair, explicative sequence and programmed cell death. Inactivation of the p53 contributes to the wide range of human tumors; including Glial neoplasm's. Due to its lot of benefits and features its need to analyze the tumor protein more critically. In this paper, the proposed algorithm is introduced to detect the motifs of different lengths with gaps and without gaps through p53 Motif Detection Algorithm. The sequence analysis of the tumor protein p53 is performed by using MEME tool. MEME tool is the only tool that finds the site-based and residue-based motifs with a high score and stood first in all of the other motif detection tools that are GLAM2, HHMOTIF, and SLIM Finder.

In this paper we have originated the problem of detecting motifs in the tumor proteins p53 that is depicted as "the guardian of the genome", referring to its role in persevering stability by preventing genome mutation and have offered a universal scheme for it. The P53MDA is purposed to detect the motifs in the tumor protein. Our formulation of the problem provides for a rigorous measure of the best fit between a given pattern and an example.

Up to our best knowledge, this work is firstly done on the basis of both data types format. To select the discovery mode discriminative then the number of protein are aligned is some order otherwise random order applies. The novelty of this work is not only the detection of motif inside the protein but also provides the detailed comparison among different available tools. The possible future direction is that to find out the best motifs and its correct alignment in a well-disciplined manner, in this way the diseases that are associated the DNA and Genome structure is easily trace out.

REFERENCES

- [1] Maass, W., Parsons, J., Puro, S., Rosales, A., Storey, V. C., & Woo, C. C. (2017). Big Data and Theory. Encyclopedia of Big Data, 1-5.
- [2] Rai, A., Pradhan, P., Nagraj, J., Lohitesh, K., Chowdhury, R., & Jalan, S. (2017). Understanding cancer complexome using networks, spectral graph theory and multilayer framework. Scientific reports, 7, 41676.
- [3] Saha, T. K., Katebi, A., Dhifli, W., & Al Hasan, M. (2017). Discovery of Functional Motifs from the Interface Region of Oligomeric Proteins using Frequent Subgraph Mining. IEEE/ACM transactions on computational biology and bioinformatics.
- [4] Lipper, C. H., Karmi, O., Sohn, Y. S., Darash-Yahana, M., Lammert, H., Song, L., ... & Jennings, P. A. (2018). Structure of the human monomeric NEET protein MiNT and its role in regulating iron and reactive oxygen species in cancer cells. Proceedings of the National Academy of Sciences, 115(2), 272-277.
- [5] Ma, W., Noble, W. S., & Bailey, T. L. (2014). Motif-based analysis of large nucleotide data sets using MEME-ChIP. Nature protocols, 9(6), 1428-1450.

- [6] Alipanahi, B., Delong, A., Weirauch, M. T., & Frey, B. J. (2015). Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning. *Nature biotechnology*, 33(8), 831-838.
- [7] Fan, C., Xiao, F., Madsen, H., & Wang, D. (2015). Temporal knowledge discovery in big BAS data for building energy management. *Energy and Buildings*, 109, 75-89.
- [8] Tran, N. T. L., & Huang, C. H. (2017). Cloud-based MOTIFSIM: Detecting Similarity in Large DNA Motif Data Sets. *Journal of Computational Biology*, 24(5), 450-459.
- [9] Schröter, M., Paulsen, O., & Bullmore, E. T. (2017). Micro-connectomics: probing the organization of neuronal networks at the cellular scale. *Nature Reviews Neuroscience*, 18(3), 131-146.
- [10] Yang, J., Jiang, B., Li, B., Tian, K., & Lv, Z. (2017). A fast image retrieval method designed for network big data. *IEEE Transactions on Industrial Informatics*.
- [11] Chen, D., Jiang, S., Ma, X., & Li, F. (2017). TFBSbank: a platform to dissect the big data of protein-DNA interaction in human and model species. *Nucleic acids research*, 45(D1), D151-D157.
- [12] Fan, C., Xiao, F., Madsen, H., & Wang, D. (2015). Temporal knowledge discovery in big BAS data for building energy management. *Energy and Buildings*, 109, 75-89.
- [13] Wang, Y., Wang, H., & Chang, S. (2018). A weighted higher-order network analysis of fine particulate matter (PM_{2.5}) transport in Yangtze River Delta. *Physica A: Statistical Mechanics and its Applications*.
- [14] Nagaraj, K., Sharvani, G. S., & Sridhar, A. (2018). Emerging trend of big data analytics in bioinformatics: a literature review. *International Journal of Bioinformatics Research and Applications*, 14(1-2), 144-205.
- [15] de França, F. O., Goya, D., & Penteado, C. C. (2018, January). Analysis of the Twitter Interactions during the Impeachment of Brazilian President. In *Proceedings of the 51st Hawaii International Conference on System Sciences*.
- [16] Hoh, J., Jin, S., Parrado, T., Edington, J., Levine, A. J., & Ott, J. (2002). The p53MH algorithm and its application in detecting p53-responsive genes. *Proceedings of the National Academy of Sciences*, 99(13), 8467-8472.
- [17] Wei, C. L., Wu, Q., Vega, V. B., Chiu, K. P., Ng, P., Zhang, T., ... & Liu, J. (2006). A global map of p53 transcription-factor binding sites in the human genome. *Cell*, 124(1), 207-219.
- [18] Li, W., Meyer, C. A., & Liu, X. S. (2005). A hidden Markov model for analyzing ChIP-chip experiments on genome tiling arrays and its application to p53 binding sequences. *Bioinformatics*, 21(suppl_1), i274-i282.
- [19] Huarte, M., Guttman, M., Feldser, D., Garber, M., Koziol, M. J., Kenzelmann-Broz, D., ... & Attardi, L. D. (2010). A large intergenic noncoding RNA induced by p53 mediates global gene repression in the p53 response. *Cell*, 142(3), 409-419.
- [20] Zhang, Y., Hu, Y., Wang, J. L., Yao, H., Wang, H., Liang, L. & Xu, J. (2017). Proteomic identification of ERP29 as a key chemoresistant factor activated by the aggregating p53 mutant Arg282Trp. *Oncogene*, 36(39), 5473.
- [21] Higgins, S. P., Tang, Y., Higgins, C. E., Mian, B., Zhang, W., Czekay, R. P., ... & Higgins, P. J. (2018). TGF- β 1/p53 signaling in renal fibrogenesis. *Cellular signalling*, 43, 1-10.
- [22] Chowdhury, K., Kumar, S., Sharma, T., Sharma, A., Bhagat, M., Kamai, A., ... & Mandal, C. C. (2018). Presence of a consensus DNA motif at nearby DNA sequence of the mutation susceptible CG nucleotides. *Gene*, 639, 85-95.

Investigating Methods of Resource Provisioning Mechanisms in Cloud: A Review

Babur Hayat Malik, Talia Anwar, Sadaf Ilyas, Farheen Jafar, Munazza Iftikhar, Maryam Malik, Noreen Islam Deen

Department of Computer Science
University of Lahore Gujrat
Campus, Gujrat

Abstract—Delivering information through cloud computing become a modern computation. For this purpose, electronic device is required to access with an active web server. For delivering different resources, the cloud supplier provides computing power for the cloud users to organize their multiple type of application at any time on different platforms. In cloud computing, the main drawback is relevant to the best use of resources as well as resource provisioning. In cloud computing there is a lack of desired resources that is why the cloud resource provision becomes a daring work. To maintain the quality of services, the provisioning of reasonable resources is need of workloads. The main problem is to find the appropriate workload that depends on the cloud user that is related to resource pair application requirements. This paper reveals the cloud resource provisioning and identification in general and in specific, respectively. In this paper, a methodical analysis of resource provisioning in cloud computing is presented, in which resource provisioning, different types of resource provisioning mechanisms and their comparisons, and benefits are described.

Keywords—Resource provisioning; resource provisioning mechanisms; cloud computing; systematic review; comparison between resource provisioning

I. INTRODUCTION

In this modern era, cloud computing plays vital role in all types of enterprises [1]. It is also considered as major and famous technology in all institutions and departments such as the different research community, governments, business and education. It also have established its top place in all the activities of the users that depends on the application of cloud computing and in their private lives as well. This medium is getting the most preferable position for presenting the applications that are related to data intensive. By using the strong ways of cloud computing like computation, flexibility, reliability and scalability, we can avoid from the problems and hurdles that creates by big data [2]. We define the cloud computing in the easiest way that provides a platform where the cloud provider store data that is accessible for the end user with the help of internet. The consequence of integrating the cloud computing that consists of multiple functionalities that are delivered to end user via internet by cloud providers, which is one of the cloud services. A software product that the end user use with the help of internet also connected with a web browser on the cloud environment is known as cloud application [3]. The basic technology of cloud computing is virtualization that implies the different operating systems to work on the similar physical type platform and also the

structures servers from the Virtual Machines. For delivering multiple resources like storage, CPU, memory, platforms and infrastructure the cloud services provider take help from the Virtual Machine [4].

In the beginning the cloud service providers start to deliver the many kinds of public cloud computing facilities. Cross breed cloud utilized by numerous endeavors and undertakings to build up their own foundation of distributed computing [5]. The main goal of the providers to achieve the strong profit, they deliver the best services and resources to their clients and accommodate with each other. The client in the cloud computing can get and also release different resources by demanding and recurring virtual machine. They have the promptly right for availing the proper and best quality of services not costly at all [6]. The allocation of resources is likely to more difficult as compared to other distributed systems like services of grid computing. The resources of arrangement or management of physical machines can be improved by using different virtual machine. Disk storage, bandwidth, memory and CPU are the multiple physical resources that are attached with virtual machine. Resource utilization is improved and multiplexed the resources in this way [5]. In the form of software application, infrastructure and platform the services are provided. The virtualization technology has become the current progress as computing standard. It represents dynamic provisioning on pay per use basis of computing services. Different pay per use stand services like software as services in information technology industry, platform as a services and infrastructure as services provides by the cloud computing [7]. Table I describes these cloud service models and their description.

TABLE I. CLOUD SERVICE MODELS [8]

Sr#	Cloud Service Models/Ref	Description
1	Software as a Service SaaS	Client's side Managed by a third party vendor It can be used for conventional cloud computing applications.
2	Platform as a Service PaaS	Development side Provides a grid or framework Best use is to develop or customize the applications.
3	Infrastructure as a Service IaaS	Provides a pool of resources of varied types Leased by the users according to their needs and requirements

Elasticity is the important component of cloud that composes it more and more attempting which also enhance the accessibility of multiple resources in the platform of cloud. Different resource provisioning methods utilized in cloud computing are more attractive if they increase the flexibility of the cloud to the utmost limit. Physical resources in the cloud computing can settle on this limit [9].

II. RELATED WORK

All the characteristics and the uses of the cloud resources are under the umbrella of resources arrangement and management [10]. The discovery process is the basic part of resource management. This process includes the finding for the reasonable resource that makes the way with the application requirements. The discovery process is organized from the cloud service provider. The physical computing is delivered by the infrastructure provider [11]. Service provisioning depends on service level agreements, these are the paired of non-functional matters and settlement that falls between the clients and service providers. It describes the term for the service that have the qualities of the service i.e. duties, prices and fines if there occurs any kind of violation in agreement. It is important for both the user and cloud services provider the reliable and flexible type of management of the service level agreements. In some cases the avoidance of service level agreements desecration evades the fines that are more costly to the providers. But in some other cases on the origin of timely and flexible responses to the suitable and possible service level agreements like desecration threat, minimum interaction with the system to make the cloud computing more reliable to attain the root and flexible type of on demand computing application. For the assurance of agreement of the service level agreements in the cloud computing, the cloud services provider should be able to observe its infrastructure and also the resource metrics to impel the desired services level purposes and objectives [12].

A. Cloud Workload Management

The set of examples which about to perform related to the cloud work that is a theoretical work. The effective and legal workload is relevant to running a web, which is one of the examples. In the perspective or framework of the cloud, the project of the workloads is offered. Various cloud workloads and the quality of services demands are not considered by them such as performance, price and time respectively. The users create or generate the workload requests that have been kept on the VI category which should be applied mechanism. Primarily, in the track of the purposes of the cloud users, the cloud users carefully determine that the assigned workload request can be executed. When the workload is admitted, the most important work or project is communicated with the preparations of the parts of the applications which implies the performance that once more lie in the arrangements' objectives of the cloud user. For creating the best and appropriate provisioning system, it is essential to sort or pinpoint the assortment of the cloud workloads and furthermore the nature of the administrations included too [7].

B. Need of Resource Provisioning

To increase the gratification and the chances or possibility of the users reaching the cloud, there is needed to increase the

large number of the requests or feedbacks that gratified from the cloud. Therefore, because of these perspectives the profit becomes so higher to the cloud in the consequence. There is possibility to appeal the customers of the cloud computing application to the cloud is to merge or short time for the responding. To make the attraction of the customers or users with cloud, the cloud is to need for accepting the resource provisioning technique which creates or generates the highest rate of the business deal. For the developing of the higher qualities of the business deal, there is not needed to be settled with the lack of period of the time. By giving the preference to the last, the trade-offs is to be sorted among the transaction success and U-turn time. Hence the shortage of the time can be created as far as it is conceivable for having the main goal to keep the high rate of the dealing success [9].

The applications can be used properly by applying the purpose of resource provisioning which implies that to discover the reasonable resources for the appropriate workloads in time. The best consequences can get by using the more effective resources. The reasonable and appropriate workload discovery is one of the main goals that maintain the program of different workloads. For making the quality of services more effective there is needed to satisfy the parts or units like utility, availability, reliability, time, security, price and CPU etc. So the resource provisioning reflects the performance of the time for the various workloads. All the presentations depend upon the kind or type of workload.

There are entirely two generic way of resource provisioning

- Static Resource Provisioning
- Dynamic Resource Provisioning

C. Static Resource Provisioning

For an application all types of desired resources are required in the peak time normally. Mostly this type of cloud provisioning the misuse of resources and wastage of resources because of workload is not considered in the peak time. Despite of this the resource provider offer the maximum desired resource for the purpose of avoids the service level application violation [13].

D. Dynamic Resource Provisioning

The customer demand, requirements and workloads are changed rapidly so that the cloud computing contain the elasticity element to the level of advanced automation adaption in the way of resource provisioning. This aim can be achieved through making the automatically scaling up and down of the resources that are assigned to a particular customer. This method is used to match the existing resources with the consumer current needs and demands with more good and reasonable way. In this way the element of elasticity is helpful to overcome the problem of under and over provisioning and also helpful in good and appropriate dynamic resource provisioning [14].

E. Parameters of Resource Provisioning

1) Response time: The algorithm of resource provisioning is designed to give response in minimum time after completing any task.

- 2) Minimize Cost: The cloud services cost should be less for the cloud consumer.
- 3) Revenue Maximization: The cloud services provider should be earned maximum revenue.
- 4) Fault tolerant: The algorithm provide services continuously in spite of collapse of nodes.
- 5) Reduced SLA Violation: The design of algorithm should be capable to decrease SLA violation.
- 6) Reduced Power Consumption: The placement & migration methods of virtual machine should be consume low power [13].

III. RESOURCE PROVISIONING

The resource provisioning term was commerce in the context of framework and grid computing. Due to lack of required and appropriate resources the cloud resource provisioning becomes a complicated task. In different distributed system the method of resource provisioning frequently contain the objectives and ways of share the workload on different resources and also enhance the amount of resource consumption and also minimize the workload execution time [7]. In cloud application the quality of services contains the provisioning of suitable resources for cloud workload. For the provision of appropriate resources that are used to workloads is a complex work and on the other hand based on quality of services in requirements and also the recognition of suitable resource pair workload in cloud is a hot research issue.

Fig. 1 explains the basic model of resource provisioning in cloud. Cloud user sends their workload like cloud application to the resource provisioning agents and establish good interaction with them. Resource provisioning agent (RPA) does resource provisioning and provide most suitable resource according to the customer requirements. When resource provisioning agent received the workload from user, his connection and access with the resource information centre (RIC) that have all the desired information about all type of resources with a resource pool. After that output can be achieved depend on the workload requirements as précised by consumer. Through resource discovery we know about the available resources and desired resources list can be generated. On the other hand the selection of resources is a procedure of choosing the most appropriate workload resource competition and match depended on the quality of services need expressed by the cloud user in tenure of services level application from the catalog and list which is created by the resource provisioning.

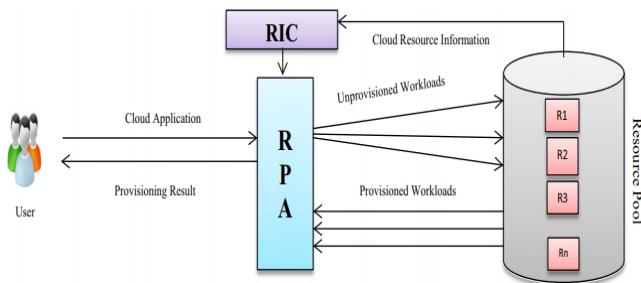


Fig. 1. Basic model of resource provisioning [10].

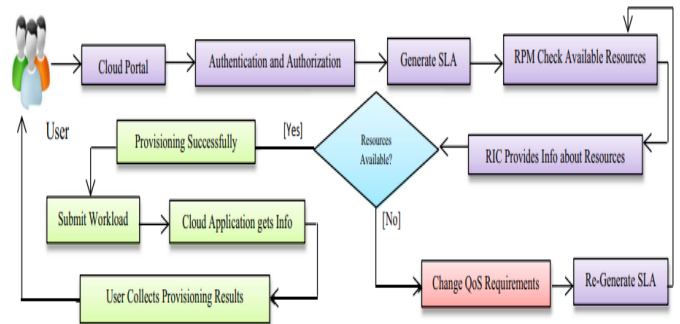


Fig. 2. Flow of resource provisioning [10].

Fig. 2 explains the procedure of cloud resource provisioning. Cloud consumer can interact with the help of cloud portal and also presents the quality of services (QoS) requirements of workload after the complete authentication procedure. Resource information centre (RIC) delivered the information that is based on customer requirements and available resources are checked by the resource provisioning agent (RPA). It is also helpful in provision of desired resources to the cloud application workload for running or execution in the environment of cloud but after fulfill the condition that is demanded resources are present in the resource pool. Resource provisioning agent (RPA) applications for again submit the workload with another quality of services requirement like a service level application article or document in the condition of if the desired resource is unavailable as the requirements of quality of services. Resource scheduler submitted the workloads after the appropriate provisioning of available resources. Resource scheduler requests to submit and present the whole workloads for all the available provisioned resources. Then resource provisioning agent again received the results and send these provisioning outputs and results to the cloud consumer.

IV. RESOURCE PROVISIONING MECHANISEMS

Some of the extensively and widely used cloud resource discovery and resource provisioning methods or mechanisms are based on the dynamic or distributed resource provisioning. Table II describes some resource provisioning mechanisms.

TABLE II. RESOURCE PROVISIONING MECHANISM

Sr#	RPM	Description
1	QoS Based RPM	The major objective of such work is to provide provision on different resources before managing in an appropriate manner or way and then execute this application for getting optimal results to the end user.
2	Cost Based RPM	Minimize the total amount of resource provisioning cost like over provisioning cost and under provisioning cost. Cost reduction can assure the double capacity of application
3	SLA Based RPM	SLA provisioning method depend on the admission control that can be maximizes the revenue and also the utilization of resources resource utilization and also pay attention on multiple type needs of SLA that are consumer described.
4	Time Based RPM	Minimum execution time can double the application capacity as well as minimize the overhead cost of switching servers.

5	Energy Based RPM	Enhance the resource utilization and must be reduce the consumption of power.
6	Dynamic Based RPM	Decision related to different and changing environment such as electricity bills and user requirements. Fully and partially distributing the clouds computing services facilities with other consumers.
7	Adaptive Based RPM	Methodologies which can be based on the virtualization for the purpose of resource provisioning depend on the need of application dynamically as well as minimize the consumption of power and energy by maximizing the server's usage.
8	Optimization Based RPM	The running cost of consumer application can be reduced through advancing the energy resources and also meet the required deadline on time make sure that SLA objectives cannot be violated.

A. QoS based Resource Provisioning

The achievement of cloud administrations depends immensely on the level of fulfillment of cloud clients as far as execution and nature of-benefit (QoS) they get from cloud specialist co-ops. QoS alludes to an arrangement of characteristics or qualities of an administration, for example, accessibility, security, reaction time, throughput, inactivity, unwavering quality, and notoriety. Asset provisioning research work in view of QoS has been finished by following creators.

Xiaoyong Xu et al. [15] propose an occasion driven asset provisioning structure. This system recognizes all occasions that conceivably cause any Map Reduce calculation hard due date absent and pointless asset (Virtual Machine) squander, and instantly handles those occasions. Along these lines, this structure can ensure that the due dates of those Map Reduce calculations running in system are met while limiting the running cost of the structure.

Bahman Javadi et al. [16] considered the issue of QoS-based asset provisioning in a cross breed Cloud figuring framework where the private Cloud is disappointment inclined and built up an adaptable and versatile half breed Cloud design to take care of the issue of asset provisioning for clients' solicitations. The proposed engineering uses the Inter Grid ideas which depend on the virtualization innovation and embrace an entryway (IGG) to interconnect distinctive asset suppliers. The creator proposed facilitating techniques in the half breed Cloud framework where an association that works its private Cloud plans to enhance the QoS for the clients' solicitations by using the general population Cloud assets.

Sukhpal Singh Gill et al. [17] show a keen QoS-mindful autonomic asset administration approach named as CHOPPER (Configuring, Healing, Optimizing and Protecting Policy for Efficient Resource administration). CHOPPER offers self-setup of uses and assets, self-mending by taking care of sudden disappointments, self-protection against security assaults and self-advancement for greatest asset utilization. Author assessed the execution of the proposed approach in a genuine cloud condition and the trial comes about show that the proposed approach performs better as far as cost, execution time, SLA infringement, asset dispute and furthermore gives security against assaults.

B. Cost based Resource Provisioning

Cost provisioning research work has been finished by following creators.

Aarti Singh et al. [7] proposes another Agent based Automated Service Composition (A2SC) calculation containing demand handling and computerized benefit organization stages and isn't in charge of seeking thorough administrations yet in addition considers lessening the cost of virtual machines which are devoured by on-request benefits as it were.

Adel Nadjaran Toosi et al. [18] propose another asset provisioning calculation to help the due date prerequisites of information escalated applications in crossover cloud situations. To assess this proposed calculation, creator actualizes it in Aneka, a stage for creating adaptable applications on the Cloud. Trial comes about utilizing a genuine contextual analysis executing an information concentrated application to quantify the walk capacity list on a half and half cloud stage comprising of dynamic assets from the Microsoft Azure cloud demonstrate that the proposed provisioning calculation can all the more productively dispense assets contrasted with existing techniques.

Smita Vijayakumar et al. [19] consider versatile spilling applications where a client needs to accomplish the base asset costs while keeping up a predetermined exactness objective. Creator shows a dynamic and robotized system which can adjust the versatile parameters to meet the particular exactness objective, and afterward powerfully focalize to close ideal asset distribution. This arrangement can deal with surprising changes in the information appropriation qualities and additionally rates. Creator assesses our approach utilizing two gushing applications and exhibits the adequacy of our structure.

Safiye Ghasemi et al. [2] proposed a novel learning based asset provisioning approach that accomplishes cost-decrease certifications of requests. The commitments of this upgraded asset provisioning (ORP) approach are as per the following. Right off the bat, it is intended to give a financially savvy strategy to proficiently deal with the provisioning of asked for applications. ORP performs in light of administrations of which applications included and thinks about their proficient provisioning completely. Furthermore, it is a learning automata-based approach which chooses the most appropriate assets for facilitating each administration of the requested application. Thirdly, a far reaching assessment is performed for three regular workloads: information escalated, process-concentrated and ordinary applications.

C. SLA based Resource Provisioning

Distributed computing depends on getting to each sort of asset through an "as-a-benefit" interface, and on the reception of a compensation for every utilization plan of action. In such a specific situation, Security Service Level Agreements (SLAs) expect a key part, as they permit, in addition to other things, to announce plainly the security level conceded by suppliers to clients, and also the imperatives postured to the two gatherings (suppliers and customers).

Yoori Oh et al. [20] propose an auto-scaling system with relating calculations to oversee assets powerfully in virtual

conditions, so as to meet client determined SLA (Service Level Agreement) given an arrangement of restricted assets. In this paper, he proposes an auto-scaling strategy for using asset of Spark groups successfully in distributed computing condition. The proposed auto-scaling technique has an objective to meet client indicated due date. Likewise perform tests to check the adequacy of the proposed scaling calculation.

Elarbi Badidi et al. [21] propose a system for SaaS provisioning, which depends on expedited Service Level assertions (SLAs), between benefit buyers and SaaS suppliers. A Cloud Service Broker (CSB) helps shoppers choosing the privilege SaaS supplier that can satisfy their useful and nature of-benefit (QoS) prerequisites. Besides, the CSB is responsible for arranging the SLA expressions utilizing a multi-characteristics transaction display with a chose SaaS supplier in the interest of the administration purchaser, and checking the consistence to the SLA amid its usage.

Obinna Anya et al. [22] present an approach for versatile administration provisioning in the Cloud in view of QoS examination. A noteworthy commitment of the approach is the improvement of an examination motor for prescient flexibility administration of Cloud benefit provisioning that incorporates top to bottom mining of SLA consistence history with learning of business setting, e.g. workload fluctuation, a client's business objectives, application execution, and administration operational setting. In this work-in-advance report, creator portrays the proposed system and talks about conceivable usage and arrangement situations.

Valentina Casola et al. [23] introduce the SPECS system, which empowers the improvement of secure cloud applications secured by a Security SLA. The SPECS structure offers APIs to deal with the entire Security SLA life cycle and gives every one of the functionalities expected to automatism the implementation of appropriate security systems and to screen client characterized security highlights. The improvement procedure of SPECS applications offering security-upgraded administrations is represented, exhibiting as a certifiable contextual investigation the provisioning of a safe web server.

D. Time based Resource Provisioning

Time provisioning research work has been finished by following creators.

Lakshmi Ramachandran et al. [24] propose a novel User Interface-Tenant Selector-Customizer (UTC) model and approach, which empowers cloud-based administrations to be methodically demonstrated and provisioned as variations of existing administration occupants in the cloud. This approach thinks about utilitarian, non-practical and asset allotment prerequisites, which are unequivocally indicated by the customer through the UI segment of the model. This is the primary such coordinated approach that represents the thoughts utilizing a practical running case, and furthermore exhibit a proof-of-idea model manufactured utilizing IBM's Rational Software Architect displaying device.

Izzet F. Senturk et al. [25] propose BioCloud as a solitary purpose of passage to a multi-cloud condition for non-PC adroit bio analysts. They talk about the design and segments of BioCloud and present the planning calculation utilized in

BioCloud. Trials with various utilize cases and situations uncover that BioCloud can diminish the worked execution time for a given spending plan while typifying the multifaceted nature of asset administration in numerous cloud suppliers.

E. Energy based Resource Provisioning

Regardless of the achievement of some outstanding CSPs, for example, Google App Engine (GAE) and Amazon Elastic Compute Cloud (EC2), the huge vitality costs as far as power devoured by server farms is a genuine test. Vitality use of server farms has two critical highlights: (i) servers have a tendency to be more vitality wasteful under low usage rate, and (ii) servers may expend a lot of energy out of gear mode.

YaGao et al. [26] create measurable nature of administration (QoS) driven power control approaches to expand the successful vitality proficiency (EEE), which is characterized as the range effectiveness under given determined QoS imperatives per unit reaped vitality, for vitality gathering based remote systems. Specifically, to begin with, break down the long haul accessible vitality imperatives and detail the EEE amplification issue. At that point, determine the shut frame arrangements of ideal power control strategies to the EEE augmentation issue under the battery limit overwhelmed imperative, the normal collected vitality requirement, and both the battery limit and normal gathered vitality limitations, individually.

Mingxi Cheng et al. [4] presents DRL-Cloud, a novel Deep Reinforcement Learning (DRL)- based RP and TS framework, to limit vitality cost for vast scale CSPs with substantial number of servers that get colossal quantities of client demands every day. A profound Q-learning-based two-organize RP-TS processor is intended to naturally produce the best long haul choices by gaining from the changing condition, for example, client ask for designs and reasonable electric cost. With preparing strategies, for example, target arrange, encounter replay, and investigation and misuse, the proposed DRL-Cloud accomplishes astoundingly high vitality cost productivity; low reject rate and in addition low runtime with quick merging. Contrasted and one of the best in class vitality effective calculations, the proposed DRL-Cloud accomplishes up to 320% vitality cost proficiency change while keeping up bring down reject rate all things considered.

F. Dynamic based Resource Provisioning

The one element of cloud that makes it speaking to its clients is Elasticity which expands the accessibility of assets in the cloud. The asset provisioning strategy utilized as a part of a cloud is said to be sound in the event that it improves the cloud's versatility to as far as possible. This point of confinement is controlled by the measure of physical assets in the cloud.

Kirthica S. et al. [9] give a proposition. That proposition means to supplant the existing asset provisioning strategy in an open system, Cloud Inter-task Toolkit (CIT), to make an expanded exchange progress rate which is apparent from the experimental comes about acquired from an ongoing heterogeneous cloud condition set up utilizing Eucalyptus, OpenNebula and OpenStack. Notwithstanding fulfilling a demand with assets from numerous mists, these assets are

totalled and given in a way that is effectively accessible by the client.

G. Adaptive based Resource Provisioning

It is trying for cloud suppliers to assign the pooled processing assets progressively among the separated clients as to amplify their income. It isn't a simple errand to change the client arranged administration measurements in to working level measurements, and control the cloud assets adaptively in light of Service Level Agreement (SLA) [27].

Guofu Feng et al. [27] addresses the issue of amplifying the supplier's income through SLA-based dynamic asset designation as SLA assumes an essential part in distributed computing to connect specialist organizations and clients. Creator formalizes the asset portion issue thinking about different Quality of Service (QoS) parameters.

Ayoub Alsarhan et al. [28] propose a novel Service Level Agreement (SLA) structure for distributed computing, in which a value control parameter is utilized to meet QoS requests for all classes in the market. The structure utilizes fortification learning (RL) to determine a VM procuring arrangement that can adjust to changes in the framework to ensure the QoS for all customer classes. These progressions include: benefit cost, framework limit, and the interest for benefit. This approach incorporates processing assets adjustment with benefit confirmation control in light of the RL show.

Abiola Adegboyega et al. [29] built up a model to anticipate transfer speed usage pertinent in keeping up SLAs for various activity streams at the cloud organize edge and center. The created univariate estimate show utilizes the Auto-Regressive Integrated Moving Average (ARIMA) demonstrate expanded with a general class of Adaptive Conditional Score Models (ACS). Creator inspiration for utilizing the ACS comes from its powerful adjustment to exceptions and drifters more proficiently with expanded computational exactness than current strategies; one of such techniques being the as of late embraced Generalized Auto-Regressive Conditional Heteroskedasticity (GARCH) to show instability.

H. Optimization based Resource Provisioning

Enhancing provisioning research work has been done by following makers.

YaGao et al. [26] create factual nature of administration (QoS) driven power control strategies to augment the compelling vitality effectiveness (EEE), which is characterized as the range proficiency under given indicated QoS limitations per unit collected vitality, for vitality gathering based remote systems. Specifically, in the first place, creator breaks down the long haul accessible vitality imperatives and defines the EEE augmentation issue. At that point, infer the shut frame arrangements of ideal power control strategies to the EEE boost issue under the battery limit overwhelmed imperative, the normal collected vitality limitation, and both the battery limit and normal gathered vitality requirements, individually.

Marcus Lemos et al. [30] ACOSIM, an approach ACOSIM, to limit the general sensor cloud vitality utilization by choosing just a subset of sensor hubs to create the virtual sensors. Results from starting investigations demonstrate that the

approach decreases the sensor cloud vitality utilization the by 73.97%, giving an answer for be considered in sensor cloud situations.

V. COMPARISON OF RESOURCE PROVISIONING MECHANISEMS

Examination of asset provisioning systems is a troublesome assignment because of various sorts of asset provisioning components and the absence of benchmarks. We considered distinctive characteristics of asset provisioning components and look at them.

A. Traits of Resource Provisioning

RPM in cloud frameworks can be looked at in view of some normal qualities for taking care of provisioning issues. Searching mechanism, objective function, resource provisioning strategy, merits and demerits are a portion of the normal and essential qualities that ought to be inspected in every RPM as depicted in Table III. Table IV demonstrates the difference of asset provisioning components in view of these qualities.

TABLE III. RESOURCE PROVISIONING TRAITS

Sr#	Traits	Description
1	Provisioning Mechanism	Strategy for give assets
2	Searching Mechanism	Finding the best workloads and assets relies upon seeking speed..
3	Objective Function	A target capacity of each RPM is particularly intended for a particular reason for the system..
4	Resource Provisioning Strategy	The methodology of giving assets to workloads execution is called Resource Provisioning Strategy (RPS) scientific classification. Two kinds of RPS scientific classification are depicted beneath: • Dynamic • Distributed
5	Merits	The benefits of Resource Provisioning Mechanism are depicted in this segment.
6	Demerits	The disservices of Resource Provisioning Mechanism are depicted in this segment.

B. Resource Provisioning Comparison

TABLE IV. RESOURCE PROVISIONING COMPARISON

Sr #	Provisioning Mechanism/Ref	Topic	Searching Mechanism	Objective function	RPS	Merits	Demerits
1	Cost Based RPM/[30][19]	Competent resource provisioning and distribution techniques for cloud	KFCM algorithm was used to cluster the available resource. Particle swarm optimization algorithm is used to select the	To distribute the resources in a powerful way	Distributive	This system accomplishes least execution time, least cost esteem and low	Calculation Difficult

		computing environment	optimal resource with minimum cost			memory.													provisioning for cloud-based MapReduce in dynamical environments	framework. Scaling Up Algorithm (SUA) is applied to handle the computation falling behind and latest intervention events.	base asset cost.			oning structure not just certifications the QoS of those MapReduce calculations, yet in addition decreases the running expense of MapReduce calculations	factor which influence the execution time of MapReduce Calculations.
		Automated and Dynamic Application Accuracy Management and Resource Provisioning in a Cloud Environment	Framework that dynamically achieves the user-specified accuracy level by adapting an adaptive parameter at runtime. Main Processing Loop Algorithm.	The objective is to precisely apportion assets in order to coordinate the rate of information entry, for the picked estimation of versatile parameters	Adaptive	The structure is viable. The CPU designations perfect esteem and the overhead of the general structure is very little.	Structure is touchy to the adjustments in input information attributes as well as information landing rates, which could expect changes to the versatile parameters further more, the CPU assignment, separately												Failure-aware resource provisioning for hybrid Cloud infrastructure	A failure-aware resource provisioning algorithm that is capable of Attending to the end-users quality of service (QoS) requirements. developed a flexible and scalable hybrid Cloud architecture	To understand the full capability of the half breed Cloud stage, a building structure for proficiently coupling open and private Clouds is vital..	Dynamic	Ready to enhance the clients' QoS about 32% as far as due date infringement rate and 57% as far as log jam with a constrained cost on an open Cloud.	Try not to consider the tradeoff among st cost and performance in instance of asset disappointments on the nearby group.	
2	QoS Based RPM / [17][15][16][31][32]	CHOPPER: an intelligent QoS-aware autonomous resource management approach for cloud computing	QoS-aware autonomous resource management approach named as CHOPPER	CHOPPER offers self-design of uses and assets, self-recuperating by taking care of sudden disappointments, self-protection against security assaults and self-enhancement for most extreme asset use	Distributed	makes strides security, vitality productivity, dependability and accessibility of cloud based administrations in genuine cloud stages	Not give versatility												Self-Tuning Service Provisioning for Decentralized Cloud Applications	An auction-based resource allocation approach. Vickrey Auction: A well researched Mechanism .	To give calculation assets to clients which are significantly nearer to them, potentially inside switches.	Dynamic	The proposed instrument is very adaptable, effective and approved by broad Reenactments.	It is unpredictable and resistant to false-name assaults.	
		QoS-guaranteed resource	Event-driven resource provisioning	To ensure Quality of Service (QoS) with the	Dynamic	Occasion driven asset provision	Systems administration												QoS based resource provisioning and schedule	Resource provisioning framework	Formal detail and confirmation of the structure helps in foreseeing conceivab	Static	Outlined and tried for asset provisioning and	Adaptability as a metric isn't considered in this	

		ing in grids		le blunders previously the planning procedure itself, and subsequently brings about proficient asset provisioning		bookin g challenges	structu re			ks						
3	Dyna mic Base d RPM /[9][14]	Horizo ntal scaling and aggregation across heterog eneous clouds for resourc e provisi oning	Open framework Cloud Inter-operation Toolkit (CIT) use	Supplant the existing asset provisioning strategy in an open system to make an expanded exchange progress rate which is apparent from the experimen tal comes about acquired from a constant heterog eneous cloud condition	Inde pendent	Fulfilling a demand with assets from numero us mists, these assets are totaled and given in a way that is effortle ssly accessi ble by the client.	Just for a solitar y VM, and not for a total of VMs.			DRL- Cloud: Deep Reinforcement Learnin g- Based Resourc e Provisi oning and Task Schedu ling for Cloud Service Provide rs	DRL- Cloud, a novel Deep Reinforcem ent Learning (DRL)- based RP and TS system	A profound Q- learning- based two- arrange RP-TS processor is intended to naturally produce the best long haul choices	Dyn amic	limit vitality cost for huge scale CSPs	Expan sive scale server farm with condit ions	
		Dyna mic provisi oning in multi- tenant service clouds	User Interface- Tenant Selector- Customizer (UTC) model and approach	A coordinate d calculatio n for coordinati ng occupant functional ities with a customer's prerequisi tes	Stati cs	Dispos al of excess inhabit ant functio nalities all togethe r to prune the inquiry space	Just upper bound dispen sing the assets		5	Time Base d RPM /[25]	A resourc e provisi oning framew ork for bioinfo matics applica tions in multi- cloud environm ents	Workflow improveme nt mechanism	Improves submitted theoretical work processes by misusing parallelis m.	Dyn amic	Decline s the work process executi on time for a given spendin g plan.	Requir es dynam ic bunch arrang ement and dynam ic scalin g of the proces s hubs
									6	Adap tive Base d RPM /[27][29][1]	Revenu e Maxim ization Using Adap tive Resourc e Provisi oning in Cloud Compu ting Enviro nments	Resource allocating using Queuing Theory and propose optimal solution	Boosting the supplier's income through SLA- based dynamic asset portion	Dyn amic	Expand supplie r's income in light of executi on	Compl exity
4	Ener gy Base d RPM /[26][4]	Energy Efficie ncy Optimi zation With Statisti cal QoS Provisi oning for Energy Harvest ing Networ	Statistical delay- bounded QoS-driven power control policies for energy harvesting	Determine the shut shape arrangeme nts of ideal power control strategies to the EEE expansion issue under imperativ es	Stati cs	Boost the success ful vitality effectiv eness	EEE augme ntatio n issue is unrave led under just the batter y limit requir ement			An Adap tive Score Model for Effecti ve Bandwi dth Predicti on & Provisi oning in the Cloud Networ k	Develop prediction- based resource measureme nt and provisionin g strategies	Forecast system utilizes factual models which can conjecture the future surge in asset prerequisi te; in this manner empoweri ng proactive scaling to deal with transient	Dyn amic	Improv e the adequacy of versatil e asset allotme nt regardi ng both executi on and cost.	For the most part for web based busine ss applica tions	

				bursty workload in a controllable way.			
		Empirical prediction models for adaptive resource provisioning in the cloud	Auto-Regressive Integrated Moving Average (ARIMA) model augmented with a general class of Adaptive Conditional Score Models (ACS).	Models offers enhanced prescient capacity with bring down gauge blunders	Dynamic	Adjust adequately to homeless people while ready to keep up application SLAs	Keep on testing for substantial scale situation

C. Benefits of Resource Provisioning

- Effective cloud asset provisioning lessens execution time of cloud workloads.
- Better asset usage under various prerequisites of need and maintains a strategic distance from over provisioning and under provisioning.
- No provisioning delay and lesser odds of asset disappointment because of productive administration of assets.
- No long VM start-up delay gives provisioned assets promptly in compelling cloud asset provisioning.
- Increase the vigor and limit make traverse of work process at the same time.
- Meet even strict application due date with least spending consumption and increments worldwide benefit.
- Power utilization lessened without infringement of SLA in powerful cloud asset provisioning.
- Efficient adjusting of load by proficient dissemination of the workloads on accessible assets.
- Improve client due date infringement rate because of assets provisioning before asset planning.
- Effective cloud asset provisioning decreases lining time in workload line.

VI. CONCLUSION

The examinations which were contemplated above are attempting to improve and use the assets. A few techniques were said here which utilized diverse parameters as an objective for asset provisioning, for example, reaction time, dismissal rate, benefit level assertion (SAL) infringement rate, cost and so forth. For provisioning arranging should take proper provisioning times, Provisioning assets too early will squanders our assets and in this way our cash, on the opposite side provisioning assets past the point of no return will cause

possibly SLA infringement and makes the clients furious. This paper displays an exhaustive audit on successful assets provisioning in cloud. We have talked about asset provisioning when all is said in done. We outlined asset provisioning component and correlation between various assets provisioning instrument as far as a superior execution, focused and productivity to meet the required SLA enhanced the asset execution and brought down the power utilization. So we reason that portion of assets in light of kind of workload. Appropriate coordinating of workload and asset can enhance the execution significantly. It is extremely troublesome for supplier to recognize the quantity of assets required precisely for given workload from asset pool, since assets might vary in one or other criteria for example, asset limit, cost and speed. User can choose fitting asset provisioning component based on QoS necessities of workload/application portrayed through assessment and correlation of asset provisioning in cloud. Differentiation and evaluation of asset provisioning systems in cloud can help to choose the asset provisioning instrument in view of workload's QoS requirements. Cost can be lessened in the conveyed cloud benefit if assets are held ahead of time. We trust this paper will spur specialists to investigate and figure another system to explain issues in designating and checking assets in distributed computing and this research work will be beneficial for researchers who want to do research in area concerning to resource management such as cloud resource provisioning and resource provisioning mechanism.

REFERENCES

- [1] "Empirical prediction models for adaptive resource provisioning in the cloud", Sadeka Islam and et al. 2011.
- [2] "A cost-aware mechanism for optimized resource provisioning in cloud computing", Safiye Ghasemi and et al.2017.
- [3] "An autonomic resource provisioning approach for service-based cloud applications: A hybrid approach", Mostafa Ghobaei-Arani and et al. 2017.
- [4] "DRL-Cloud: Deep Reinforcement Learning-Based Resource Provisioning and Task Scheduling for Cloud Service Providers", Mingxi Cheng and et al. 2018.
- [5] "Resource Provision Algorithms in Cloud Computing: A Survey", Jiangtao Zhang and et al. 2016.
- [6] "Combinatorial Double Auction-based Resource Allocation Mechanism in Cloud Computing Market", Seyede Aso Tafisiri and et al. 2017.
- [7] "A Novel Agent Based Autonomous and Service Composition Framework for Cost Optimization of Resource Provisioning in Cloud Computing", Aarti Singh and et al. 2015.
- [8] "Review on: Resource Provisioning in Cloud Computing Environment", Sagar Girase and et al. 2013.
- [9] "Horizontal scaling and aggregation across heterogeneous clouds for resource provisioning", Kirthica S and et al. 2017.
- [10] "Cloud resource provisioning: survey, status and future research directions", Sukhpal Singh and et al. 2016.
- [11] "A Survey on Resource Allocation and Monitoring in Cloud Computing", Mohd Hairy Mohamaddiah and et al. 2014.
- [12] "Cloud resource provisioning and SLA enforcement via LoM2HiS framework", Vincent C. Emekaroha and et al. 2012.
- [13] "Resource Provisioning in Single Tier and Multi-Tier Cloud Computing: —State-of-the-Art", Marwah Hashim Eawna and et al. 2015.
- [14] "Multi-agent based dynamic resource provisioning and monitoring for cloud computing systems infrastructure", Mahmoud Al-Ayyoub and et al. 2015.
- [15] "QoS-guaranteed resource provisioning for cloud-based MapReduce in dynamical environments", Xiaoyong Xu and et al. 2018.

- [16] "Failure-aware resource provisioning for hybrid Cloud infrastructure", Bahman Javadi and et al. 2012.
- [17] "CHOPPER: an intelligent QoS-aware autonomic resource management approach for cloud computing", Sukhpal Singh Gill and et al. 2017.
- [18] "Resource provisioning for data-intensive applications with deadline constraints on hybrid clouds using Aneka", Adel Nadjaran Toosi and et al. 2017.
- [19] "Automated and Dynamic Application Accuracy Management and Resource Provisioning in a Cloud Environment", Smita Vijayakumar and et al. 2010.
- [20] "A SLA-based Spark Cluster Scaling Method in Cloud Environment", Yoori Oh and et al. 2016.
- [21] "A Cloud Service Broker for SLA-based SaaS Provisioning", Elarbi Badidi. 2013.
- [22] "SLA Analytics for Adaptive Service Provisioning in the Cloud", Obinna Anya and et al. 2016.
- [23] "SLA-based Secure Cloud Application Development: the SPECS Framework", Valentina Casola and et al. 2016.
- [24] "Dynamic provisioning in multi-tenant service clouds", Lakshmi Ramachandran and et al. 2012.
- [25] "BioCloud: A resource provisioning framework for bioinformatics applications in multi-cloud environments", Izzet F. Senturk and et al. 2018.
- [26] "Energy Efficiency Optimization With Statistical QoS Provisioning for Energy Harvesting Networks", Ya Gao and et al. 2017.
- [27] "Revenue Maximization Using Adaptive Resource Provisioning in Cloud Computing Environments", Guofu Feng and et al. 2012.
- [28] "Adaptive Resource Allocation and Provisioning in Multi-Service Cloud Environments", Ayoub Alsarhan and et al. 2017.
- [29] "An Adaptive Score Model for Effective Bandwidth Prediction & Provisioning in the Cloud Network", Abiola Adegboyega. 2015.
- [30] "An Algorithm Based On Ant Colony Optimization for Provisioning Virtual Sensor in Sensor Cloud", Marcus Lemos and et al. 2017.
- [31] "Self-Tuning Service Provisioning for Decentralized Cloud Applications", Raul Landa and et al. 2013
- [32] "QoS based resource provisioning and scheduling in grids", Rajni Aron and et al. 2013

University Notification Subscription System using Amazon Web Service

Babur Hayat Malik, Zaheer Mehmood Dar, Sabah Mubarak Kayani, Mahnoor Dar, Muhammad Hassan Shafiq, Imran Kabir, Fatima Masood, Hamna Zakriya, Asad Ali
Department of Computer Science and Information Technology,
University of Lahore
Gujrat Campus, Pakistan

Abstract—Publish-Subscribe (Pub-Sub) system is an asynchronous communication service widely used in server-less and micro-services architecture. In a Pub-Sub system, publisher publish message to a topic that is immediately received by all of the subscribers of that topic. Nowadays, students face number of problems regarding admission details, assignments, offered courses, fee schedule, etc. Many a times, they missed the deadlines and it affects their studies. This paper is focused on issues faced by students regarding message delivery, duplication of data and heavy traffic, etc. It should be overcome by using amazon web services to make optimize product and to make it flexible for university Pub-Sub system. Implement the cloud services by using hybrid technique, i.e., content based and topic based architecture. It also explained the multitude use-case of university notification system which leads to make it more adaptable as subscriptions are identified with particular data content.

Keywords—Publish-Subscribe system; content and topic-based; university notification; Amazon web services

I. INTRODUCTION

A Publisher-Subscriber system defined as Pub-Sub is in a correspondence worldview generally utilized to give occasion scattering between inexact publishers (distributors) and subscribers (follower) [1]. Distributors must distribute the message which are coordinated and conveyed by Pub-Sub operators (called brokers) to subscribers in light of their enrolled subscriptions. The main focus is to put on coordinating and conveying a high throughput of occasions to these steady subscribers [13]. The representative speaks with various substances (e.g., distributors and follower), coordinates the reasonable client prerequisite, and transmits clients' information [14]. Pub-Sub system implements on application layer of OSI model.

In Pub-Sub system, the distributed applications provide instant event notifications. This model allows event-driven architecture and asynchronous processing for performance improvement, reliability, scalability and versatility [2]. Gregor Hohpe and Bobby Woolf, defines the Competing Consumers design as “Competing Consumers are numerous clients that are altogether to get messages from a Point-to-Point Channel. At the point when the channel conveys a message, any of the buyers could possibly get it. The information system figures

out which customer really gets the message, yet as a result the buyers rival each other to be the collector. Once a customer gets a message, it can delegate to whatever is left in its application to help process the message.”

There are different types of Pub-Sub system like, content based, type based and topic based. A more adaptable yet additionally complex worldview in the Pub-Sub conspires is content-based membership. It gives greater adaptability to the supporter by giving more control in buying in an occasion in view of the genuine substance of the occasion. It enables endorser of force set of limitations as condition in shaping an inquiry on an occasion warning (otherwise called channel). Making a notice utilizing a channel gives supporters a more refined route for buying in occasions.

The features of Pub-Sub system are Push delivery by using multiple delivery protocols, fan-out, filtering, durability and security. The main purpose of using Pub-Sub system is push delivery and message durability. To develop Pub-Sub system, need messages from authentic sources. Students can receive the reliable data. This paper will discuss the categories of Pub-Sub system, problem statement; propose system and its working and discussion and conclusion.

II. CATEGORIES OF PUB-SUB SYSTEM

There are three main categories of Pub-Sub system used to implement any system, i.e, named as: content based, type based and topic based.

A. Content-Based

In content based Pub-Sub system, publisher can publish the content and subscriber follows the content as per need. In this system, publisher publishes the content on the message system. After this, message broker broke the message to know message content. Then send this message to the particular subscribers who want to know about this content. Message brokers use filtering pattern that applied on consumers subscription to elect events by using a subscription language (constraints <>) [3].

In the below Fig. 1, publisher publishes the message over Pub-Sub system [25]. Then message delivers to the concerned subscriber [22].

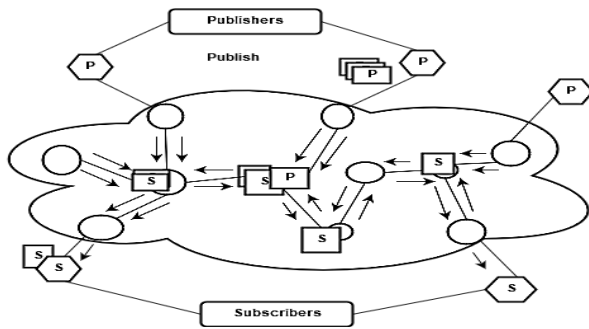


Fig. 1. Content based Pub-Sub system [25].

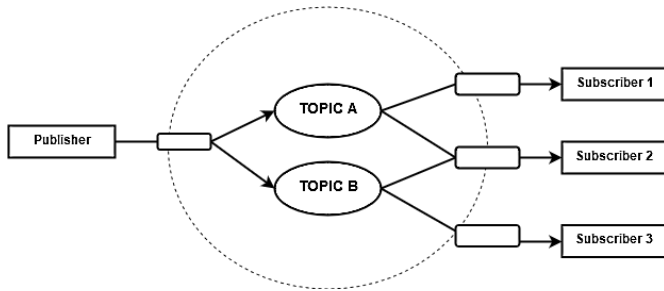


Fig. 2. Topic based Pub-Sub system [12].

B. Topic-Based

Topic based Pub-Sub design; sender sends messages with the name of a topic. The message is sent to the message system and then delivered to all the applications that want to receive messages on that topic [6].

In Fig. 2, publisher publishes topic A and topic B and subscribers subscribe the topic according to their interest. One subscriber can subscribe many topics as per need [12].

C. Type-Based

In type based Pub-Sub defined events in its interface [12]. It is based on static and dynamic schemas as shown in Fig. 3. To implement this Pub-Sub system, need languages which support structural reflection. For this, no need for specific events (e.g. Java introspection). In other languages, event can subtype an introspective event type.

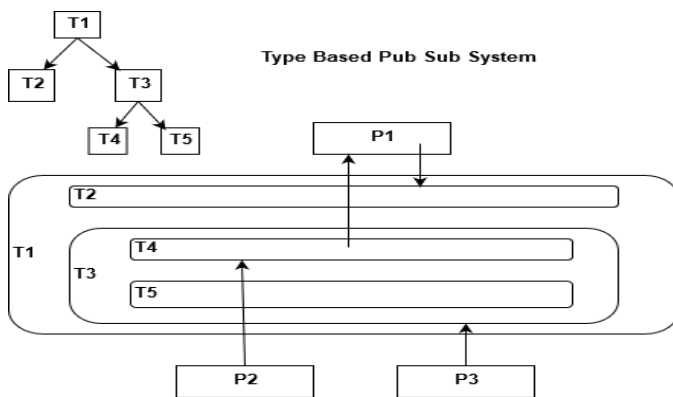


Fig. 3. Type based Pub-Sub system [12].

III. PROBLEM STATEMENT

University Pub-Sub system is particularly dealing with all major entities of system e.g. Admin, Teacher and students. The propose application supports entire range of notification data to facilitate our subscribers on other side they also face some technical issues regarding the implementation of Pub-Sub architecture on university subscription system. Some of them are mentioned below.

A. Delivery Notification Unguaranteed

Perfect knowledge is obvious for publisher but in this system publisher message status is not guaranteed. As publisher is unaware by the system delivery service, that may cause loss of some important notification for both subscribers and publishers [24].

B. Decreased Performance

Public systems accessibility cause high risk of unattended attacks, broker can be vulnerable to attack this system easily. In propose system use case there is a large amount of register subscribers (HR, Teachers and Students) that makes system overloaded. Highly communicative, but on runtime it requires complex protocols to implement subscription functionality [2].

C. Redundant Data

Publish-Subscribe system offers multiple instances of loggers that can run concurrently looks identical. However, in system designs it allows for a high level of redundancy. Such replications in data make it persistent [17].

D. Inflexible Data

The propose system then firstly make its paper prototype then after analysis, propose the structure then it is become difficult to change when system architecture already established. With a specific end goal to change the structure of the messages, the greater part of the system must be adjusted to acknowledge the changed arrangement.

IV. PROPOSED SYSTEM

In this section, propose a Pub-Sub system for educational institutes. At first, describe simple working of Pub-Sub system then explains the architecture how it would be implemented in any use case or helps users to follow proper university subscription system.

A. Pub-Sub system

Publish-Subscribe is a software design pattern that describes the relationship between the flow of users, services or services of all publisher messages as shown in Fig. 4. The so-called Pub sub usually works this way: the publisher (i.e. any data source) pushes the user in the message (i.e. the data recipient) by streaming interest in a real-time feed called a channel (or topic). When a new message is posted on this channel, the subscribers of all the specific publisher channels are notified immediately and the message data (or payload) is received along with the notification. In daily use, especially in the Internet of Things, automation, network operations or distributed cloud environments, an intermediate layer called a message broker is usually required to handle the distribution

and filtering of messages, and also provides a low latency messaging private network infrastructure [18].

In proposed Pub-Sub system used topic based architecture. In topic based architecture, publisher publish message on any topic using Pub-Sub system. Then message broken works and sends it to the concerned subscriber.

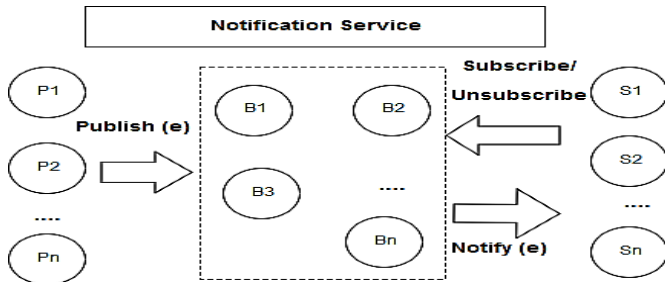


Fig. 4. Subscription model [5].

Centralized event breaking system is a key factor of current Pub-Sub systems framework that relies on a single event broker. If it working well and its working is down then the event propagation will be negotiated within the current framework whereas system vulnerability of whole system would be enhanced if depending on a single event broker [4].

In order to decrease the liability additional architecture can be made, such as received messages receipts receiving. With that component included, feedback can be given to the distributor with regards to the status of the subscriber. It is more adaptable as subscriptions are identified with particular data content and, thus, every packet of information can really be viewed as a solitary dynamic consistent channel. This exponential enlargement of Potential logical channels enlarges exponentially has changed the implementation level working of Pub-Sub system. The pub- sub operational models description is given in Fig. 5.

Event notification Pub-Sub communication system and the compatible web service technology giving a combination of emerging technology named as web service based notification system [7]. Event driven and service oriented architecture both are emerging technology giving attention to web service based notification system [11]. It helps in integrating applications either within or outside the organization. Web service Event [20] and Web service-Notification [21] are two major specifications for such systems.

Create Channel	POST/channel
Subscribe Channel	Subscribe
Publish Events	Publish
Read Events	Get Event Messages
Unsubscribe Channel	Delete or Unsubscribe

Fig. 5. Pub-Sub operational model [5].

B. University Pub-Sub Architecture

The proposed system is implemented by the Amazon Simple Notification Service (Amazon SSN). This is a web service that enables end quote or managing or sending messages to subscribe to customers In Amazon SMS, two types of customers - publishers and consumers - also known as producers and consumers [15]. Publishers send message asynchronously with subscriptions to create messages and send messages on logical access points and topics of communication channels. Subscriber (i.e. web servers, email addresses, Amazon SQS queues, AWS functions) one of the support protocols for message or notification (such as MMS, SMS, HTTP/S, Email). Subscribers to a topic Publish and receive messages on a specific topic through a message service provider. Extensibility is achieved by distributing topic sets in a number of message providers or through a cluster of providers. Different ways of communicating are topic-based university Pub-Sub middleware, such as Amazon SNS server [16]. Among these, the filter group that specifies the event subject by the subscriber connected to the publisher through the broker network.

Then, the Pub-Sub middleware is responsible for forwarding the publisher’s events to relevant users throughout the network. Event filters distribution, matching process to achieve high scalability, among a large number of brokers [8]. A diagrammatic view of university pub-sub architecture is shown in Fig. 6.

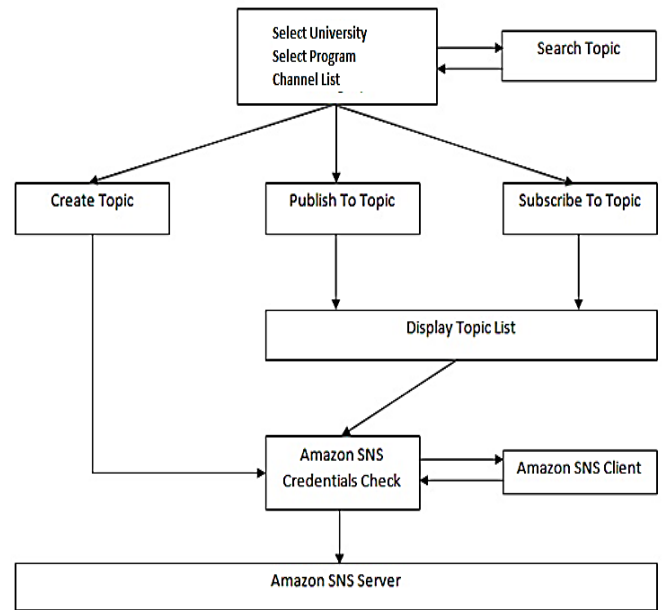


Fig. 6. University Pub-Sub architecture [19].

V. UNIVERSITY PUB-SUB PROPOSED DESIGN

Here is the detailed design of system working as shown in Fig. 7 It is actually a hybrid (content and topic based) publisher-subscriber system.

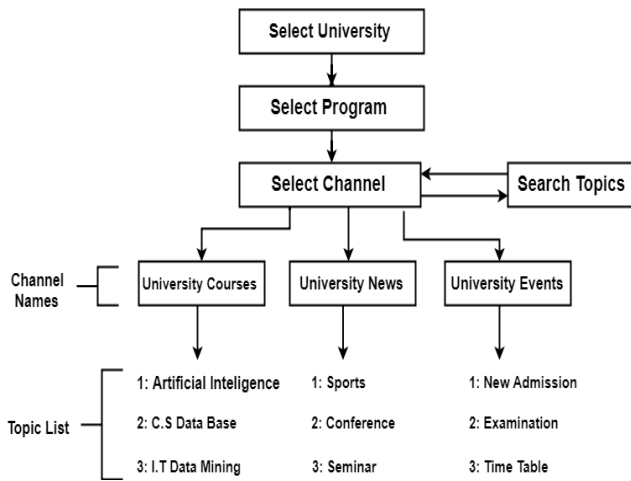


Fig. 7. University Publish-Subscribe proposed design [19].

VI. UNIVERSITY PUB-SUB SYSTEM USE CASE

It depicts overall university notification system working. Defining proposed system scope where two main actor's publisher and subscriber play an important role. System is divided into two main modules, i.e. publisher as an admin module or subscriber as a user module. Fig. 8 depicts the use case for university pub-sub system use case.

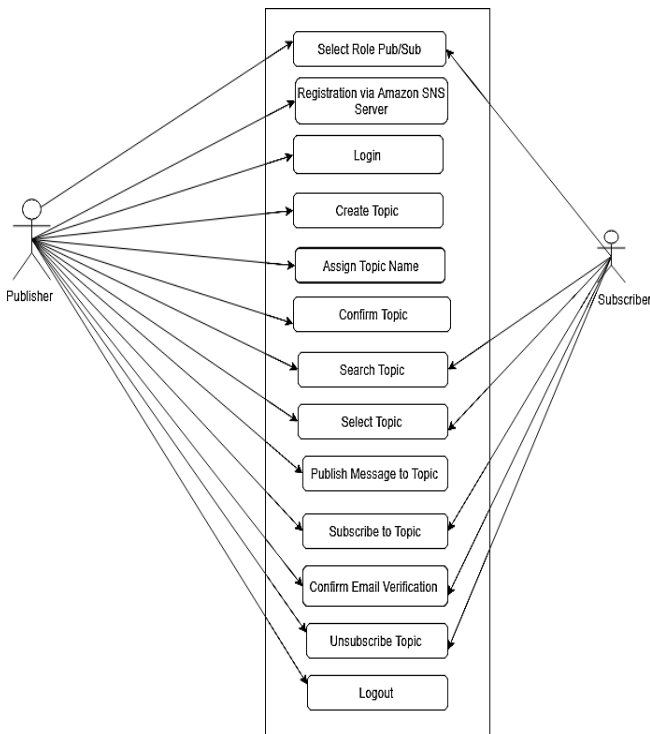


Fig. 8. University Pub-Sub use case diagram [9].

A. University Pub-Sub Actor's Use Cases

There are number of actors in proposed system which plays different roles as a publisher or subscriber.

1) New publisher use case diagram

Fig. 9 shows the pictorial representation of new publisher use case.

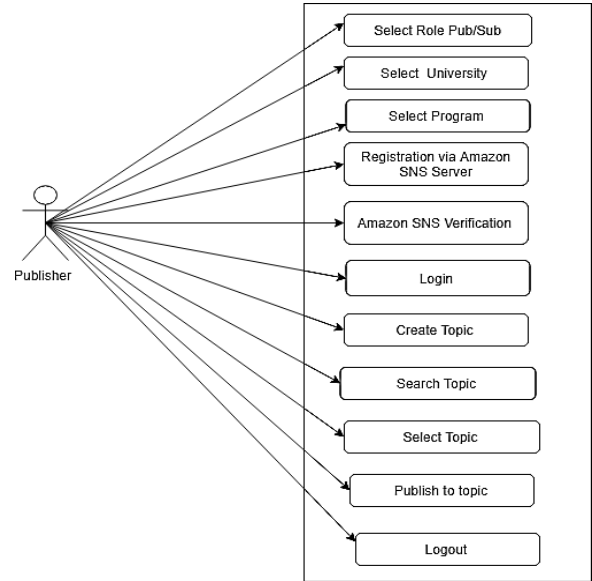


Fig. 9. New publisher use case diagram [10].

2) Teacher Use Case

In publisher module user has a right to publish any content to relevant topic or also subscribe channels within same frame. All use case activities described in the following Fig. 10.

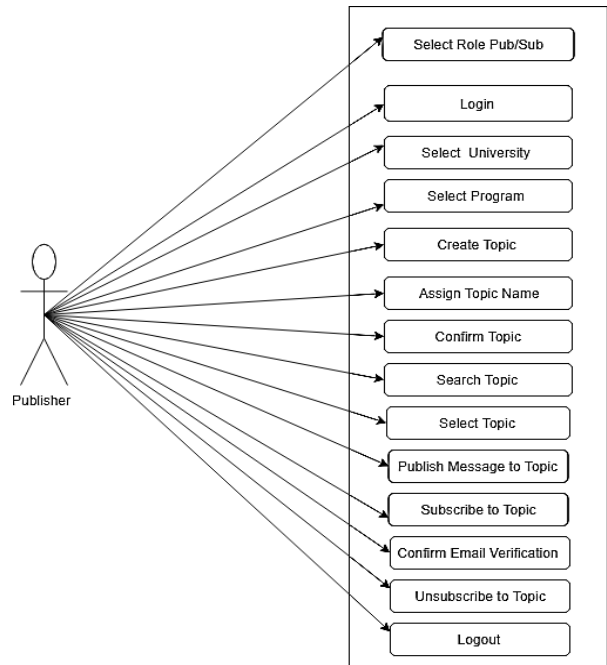


Fig. 10. Teacher as publisher role use case diagram [9].

3) Student Use Case

Student would always be in the subscriber module unless university management authorize him any rights as shown in Fig. 11.

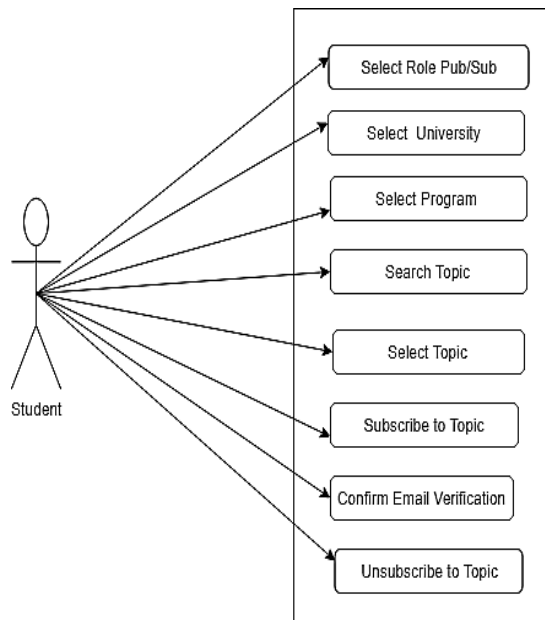


Fig. 11. Student as subscriber use case diagram [10].

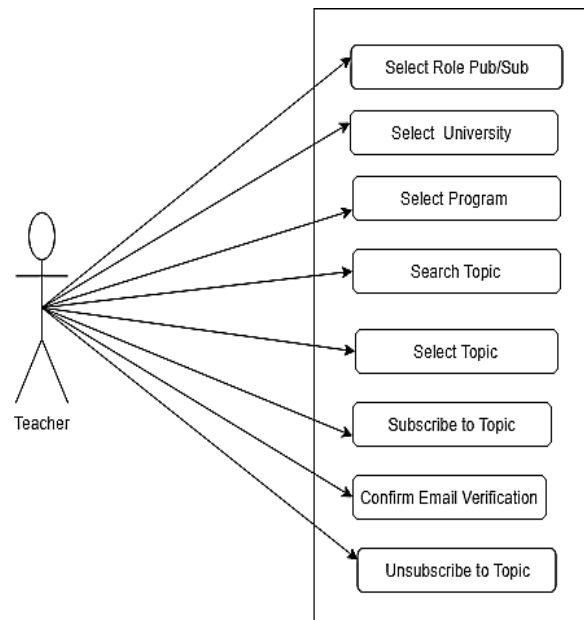


Fig. 13. Teacher as subscriber use case diagram [10].

4) Examiner Use Case

Here examiner is a person who has a right to publish any examination related news to relevant channel. If he is new registered user then he will first create the topic then publication should be done. The roles of examiner as publisher are shown in Fig. 12.

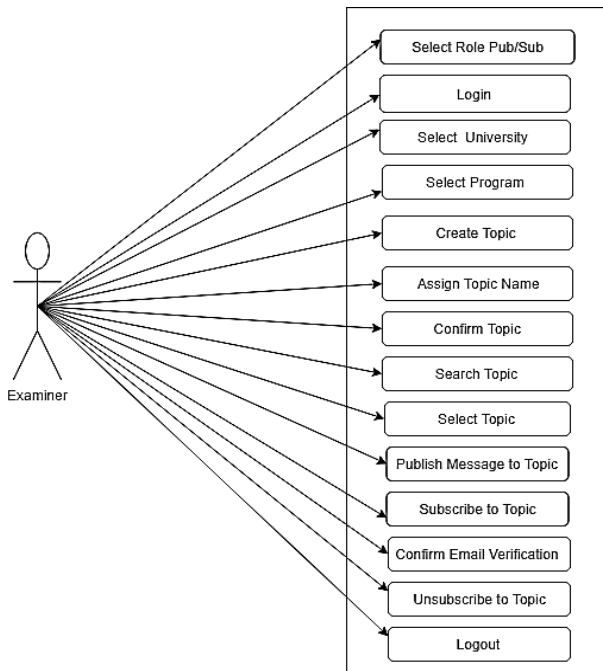


Fig. 12. Examiner as publisher use case diagram [9].

5) Teacher as a Subscriber Use Case

As mentioned earlier in the use case of teacher role. He can also subscribe to channel number of activities mentioned in the use case in Fig. 13.

VII. DISCUSSION

After extensive analysis, some findings that loosely coupled modules, architecture flexibility and reliability of system should lead to a model Pub-Sub system to smoothly run university affairs that provide quality time delivery of event notification.

A. Flexibility in Architecture

You can definitely include more endorsers and add group extension if message creation supplants message utilization without code change. That's why in Pub-Sub model there should be low coupling between each module as message passed to the subscribers by the publishers without any trouble because it makes system highly independent.

B. Configuration Ease

Pub-Sub system is easy to configured and also support different models of subscriptions. It helps publishers and subscribers to distribute and receive the messages using fewer resources by providing offline notification [23].

C. System Service Availability

The communication framework transports the distributed messages just to the applications that are bought in to the relating subject. Since various physical endorsers can have a similar membership, there's ensured bolster for high accessibility, for the subject itself as well as for the theme supporters.

VIII. CONCLUSION

In order to decrease the liability additional architecture can be made, such as received messages receipts receiving. With that component included, feedback can be given to the distributor with regards to the status of the subscriber. Because subscriptions with specific data content are identified, they are more adaptable which helps to make system more flexible, highly configured and proper distributed system for reliable

communication thus, every packet of information can really be viewed as a solitary dynamic consistent channel. This exponential enlargement of Potential logical channels enlarges exponentially has changed the implementation level working of Pub-Sub system.

REFERENCES

- [1] Joao Paulo de Araujo, Luciana Arantes, Elias P. Duarte Jr., Luiz A. Rodrigues, Pierre Sens, "A Publish/Subscribe System Using Causal Broadcast Over Dynamically Built Spanning Trees", 29th International Symposium on Computer Architecture and High Performance Computing, 2017.
- [2] Tarek R. Sheltami, Anas A. Al-Roubaiey, Ashraf S. Hasan Mahmoud, "A survey on developing publish/subscribe middleware over wireless sensor/actuator networks", Springer Science+Business Media New York 2015
- [3] César Canas, et. Al, "Self-Evolving Subscriptions for Content-Based Publish/Subscribe Systems", IEEE 37th International Conference on Distributed Computing Systems, 2017.
- [4] Hiroki Nakayama, Dilawaer Duolikun, Tomoya Enokido, Makoto Takizawa, "Causally Ordered Delivery of Event Messages with Keyword Vectors in P2P Publish/Subscribe Systems", IEEE 29th International Conference on Advanced Information Networking and Applications, 2015.
- [5] Muhammad Agus Triawan, Hilwadi Hindersah, Desta Yolanda, Febrian Hadiatna, "Internet of Things using Publish and Subscribe Method Cloud-based Application to NFT-based Hydroponic System", IEEE 6th International Conference on System Engineering and Technology (ICSET), Oktober 3-4, 2016 Bandung – Indonesia
- [6] Ryohei Banno, Susumu Takeuchi, Michiharu Takemoto, Tetsuo Kawano, Takashi Kambayashi, Masato Matsuo, "A Distributed Topic-based Pub-Sub Method for Exhaust Data Streams Towards Scalable Event-driven Systems", IEEE 38th Annual International Computers, Software and Applications Conference, 2014.
- [7] A. Carzaniga, D. S. Rosenblum, et al. Design and Evaluation of a WideArea Event Notification Service. ACM Transactions on Computer, 2004, pp.343-356
- [8] G. Cugola, E. D. Nitto. The JEDI event-based infrastructure and its application to the development of the OPSS WFMS. IEEE Transactions on Software Engineering (TSE), 2001, pp.827-850
- [9] S. Shukla, P. Saikia, M. Cheung, J. She and S. Park, "Effectiveness of Mobile Notification Delivery", Mobile Data Management (MDM)", 18th IEEE International Conference, pp. 21-29, 2017.
- [10] Carzaniga, A., D.S. Rosenblum, and A.L. Wolf, Achieving scalability and expressiveness in an internet-scale event notification service. Proceeding of Nineteenth ACM Symposium on Principles of Distributed Computing (PODC 2000), 2000.
- [11] Banavar, G., et al., An efficient multicast protocol for content-based Publish-Subscribe systems. Proceedings of the 19th International Conference on Distributed Computing Systems (ICDCS'99), 1999. 11. Pietzuch, P. and J. Bacon. Hermes: A Distributed Event-Based Middleware Architecture. in Workshop on Distributed Event-Based Systems (DEBS). 2002.
- [12] G. Hohpe and B. Woolf, Enterprise integration patterns. Boston: Addison-Wesley, 2004.
- [13] Altinel, M. and M.J. Franklin., Efficient Filtering of XML Documents for Selective Dissemination of Information. Proc. of VLDB 2000, 2000.
- [14] Diao, Y. and M.J. Franklin, High-Performance XML Filtering: An Overview of YFilter. IEEE Data Engineering Bulletin, 2003(March, 2003).
- [15] Barton, C.M., et al., Streaming XPath Processing with Forward and Backward Axes. Proc. of ICDE, 2003.
- [16] Feng Peng, S.S.C., XPath Queries on Streaming Data In Proc. of SIGMOD, 2003.
- [17] E.Onica,P.Felber,H.Mercier,andE.Rivière, "Confidentiality-preserving publish/subscribe: A survey," ACM Comput. Surv., vol. 49, no. 2, p. 27, 2016
- [18] Baldoni, R. and Virgillito, A. 2005. Distributed event routing in publish/subscribe communication systems: a survey. Technical Report TR-1/06. The Computer Journal, vol.50 (2), pp.444 -459
- [19] Klein, A., Mannweiler, C., Schneider, J., Schotten, H.D.: Access Schemes for Mobile Cloud Computing. In: Proceedings of the 2010 Eleventh International Conference on Mobile Data Management, pp. 387–392, 2010.
- [20] P. Eugster, "Type-based publish/subscribe", ACM Transactions on Programming Languages and Systems, vol. 29, no. 1, p. 6-es, 2007.
- [21] S. Lakhota, "Why we Publish, Where we publish and What we Publish?", Proceedings of the Indian National Science Academy, vol. 80, no. 3, p. 511, 2014.
- [22] Y. Zhao and J. Wu, "Building a reliable and high-performance content-based publish/subscribe system", Journal of Parallel and Distributed Computing, vol. 73, no. 4, pp. 371-382, 2013.
- [23] S. Oh, J. Kim and G. Fox, "Real-time performance analysis for publish/subscribe systems", Future Generation Computer Systems, vol. 26, no. 3, pp. 318-323, 2010.
- [24] L. H., R. S. and N. R., "Review for Event Delivering Techniques in Publish/Subscribe Scheme", International Journal of Computer Applications, vol. 132, no. 16, pp. 1-5, 2015.
- [25] P. Narasimhan and P. Triantafillou, Middleware 2012. Heidelberg: Springer, 2012.

Performance Measurement Model of Mobile User Connectivity in Femtocell/Macrocell Networks using Fractional Frequency Re-use Scheme

Mehrin Anannya

Information & Communication Technology
Bangladesh University of Professionals
Mirpur Cantonment, Mirpur-12, Dhaka, Bangladesh

Riad Mashrub Shourov

Skills for Employment Investment Program (SEIP)
Prabashi Kallyan Bhaban, Ramna,
Dhaka, Bangladesh

Abstract—Technologies are traversing to its new dimensions every day. As part of this progression, mobile cellular system is at the summit of its constant advancement. The usage of Femtocells in mobile cellular system has created a massive impact on its architecture. Likewise, the incorporation of femtocells in macrocells for 4G mobile network communication services (like-voice calls, data services, etc.) among mobile stations within few meters has been one of the promising approaches. There is a femto access point (FAP) in Femtocell which handles the authorization of the user around it. Among the three various access methods, FAP allows only the authorized users except the macro cell users in Closed Access Method (CAM). But for Open Access Method (OAM), any type of crossing macrocell user within the radio coverage of femtocell and the femtocell users can get FAP access. To reduce the cross-tier interferences OAM is more efficient, because it deals with both type of users within the femtocell coverage. This paper proposes a performance measurement model for mobile connection probability depending on the mobility factor of mobile users and the communication range in femtocell/macrocell networks. Furthermore, a derivation has been done to get the optimum result from the outage and connectivity probability under different number of femtocells and mobile users. Finally, to maximum the spectral efficiency for the probable frequency allocation, a Fractional Frequency Re-use scheme among the networks has been proposed.

Keywords—Femtocell; macrocell; cross-tier interferences; co-tier interferences; closed access methods; open access methods; connectivity probability; mobility factor; outage probability; fractional frequency re-use scheme

I. INTRODUCTION

With the advancement of technology, wireless mobile communication is on high demand around the world. This demand arises not only for the voice communication but also for data services. The need for consistent connectivity between the users at two ends providing high speed communication with low cost and without any loss of data and interruption are increasing. So, a high-speed voice and data transmission such as, voice calls, video calls and rapid and faster internet facilities having clear video images without any kind of interruption over the network and without any loss, needs to be ensured to satisfy consumers by enhancing the system coverage and capacity and reducing the capital. Depending on the different operators, the mobile cells

communicate with their respective base stations. In this case, the base stations are considered as the Macrocell Access Point (MAP), which provides an area of coverage to the mobile users. A Macrocell is nothing but a cell that provides wide radio coverage with high power cellular base station in a mobile phone network (tower). Due to substantial wireless communication between the mobile users, its workload climbs along with the increase of mobile users.

A. Necessity of Femtocell Considering the Problems of Macrocell in Increasing the Overall Performance and Signal Quality of Mobile Communication

To serve this vast number of mobile users, the performance and quality of the macrocell degrades due to slow connectivity, call drops, reduction of sound intensity, etc. Owing to this, co-tier and cross-tier interference also increases.

This workload of macrocells can be subdivided among some femtocells which work within the common home range. A femtocell is a short-range home area broadband network which gives coverage of about 10 meters providing better indoor voice and data communication. A femtocell has a base station named Femtocell Access Point, which mainly communicates with the Macrocell to build its own network under that Macrocell. Since, in this kind of network transmitter and receiver remain very close to each other, a high-quality link is created between them providing large number of spatial reuse with low power transmission and power wastage. It is seen that, 2/3 of voice communication and 90% of data communication occur within the home range of a mobile user. So, for enterprise or home environments, to offload the traffic on the macrocell by using increasing system capacity, providing high probability of connectivity and utilizing spatial reuse of resources which provides highest spectrum efficiency, so that every user located at every corner of the macrocell can have clear connectivity, femtocell needs to be used. Due to femtocell, the cross-tier interference and handoffs under the macrocell decreases.

The users can access femtocells depending on three access mechanisms, namely, Closed Access Mechanism (CAM), Open Access Mechanism (OAM) and Hybrid Access Mechanism (HAM). In this research, we will utilize open access mechanism for Femtocells.

Macrocell and femtocell has been shown in Fig. 1 and 2, respectively.

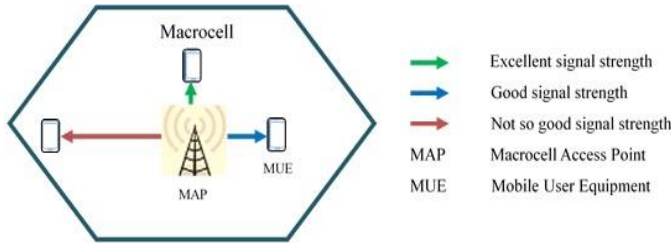


Fig. 1. Macrocell.

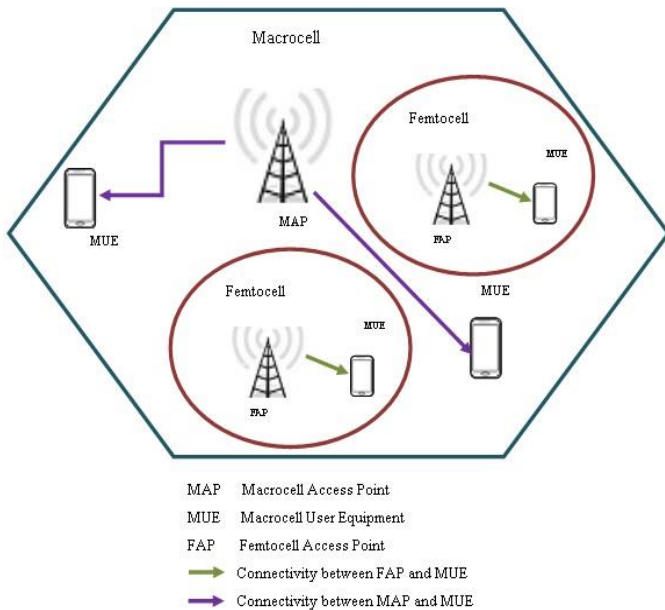


Fig. 2. Femtocell.

II. LITERATURE REVIEW

A. Background

Firstly, the connectivity with the femtocell needs to be confined which is maintained by access control of the femto cell. Regarding this, the merits and demerits of existing access methods of femtocells emphasizing on its technical impact and also a business model is described in [1] which headed towards the necessity of hybrid access methods along with several other models. In [2] a stochastic geometric model has been employed to enhance the spatial reuse of the femtocells meeting the per-tier outage constraint using cognitive radio. At the Primary User (PU) transmitter, a capacity-outage probability to the PU better than the beacon transmitter has been examined in [3] with the statistical model that has been developed so far which imposes less interference. Furthermore, the model is prolonged in investigating the cooperative sensing effect and capacity-outage performance for maximizing the likelihood cooperative detection techniques. A brief idea about the requirement of the femtocells along with the technical and business arguments for femtocells and its challenges, way of overcoming the challenges and the state-of-the-art on each front of it has been discussed in [4]. Not only for wireless communication but also in case of wired communication, it is seen that the performance of femtocells is better than the macrocells. The simulation result in [5] shows that via vehement wireless spectrum spatial reuse, the sententious capacity gain of the areal (per single area) can be achieved giving a strong signal strength in case of femtocell and weak signal strength in macrocell. In [6] an interpretive model framework has been demonstrated for computing the plausibility of mobile user connectivity depending on the femtocell density, communication range, mobility factor, and user density in femtocell/ microcell, which by examining the performance of outage probability and spectral efficiency in that networks found to be efficient during planning of Macro cellular networks integrated with Femto cellular networks. But the consummation of the transmission quality largely bets on the interference for different distances. For different measurements of distance, the interference aware and SINR estimation of femtocell networks by using frequency reuse mechanisms for LTE has been done in [7,1]. In [8], an algorithm on sub-channel allocation was developed using graph-theoretic approach in which a grouping of femtocell users is made into disparate clusters for subduing the interference within them and optimizing the femtocell throughput in densely populated femtocells deployment with adaptive power allocation to enhance the system throughput. In [9] a discussion has been done about a survey on different level of development approaches, qualitative comparison and open challenges of interference along with asset management in orthogonal frequency-division multiple access (OFDMA) based femtocell networks. In [10] discussions have been made on evolution, characteristics, design and deployment aspects for femtocell discovery by active call hand-in and idle mobiles switching from macrocell to femtocell in cdma2000-based femtocell systems to emend the performance and enrich user experience with femtocells. In [11] numerical results have been found out by making a co-channel existence for

B. Objectives

- Observe the disadvantages of macrocells and focus on the advantages of femtocells based on the co-tier and the cross-tier interference.
- Develop a performance measurement model in order to observe the variation of connectivity probability against the number of mobile users in Open Access Method environment.
- Derive the outage probability in terms of mobile users and threshold of detection (SINR).
- Evaluate the variation of spectral efficiency against SINR.
- Finally, introduce a proposal of a fractional frequency re-use scheme for Self-Organization Network (SON) based on femtocell architecture.

C. Scope

To increment the connectivity probability by reducing outage probability beyond the level that has been proposed.

proximate indoor femtocells and outdoor macrocell users demonstrating that the desired femtolink provides a robust performance than the macrolink having some equal interference by other femtocells and macrocells. In [12] some simulation results have been shown which interrogates regarding uplink and downlink capacity of macrocell and femtocell demonstrating that the macrocell pursuance derogates due to the locations of femtocell BS and macrocell UE because of femtocell transmit power which is proportional to the number of femtocells. But building these femtocells infrastructure has a huge budget dependency. In [13] a solution to this massive problem has been given by creating FemtoHaul system architecture which efficiently uses relays in the femtocells for bearing the macrocell backhaul traffic and reinforced immense data rates for the cellular subscribers by serving more users with the extant macrocell backhaul capacity. In [14] a proposal has been made on a novel algorithm for creating a neighbour cell list during handover with a minimal but exact number of cells when there are dense femtocells and it is also proposed that CAC effectively handles various calls. In [15] three integrated network architectures have been introduced which has the ability to increase the access capacity by reducing the deployment and operational costs depending on interference management, efficient frequency, xDSL-based backhaul networks quality of service provision, and ingenious handover control issues to apply it in real life scenarios. Interference reduction is one of the main concerns in wireless communication in increasing system performance. These has been discussed in [16] proposing frequency reuse mechanisms which maximizes throughput utilizing different combinations of inner cell radius and allocating frequency depending on the position of the users and the femtocells. Considering random number of cognitive radio, path loss, Raleigh fading aggregate interference using Gamma distribution approximation to perfect closed-form moment generating function with an accurate approximation is derived in [17]. In [18] the research shows a presentation of a mathematical simulation on OFDMA or TDMA based femtocells depending on open and closed access methods according to mobile user density, which suggests OFDMA to be adaptive for the average cellular user connectivity.

B. Summary

There are different access methods, but among them the open access method has been chosen to find the probability of mobile user connectivity using femtocell which provides robust connectivity than the macrocell users reducing the cost effectiveness of the femtocells. Femtohaul has been used to reduce the traffic of macrocell backhaul.

III. DIFFERENT TYPES OF ACCESS MECHANISMS OF FEMTOCELLS AND MACROCELLS

In wireless communication, the users at one end communicate with the person at the other end by building a connection between them. This connection is mainly built by the “Access Points”. An access point is mainly a device, with which different wireless devices are connected to a network, e.g., a wireless router. Usually there are built-in routers in most cases in the access points, whereas others must have a connection to a router to serve network access. In each of the two cases the access points are hardwired to network switches or broadband modems. This access control mechanism can be divided into three types:

- A. Closed access method.
- B. Open access method.
- C. Hybrid access method.

A. Closed Access Method

In closed access method, there will be some authorized and unauthorized subscribers of the femtocells. Only the authorized subscribers would be able to access the femtocells, shown in Fig. 3.

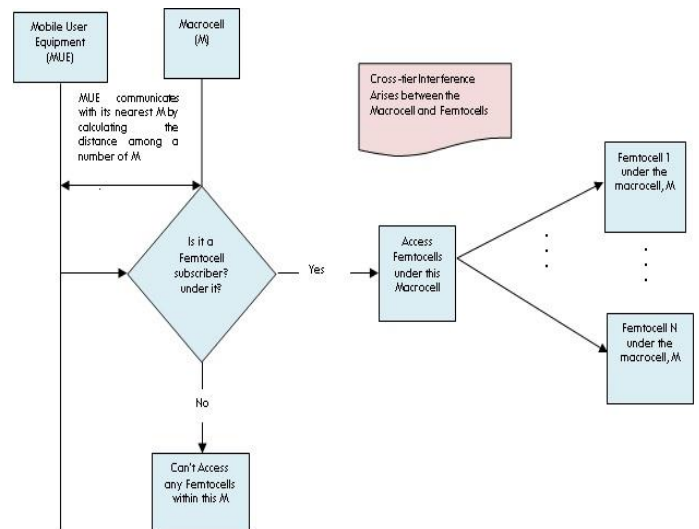


Fig. 3. Flowchart of closed access method.

B. Open Access Method

In open access method, the subscribers and nonsubscribers of the femtocells would be able to get the service on different condition which is shown in Fig. 4.

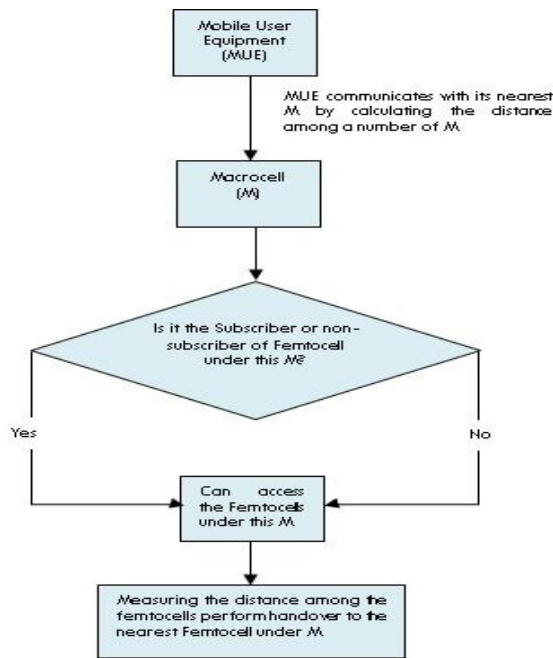


Fig. 4. Flowchart of open access method

C. Hybrid Access Method

In hybrid access method, the subscribers and limited number of nonsubscribers of the femtocells would be able to get the service on different conditions, as shown in Fig. 5.

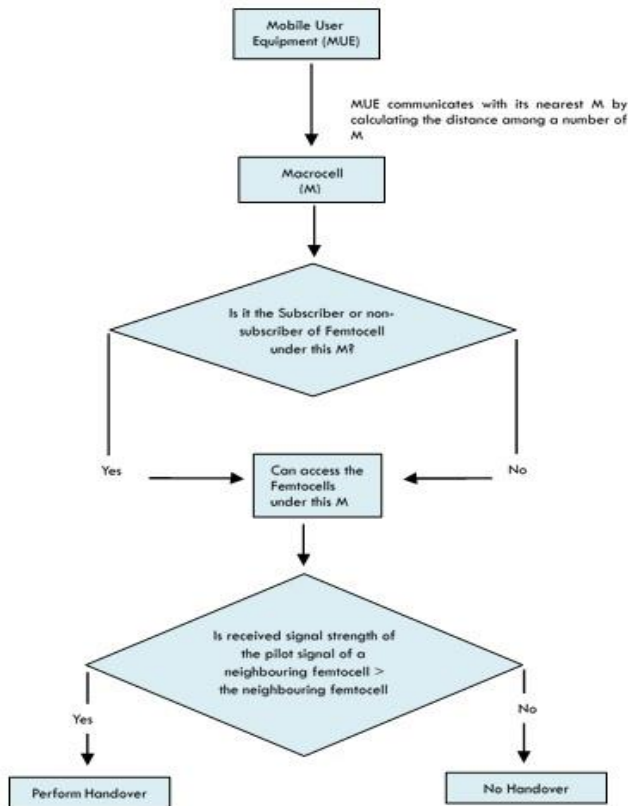


Fig. 5. Flowchart of hybrid access method.

IV. METHODOLOGY

A. System Model

The signal to noise plus interference ratio (SINR) of a Femto user at the cell boundary is given by (1):

$$SINR = \frac{P_f G_f R_f^{-\alpha}}{N_0 + I} \quad (1)$$

Where, α is the path loss exponent

and I , the total interference made by a Femto BS [3] is given by:

$$I = \sum_{Y_i \in \phi_f^p} G_{Y_i} P_f d^{-\alpha}$$

From (1) the advantage of femtocells can be easily realized. For smaller indoor coverages of femtocell, reduced distance between the femtocell and the user leads to higher signal strength. From the Shannon capacity formula,

$$C = B \log(1 + SINR) \quad (2)$$

The above formula indicates that the increment in signal strength and reduction in interference causes the progression in capacity of the signal. As the femtocell serves only around 2-4 users, it can assign a large portion of these resources (transmit power P_f and bandwidth B) to each subscriber. So, by deploying femtocells, more efficient usage of power and frequency resources can be enabled. As the femtocell is connected to the wired broadband network, the broadband operator provides sufficient QoS over the backhaul. Moreover, the femtocell networks can assist the areas where there is not that much signal coverage by degrading the comparative traffic volume on the macrocell.

Under open access method, a mathematical model for connectivity probability of mobile user in case of communication range as well as mobility factor has been introduced. Let us assume a test femtocell network is deployed within the communication range (r) of a macrocell base station and has allowed open access method in order to enhance the mobile users' connectivity by selecting the closest Femto Access Points. Let us consider, femtocell coverage to be a circle having unit radius. The distance (d) between the centres of femto access point and the macrocell user can be defined by the following equation [6]:

$$d = 1 + \beta r \quad (3)$$

Where, parameter β is called the mobility factor of macrocell user that plays an important role to avoid the users from disconnectivity. The area of intersection between a circle of unit radius and circle of radius r as follows [6]:

$$A(r, \beta) = \begin{cases} 0, & \beta \geq 1 \\ \alpha(r, \beta), & -1 < \beta < 1 \\ \pi r^2, & \beta \leq -1 \end{cases} \quad (4)$$

Where, [20]

$$a(r, \beta) = \frac{r^2}{\cos\left(\frac{d^2+r^2-R^2}{2dr}\right)} + \frac{R^2}{\cos(d^2+r^2-R^2)} - \frac{\sqrt{(-d+r+R)(d+r-R)(d+r+R)}}{2}$$

Considering $R=1$.

Let, there are N number of femtocells are overlaid in the macrocell mobile network. So, the probability of macrocell users are not able to connect to any femtocell among $N-1$ is given by [6],

$$\begin{aligned} \alpha(r, \beta) &= \left(1 - \frac{A(r, \beta)}{\pi r^2}\right)^{N-1} \\ &= \left(1 - \frac{A(r, \beta)}{\pi \cdot 1^2}\right)^{N-1} \\ &= \left(1 - \frac{A(r, \beta)}{\pi}\right)^{N-1} \end{aligned} \quad (5)$$

So, the probability of connectivity is obtained as [6]

$$P_c = 1 - \sum_{j=1}^N \left(1 - \frac{A(r, \beta)}{\pi}\right)^{N-1} \quad (6)$$

Since a consideration has been made on femtocell network having a scenario of open access method, the femtocells networks are highly and densely populated femtocells. A simplified wireless model neglecting fading has been considered.

Therefore, the probability of a macrocell user which is fully connected in femtocells [6]:

$$P_c = \frac{D_f}{D_u} \left[1 - e^{\left(1 - e^{-D_f \pi r^2}\right) \frac{D_u}{D_f}}\right] \quad (7)$$

Where, D_f is the density of femto access point and D_u is the density of mobile users.

The outage probability of a user under SINR constraint will be [6],

$$P(\text{SINR} \leq \Gamma_{th})$$

The above equation can be rewritten according to the active femto access points [6].

$$1 - P(\text{SINR} \geq \Gamma_{th}) = 1 - \frac{D_f \left[1 - e^{-\left(D_{f,active} \gamma^\alpha + D_f\right) \pi r^2}\right]}{\left(D_{f,active} \gamma^\alpha + D_f\right) \left[1 - e^{-D_f \pi r^2}\right]} \quad (8)$$

$$\text{Where, } D_{f,active} = P_c * D_u$$

Since channel model is distance dependent. So, considering path loss exponent, the probability of SINR greater or equal to the threshold SINR Γ_{th} [1],

$$\begin{aligned} P(\text{SINR} \geq \Gamma_{th}) &= P\left(\frac{P_f G_f R_f^{-\alpha}}{N_0 + I} \geq \Gamma_{th}\right) \\ &= P\left\{G_f \geq \frac{\Gamma_{th} (N_0 + I)}{P_f R_f^{-\alpha}}\right\} \\ &= e^{-\Gamma_{th} N_0 / P_f R_f^{-\alpha}} \cdot M_I\left(\frac{\Gamma_{th} I}{P_f R_f^{-\alpha}}\right) \end{aligned} \quad (9)$$

Where,

$$M_I(s) = \exp\left[-2\pi\lambda \int_0^\infty \frac{udu}{1 + \frac{u^\alpha}{sP_f}}\right] \quad [17]$$

So, the outage probability is

$$1 - P(\text{SINR} \geq \Gamma_{th}) = 1 - e^{\Gamma_{th} N_0 / P_f R_f^{-\alpha}} \cdot M_I\left(\frac{\Gamma_{th} I}{P_f R_f^{-\alpha}}\right) \quad (10)$$

An analytical work has been done to solve interference by the system models named as interference model, channel models using per tier outage probability considering its coverage. A spectrum division has been used by all proposed solutions. A specific section of the available spectrum needs to be reserved for the femtocells and the rest of them for the macrocells. By avoiding the interference problems, it provides an optimum solution. In this case, among the competing entities of the shared spectrum, one of the femtocells makes a random selection of the reserved spectrum to use a small

portion. Though the cross-tier interference is reduced potentially, femtocells interference still exists. Moreover, since the femtocells are smaller indoor coverage with 2 to 4 users for each femtocell, the subcarriers reserved for the femtocells remain idle most of the time. This is because users consume their data most of the time either at home (indoor) or at work (outdoor) or rest of the time may be at other places (outdoor). So, this solution causes waste of bandwidth and reduces the spectral efficiency. On the other hand, because of its expensiveness it creates disinterest. Even though, frequency bands sharing is technically more challenging, it is also more appropriate between femtocells and macrocells. The fractional re-use scheme that has been proposed may give a potential solution which is shown in Fig. 6 [19].

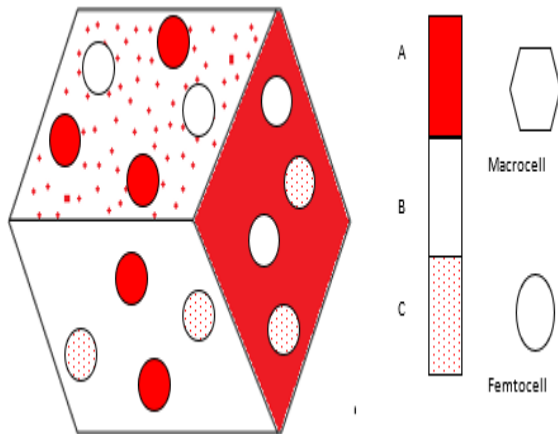


Fig. 6. A probable frequency allocation among femtocells in order to degrade interference.

The above proposed scheme could be used where the frequency re-use factor is static through macro cellular and femto cellular networks. In this case, fractional frequency re-use should be managed very carefully. Because a large number of femtocells might exist inside the macrocells, if any inappropriate frequency allocation among the femtocells has been made it might cause symbolic co-tier interference. Here a model for a small number of femto access points has been proposed giving a solution to co-tier interference and cross-tier interference.

In the proposed re-use scheme, at first the frequency band is segregated into three sub-bands A, B and C. Then, each macrocell is divided into three portions as sub-bands with respect to center and each of the sub-bands is apportioned to one portion of macrocell. The femtocell allocated in one portion of a macrocell uses one of the other two sub-bands in order to minimize the possibilities of co-tier interference. There is a suggestion for this type of fractional frequency re-use scheme that it can be utilized in SON-based femtocell architectonics. Because, it may reduce the interference by perceptive power optimization and frequency allocation.

When a femtocell is installed, the frequency allocations among the neighbouring femto access points (FAP) will auto-reconfigure to get different frequencies used by near femto access points in order to reduce the co-tier interference among the femtocells.

V. RESULTS AND DISCUSSION

In this paper, the connectivity and outage probability of mobile user under femtocell/macrocell networks has been derived. A performance measurement model setting different input parameters according to the real cellular networks has been proposed. Here, MathCAD has been used to obtain the numerical solution and graphical representations. The user connectivity probability given by (6) is the function of mobility factor (β), communication range (r) and number of femtocells (N). The proposed model is analyzed under three different number of femtocells $N=2$, $N=5$ and $N=100$.

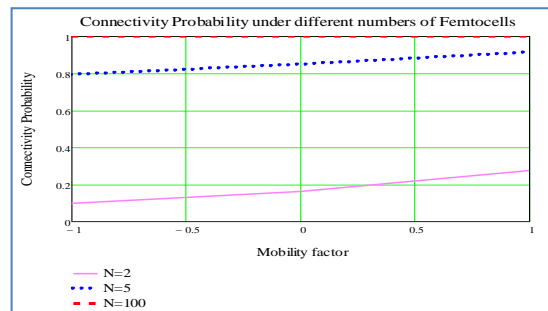


Fig. 7. Variation of connectivity probability against mobility factor.

Fig. 7 shows the probability of mobile user connectivity under different number of femtocells. Here, the mobility factor (β) has been considered first which indicates the movement of femto users. The mobility factor (β) within the range -1 to 1 has been considered to avoid the users from disconnectedness. This β decides whether the user is within the communication range or not. It is clear from the figure that the connectivity probability increases with the increase in communication range especially when the mobility factor (β) is varied from 1 towards -1. In contrast, it is shown that the connectivity probability decreases when the mobility factor is varied from -1 towards 1. Moreover, it indicates that the system's connectivity performance can be improved with the increase in number of femtocells. This is due to the fact that, when the macrocell is densely populated with femtocells, the probability of mobile users' connectivity increases.

Now, the most important part lies in the relation between connectivity probability and number of mobile users. The connectivity probability is heavily affected by the number of mobile users. Fig. 8 shows that the connectivity probability decreases with increase in the number of mobile users. Moreover, the connectivity probability is low for low density of active femto access points and high for high density of active femto access points which is shown in Fig. 8. The potential results for active femto access points with density, $D_f = 6$ has been obtained.

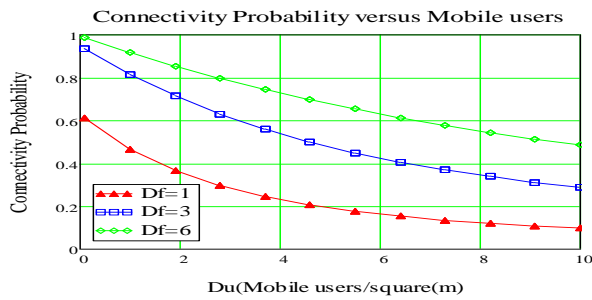


Fig. 8. Variation of connectivity probability against number of mobile.

Finally, Fig. 9 depicts the variation of the outage probability according to the mobile users per square meter under the consideration of the density of femto access points. The outage probability increases with the increase in mobile users and obtained lower value with active femto access points density, $D_f=6$. Therefore, there is a scope of reducing the outage probability by incrementing the connectivity probability beyond this level which is visualized in Fig. 8 and 9.

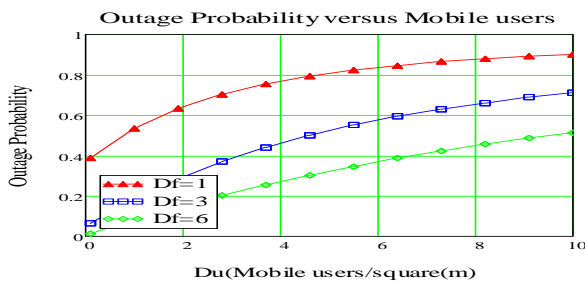


Fig. 9. Variation of outage probability against mobile users.

For Fig. 10, free space path loss exponent $\alpha = 2$ and other parameters $R_f = 15$ and $N_0 = 10^{-12}$ has been considered. The outage probability increases with the increase in density of femto BS and also with the increase in threshold SNR at receiving end.

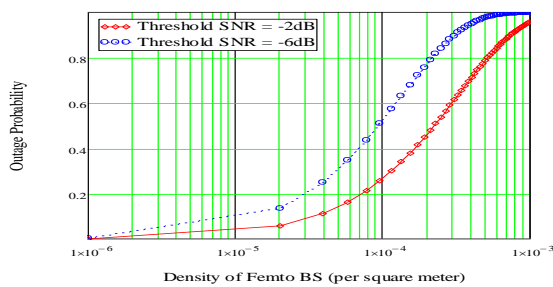


Fig. 10. Variation of outage probability against density of Femto BS taking threshold SINR as a parameter.

The Variation of Outage probability against the threshold of detection under various path loss exponents is shown in Fig. 11. It is clear from the figure that outage probability increases with the increase in threshold of detection or threshold SINR. The impact of path exponent on the outage probability is also cleared from the figure.

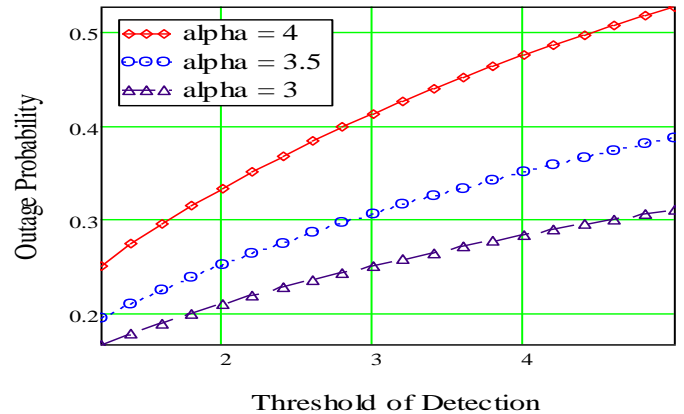


Fig. 11. Variation of outage probability against threshold of detection.

Finally, Fig. 12 shows the variation of Spectral efficiency against SINR. This indicates the maximum achievable spectral frequency through the AWGN channel under SINR.

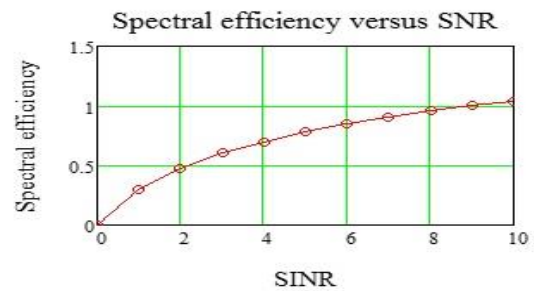


Fig. 12. Variation of spectral efficiency against SINR.

VI. CONCLUSIONS

Macrocell has a wide area of mobile network connectivity which degrades the performance with the increase in distance of the mobile users from the macrocell access point. This degradation of the connectivity probability problem can be better handled with femtocells under a macrocell. Since Open Access Method (OAM) works with both subscribed and non-subscribed femtocell users so in this paper a better mobile user connectivity probability has been proposed reducing the outage probability which was tested with different number of mobile users and different SINR threshold and showing those variations of spectral efficiency against Finally, a Fractional Frequency Reuse scheme for SON-based femtocell architecture has been proposed to optimize the interference considering the frequency allocation and power optimization for the mobile users providing improved connectivity probability.

ACKNOWLEDGMENT

The authors would like to express deep gratitude to their mentors for their impetuous guidance, enthusiastic encouragement and useful rebuke to make this research work. Lastly, the heartiest gratitude goes to their parents for continuous support.

REFERENCES

- [1] Guillaume de la Roche, Alvaro Valcarce, David L'opez-P'erez, Jie Zhang "Access Control Mechanisms for Femtocells", IEEE Communications Magazine, July 2009.
- [2] Shin-Ming Cheng, Weng Chon Ao, and Kwang-Cheng Chen, "Downlink Capacity of Two-tier Cognitive Femto Networks", 21st Annual IEEE international Symposium on Personal, Indoor and Mobile Radio Communications, Print ISBN: 978-1-4244-8017-3 2010, pp. 1301-1306.
- [3] Mahsa Derakhshani, Tho Le-Negoc, "Aggregate interference and Capacity-Outage Analysis in a Cognitive Radio Network", IEEE Transactions on Vehicular Technology, JAN 2012, Vol. 61, No. 1, pp.196 – 207.
- [4] V. Chandrasekhar and J. G. Andrews, "Femtocell networks: A survey," IEEE Commu. Mag., vol. 46, no. 9, pp. 59–67, Sept. 2008.
- [5] S. Yeh, S. Talwar, S. Lee, and H. Kim, "WiMAX femtocells: a perspective on network architecture, capacity, and coverage," IEEE Commu. Mag., vol. 46, no. 10, pp. 58–65, Oct. 2008.
- [6] Saied M.Abd El-atty and Z.M. Gharsseldien,"Analytical Model for Mobile User Connectivity in Coexisting Femtocells/Macrocells Networks", International Journal of Wireless & Mobile Networks (IJWMN) Vol. 4, No. 6, December 2012.
- [7] Kanak Raj Chaudhary, DeepeshRawat, EishaMadwal, "Interference Aware & SINR Estimation in Femtocell Networks", IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661, p-ISSN: 2278-8727 Volume 10, Issue 6 (May. - Jun. 2013), PP 64-69.
- [8] Gang Ning, Qinghai Yang, Kyung Sup Kwak, Hanzo, L., "Macro- and Femtocell Interference Mitigation in OFDMA Wireless
- [9] Nazmus Saquib, Ekram Hossain, Long Bao Le, and Dong in Kim "Interference Management in OFDMA Femtocell Networks: Issues and Approaches", Wireless Communications, IEEE (Volume:19, Issue: 3), ISSN:1536-128, pp. 86-95, June 2012.
- [10] HumbletP.Airvana, Raghothaman, B., Srinivas, A., Balasubramanian, S., Patel, C., Yavuz, M. "System design of CDMA2000 femtocells", Communications Magazine, IEEE (Volume:47, Issue: 9), ISSN: 0163-6804, pp. 92-100, September 2009.
- [11] PekkaPirinen "Co-channel co-existence study of outdoor macrocell and indoor femtocell users", Wireless Conference (EW), 2010 European, ISSN: 978-1-4244-5999-5 pp. 207 – 213.
- [12] Chang Seup Kim, Bum-Gon Choi, Ju Yong Lee, Tae-Jin Lee, HyunseungChoo, and Min Young Chung "Femtocell Deployment to Minimize Performance Degradation in Mobile WiMAX Systems".
- [13] AyaskantRath, Sha Hua and Shivendra S. Panwar "FemtoHaul: Using Femtocells with Relays to Increase Macrocell Backhaul Bandwidth", INFOCOM IEEE Conference on Computer Communications Workshops, 2010, pp. 1-5, March 2010.
- [14] Mostafa Zaman Chowdhury and Yeong Min Jang, "Handover management in high-dense femtocellular networks", EURASIP Journal on Wireless Communications and Networking, A Springer Open Journal, Vol. 2013, No. 6, pp. 1-21, 7 January, 2013.
- [15] Mostafa Zaman Chowdhury, Yeong Min Jang and Zygmunt J. Haas "Network Evolution and QOS Provisioning for Integrated Femtocell/Macrocell Networks", International Journal of Wireless & Mobile Networks (IJWMN), Vol.2, No.3, August 2010.
- [16] H Bouras, C., Kavourgiias, G., Kokkinos, V., Papazois, A., "Interference Management in LTE Femtocell Systems Using an Adaptive Frequency Reuse Scheme", Wireless Telecommunications Symposium (WTS), 2012, Print ISBN: 978-1-4577-0579-3, pp. 1-7.
- [17] Sachitha Kusaladharna, Chintha Tellambura, "Aggregate Interference Analysis for Underlay Cognitive Radio Networks",2012, Wireless Communications Letters, IEEE, Volume 1, Issue 6, pp.641 – 644.
- [18] Ping Xia, Vikram Chandrasekhar, Jeffrey G. Andrews "Femtocell Access Control in the TDMA/OFDMA Uplink", Global Telecommunications Conference (GLOBECOM 2010), IEEE, ISSN: 1930-529X, pp. 1-5, December 6-10, 2010.
- [19] Rizwan Ghaffar, Raymond Knopp "Fractional Frequency Reuse and Interference Suppression for OFDMA Networks", WIOPT 2010, 8th International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks, Avignon, France, 31 May-4 June 2010.
- [20] Wolfram Math World: <http://mathworld.wolfram.com/Circle-CircleIntersection.html>.

Heart Failure Prediction Models using Big Data Techniques

Heba F. Rammal

Information Technology Department
King Saud University
Riyadh, Saudi Arabia

Ahmed Z. Emam

Information Technology Department
King Saud University, Riyadh, Saudi Arabia
Computer Science and Math Department,
Menoufia University, Egypt

Abstract—Big Data technologies have a great potential in transforming healthcare, as they have revolutionized other industries. In addition to reducing the cost, they could save millions of lives and improve patient outcomes. Heart Failure (HF) is the leading death cause disease, both nationally and internally. The Social and individual burden of this disease can be reduced by its early detection. However, the signs and symptoms of HF in the early stages are not clear, so it is relatively difficult to prevent or predict it. The main objective of this research is to propose a model to predict patients with HF using a multi-structure dataset integrated from various resources. The underpinning of our proposed model relies on studying the current analytical techniques that support heart failure prediction, and then build an integrated model based on Big Data technologies using WEKA analytics tool. To achieve this, we extracted different important factors of heart failure from King Saud Medical City (KSUMC) system, Saudi Arabia, which are available in structured, semi-structured and unstructured format. Unfortunately, a lot of information is buried in unstructured data format. We applied some pre-processing techniques to enhance the parameters and integrate different data sources in Hadoop Distributed File System (HDFS) using distributed-WEKA-spark package. Then, we applied data-mining algorithms to discover patterns in the dataset to predict heart risks and causes. Finally, the analyzed report is stored and distributed to get the insight needed from the prediction. Our proposed model achieved an accuracy and Area under the Curve (AUC) of 93.75% and 94.3%, respectively.

Keywords—Big data; hadoop; healthcare; heart failure; prediction model

I. INTRODUCTION

In the recent years, a new hype has been introduced into the information technology field called 'Big Data'. Big Data offers an effective opportunity to manage and process massive amounts of data. A report by the International Data Corporation (IDC) [1] found that the volume of data the whole humanity produced in 2010 was around 1.2 Zettabytes, which can be illustrated physically by having 629.14 Million 2 Terabytes external hard drives that can fill more than 292 great pyramids. It has been said that 'data is the new oil', so it needs to be refined like the oil before it generates value. Using Big Data analytics, organizations can extract information out of massive, complex, interconnected, and varied datasets (both structured and unstructured) leading to valuable insights. Analytics can be done on big data using a new class of

technologies that includes Hadoop [2], R [3], and Weka [4]. These technologies form the core of an open source software framework that supports the processing of huge datasets. Like any other industry, healthcare has a huge demand to extract a value from data. A study by McKinsey [5] points out that the U.S. spends at least 600\$ - 850\$ billion on healthcare. The report points to the healthcare sector as a potential field where valuable insights are buried in structured, unstructured, or highly varied data sources that can now be leveraged through Big Data analytics. More specifically, the report predicts that if U.S. healthcare could use big data effectively, the hidden value from data in the sector could reach more than 300\$ billion every year. Also, according to the 'Big Data cure' published last March by MeriTalk [6], 59% of federal executives working in healthcare agencies indicated that their core mission would depend on Big Data within 5 years.

One area we can leverage in healthcare using Big Data analytics is Heart Failure (HF); HF is the leading cause of death globally. It is the heart's inability to pump a sufficient amount of blood to meet the needs of the body tissues [7]. Despite major improvements in the treatment of most cardiac disorders, HF remains the number one cause of death in the world and the most critical challenges facing the healthcare system today [8]. A 2015 update from the American Heart Association (AHA) [9] estimated that 17.3 million people die due to HF per year, with a significant rise in the number to reach 23.6 million by 2030. They also reported that the annual healthcare spending would reach \$320 billion, most of which is attributable to hospital care. According to World Health Organization (WHO) statistics [10], 42% of death in 2010 (42,000 deaths per 100,000) in the Kingdom of Saudi Arabia (KSA) were due to cardiovascular disease. Also, in KSA, cardiovascular diseases represent the third most common cause of hospital-based mortality second to accident and senility.

HF is a very heterogeneous and complex disease which is difficult to detect due to the variety of unusual signs and symptoms [11]. Some examples of HF risk factors are: breathing, dyspnea, fatigue, sleep difficulty, loss of appetite, coughing with phlegm or mucus foam, memory losses, hypertension, diabetes, hyperlipidemia, anemia, medication, smoking history and family history. Heart failure diagnosis is typically done based on doctor's intuition and experience rather than on rich data knowledge hidden in the database which may lead to late diagnosis of the disease. Thus, the effort to utilize

clinical data of patients collected in databases to facilitate the early diagnosis of HF patients is considered a challenging and valuable contribution to the healthcare sector. Early prediction avoids unwanted biases, errors and excessive medical costs, which improve quality of life and services provided to patients. It can identify patients who are at risk ahead of time and therefore manage them with simple interventions before they become chronic patients. Clinical data are available in the form of complex reports, patient's medical history, and electronics test results [12]. These medical reports are in the form of structured, semi-structured and unstructured data. There is no problem to use structured data for the prediction model. But, there is a lot of valuable information buried in the semi-structured and unstructured data format because those data are very discrete, complex, and noisy [13]. In our study, we collected patient's reports from a well-known hospital in Saudi Arabia: King Saud University Medical City (KSUMC). The objective of our research is to mine the useful information from these reports with the help of cardiologists and radiologist to design a predictive model that will give us the prediction of HF. The paper is organized as follows. Section II introduces the related work. Section III describes the proposed architectural model and each process involved. In Section IV, the proposed research methodology is explained. The conclusion and future work of this research are found in Section V.

II. LITERATURE REVIEW

Big Data predictive analytics represents a new approach to healthcare, so it does not yet have a large or significant footprint locally or internationally. To the best of our knowledge, no prior work has investigated the benefits of Big Data analytics techniques in heart failure prediction problem. A work by Zolfaghar K, et al. [14] proposed a real-time Big Data solution to predict the 30-day Risk of Readmission (RoR) for Congestive Heart Failure (CHF) incidents. The solution they proposed included both extraction and predictive modeling. Starting with the data extraction, they aggregate all needed clinical & social factors from different recourse and then integrated it back using a simple clustering technique based on some common features of the dataset. The predictive model for the RoR is formulated as a supervised learning problem, especially binary classification. They used the power of Mahout as machine learning based Big Data solution for the data analytics. To prove quality and scalability of the obtained solutions they conduct a comprehensive set of experiments and compare the resulted performance against baseline non-distributed, non-parallel, non-integrated dataset results previously published. Due to their negative impacts on healthcare systems' budgets and patient loads, RoR for CHF gained the interest of researchers. Thus, the development of predictive modeling solutions for risk prediction is extremely challenging. Prediction of RoR was addressed by, Vedomske et al. [15], Shah et al. [16], Royet al. [17], Koulaouzidis et al. [18], Tugerman et al. [19], and Kang et al. [20]. Although our studied problem is fundamentally different as they are all using structure data; nevertheless, our proposed model could benefit from the proposed large-scale data analysis solutions.

Panahiazar et al. [21] used a dataset of 5044 HF patients admitted to the Mayo Clinic from 1993 to 2013. They applied 5 training algorithms to the data that includes decision trees, Random Forests, Adaboost, SVM and logistic regression. 43 predictors were selected which express demographic data, vital measurements, lab results, medication, and co-morbidities. The class variable corresponded to survival period (1-year, 2-year, 5-year). 30% of the dataset were used for training and the rest 70% for testing. The authors observed that logistic regression and Random Forests were more accurate models compared to others, also among the scenarios, the best prediction accuracy was 87.12%.

Saqlain, M. et al. [22] worked on 500 HF patients from the Armed Forces Institute of Cardiology (AFIC), Pakistan, in the form of medical reports. They started by manually applying pre-processing steps to transform unstructured reports into the structured format to extract data features. Then they perform multinomial Naïve Bayes (NB) classification algorithm to build 1-year or more survival prediction model for HF diagnosed patients. The proposed model achieved an accuracy and Area under the Curve (AUC) of 86.7% and 92.4%, respectively. Even though the above model is based on some attributes extracted from the unstructured data, they used a manual approach to achieve this.

On the other hand, our model deals with unstructured data by automatically recognizing attributes using Machine Learning (ML) approaches without the need for a radiologist opinion. A scoring model for HF diagnosis based on SVM was proposed by Yang, G. et al. [23]. They applied it to a total of 289 samples clinical data collected from Zhejiang Hospital. The sample was classified into three groups: healthy group, HF-prone group, and HF group. They compared their results to previous studies which showed a considerable improvement in HF diagnosis with a total accuracy of 74.44%. Especially in HF-prone group, accuracy reaches 87.5%, and this implies that the proposed model is feasible for early diagnosis of HF. However, accuracy in the HF group is not so satisfied due to the absence of symptoms and signs and also due to the high prevalence of conditions that may mimic the symptoms and signs of heart failure.

More studies were listed in Table I, which was collected and summarized as recent analytics techniques and platform to predict heart failure. The table shows that supervised learning technique is the most dominant techniques in building HF prediction model, also Weka and Matlab are the preferable platforms to build HF prediction model.

The literature presented above shows a gap in multi-structured predictors for HF prediction and data fusion which will be our main task. It is easy to observe that our effort is orthogonal to this related work but, unlike us, none of these works deal with the problem semi-structured or unstructured

HF predictor variable. They did not generate Big Data analytics prediction model, nor do they perform on large scale or distributed data.

TABLE I. STATE OF ART FOR HF PREDICTION STUDIES

Author	Prediction Technique Used	Platform	Objective
Zolfaghar K, et al (2013)	Logistic regression, Random forest	Mahout	BD solution to predict the 30- day RoR of HF
Meadam N., et al (2013)	Logistic regression, Naive Bayes, Support Vector Machines	R	Evaluation preprocessing techniques for Prediction of RoR for CHF Patients
Yang, G. et al (2010)	support vector machine (SVM)	n/a	A heart failure diagnosis model based on support vector machine
Panahiazar et al. (2015)	Decision trees, Random Forests, Adaboost, SVM and logistic regression	n/a	Using EHRs and Machine Learning for Heart Failure Survival Analysis
Donzé, Jacques et al (2013)	Cox proportional hazards	SAS	Avoidable 30-Day RoR of HF
K. Zolfaghar et al (2013)	Naive Bayes classifiers	R	Intelligent clinical RoR of HF calculator
Bian, Yuan et al (2015)	Binary logistic regression	n/a	Scoring system for the prevention of acute HF
Suzuki, Shinya et al (2012)	logistic regression	SPSS	Scoring system for evaluating the risk of HF
Auble, T. E. et al (2005)	Decision tree	SPSS	Predict low-risk patients with HF
Pocock, S. J. et al (2005)	Cox proportional hazards	n/a	Predictors of Mortality and Morbidity in patients with CHF
Miao, Fen et al (2014)	Cox proportional hazards	R	Prediction for HF incidence within 1-year
S.Dangare et al (2012)	Decision Trees, Naïve Bayes, and Neural Networks	Weka	HD prediction system using DM classification techniques
Rupali R. Patil (2014)	Naive Bayes classifiers	MATLAB	HD prediction system
Rupali R. Patil (2012)	Artificial Neural network	Weka	A DM approach for prediction of HD
Wu, Jionglin et al (2010)	Logistic regression, SVM, and Boosting	SAS, R	HF prediction modeling using EHR
Zebardast, B. et al (2013)	Generalized Regression Neural Networks	MATLAB	Diagnosing HD
Vanisree K. & Singaraju J. (2011)	Multi-layered Neural Network	MATLAB	Decision Support System for CHD Diagnosis
Guru N. et al. (2007)	Neural network	MATLAB	HD prediction system
R, Chitra and V, Seenivasagam (2013)	Cascaded Neural Network	n/a	HD Prediction System
Sellappan Palaniappan and Rafiah Awang (2008)	Decision trees, naïve bayed and neural network	.Net	HD prediction system using DM techniques
K. Srinivas et al (2010)	Naive Bayes classifiers	Weka	DM technique for prediction of Heart Attacks
Saqlain, M. et al (2016)	Naive Bayes	n/a	Identification of HF using unstructured data of Cardiac Patients
Strove, Sigurd et al (2004)	Structured prediction (Bayesian network)	HUGIN	Decision Support Tools in Systolic HF Management
Gladence, L.M. et al (2014)		Weka	Method for detecting CHF
Liu, Rui et al (2014)		MicrosoftAzure (R & python)	Framework to recommend interventions for 30-Day RoR of HF
C. Ordonez (2006)	Association rules	n/a	HD Prediction
M. Akhil Jabbar et al. (2012)	Associative classification	Gini index, Z-statics & genetics algorithm	Decision Support System for HD prediction
K. Chandra Shekar et al (2012)		association rule mining and classification	Java

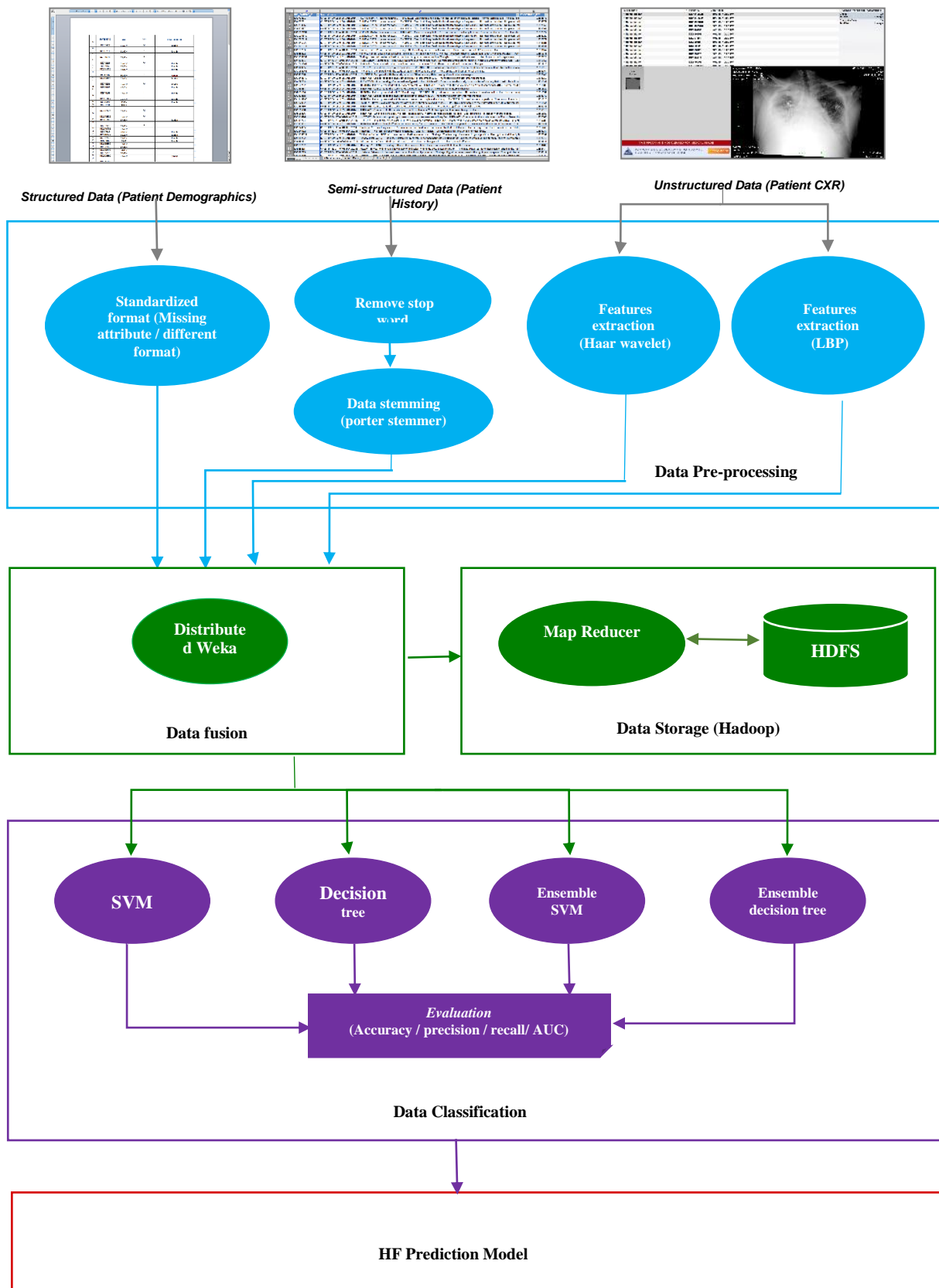


Fig. 1. HF prediction model.

III. PROPOSED ARCHITECTURE

Predictive analysis can help healthcare providers accurately expect and respond to the patient needs. It provides the ability to make financial and clinical decisions based on predictions made by the system. Building the predictive analysis model includes various phases as mentioned in the literature (Fig. 1 shows the complete architecture of proposed model).

- Layer 1: Data collection from KSUMC in the form of structure, unstructured, and semi-structured.
- Layer 2: Data pre-processing to prepare and filter the dataset to make it ready for the next step in building the model.
- Layer 3: Data fusion and storage which is an important layer that used to integrate all preprocessed data and store it in HDFS to be then fed to the next step.
- Layer 4: Data classification and evaluation are the two final steps that include training, testing then evaluating the model.

IV. PROPOSED METHODOLOGY

In the following, we will describe the adapt methodology and each step in toward our proposed model.

A. Data Collection

In our study, we collaborated with King Saud University Medical City (KSUMC) system located in Riyadh, Saudi Arabia to extract manually all needed clinical and demographic that we needed to adapt to evaluate the performance of the proposed model in identifying HF risk, from January 2015 to December 2015.

The dataset contained 100 real patient records extracted form KSUMC Electronic Health Record (EHR) and Picture Archiving Communication System (PACS), with approval from KSUMC administrative office. Due to patients' privacy, some demographic information that includes name, national ID number or iqama number, phone, address were excluded. Basic characteristics of the samples' demographic information are shown in Table II. Obviously, our sample doesn't have a uniform distribution in terms of gender. Also, patients aged from 60 years old to 70 years old account for the most part of our data. One of the major steps is the distillation of data, which responsible of determining the subset of attributes (i.e., predictor variables) that has a significant impact in predicting patient with HF from the myriad of attributes present in the dataset. In this study, parameters are selected from 3 datasets which are summarized in Table III.

The validation of the selected dataset achieved by consolidating some cardiologist and according to their evaluation all cases were labeled into two groups. The selected dataset has many noises such as missing values and misidentified attributes. The output values were categorized into two labels denoted as Non-HF (meaning HF is absent) and HF (meaning HF is present). Our dataset contains 69 predictor variable, having 1 binary variables (gender) 3 text values (place of birth, history, and symptoms) and 65 numerical variables (including age and all CXR features) and a single response variable 'Result' having only two values HF and Non-HF.

B. Data Preprocessing

In this phase structured, semi-structured, and unstructured data are accumulated, cleansed, prepared, and made ready for further processing.

TABLE II. DEMOGRAPHIC BASIC CHARACTERISTICS

Characteristic	Group		
	HF group	Non-HF group	Total
Female	21	23	44
Male	25	31	56
Age (mean ± SD)	69 ± 12	61 ± 15	65 ± 14

TABLE III. SELECTED ATTRIBUTES FROM THE DATASET

	Label	Feature	Format
Structured	Demographics	Age	Numeral
		Sex	Binary
		Place of birth	Nominal
Semi-Structured	Clinical indications / History	Hypertension, Anemia, Diabetes, Chronic Kidney Disease, Ischemic heart disease, SOB, Swilling hands, Cough, Previous CHF	String
Un-Structured	Front CXR Back CXR Side CXR	64 Features (Haar)	Numeral
		61 Features (LBP)	

- Raw structured information has some missing values and written in different formats during information entry or management. Those data with too many missing attributes were all wiped off when we selected the dataset. Also, all data formats were standardized, see Table IV.
- Apply text analysis techniques on the semi-structured dataset to get the needed information. Three steps were applied to the text to process the data, tokenizer, stop word removal, and stemming. Before any real text processing is to be done, the text needs to be segmented into words, punctuation, phrases, symbols, and other meaningful elements called tokens. Next, stop word removal, illustrated in Fig. 2, helped in removing all common words, such as 'a' and 'the' from the text. Then, Porter algorithm was used as the stemmer to identify and remove the commoner morphological and inflexional endings from words, which is part of the snowball stemmers in WEKA [24].

TABLE IV. UNSTANDARDIZED STRUCTURED DATA

	Age	Sex	P_B	Diagnosis
1	045Y	Female	Riyadh	HF
2	62	F	?	HF
.....
100	098	male	riy	Non-HF

Explanatory Data

Label

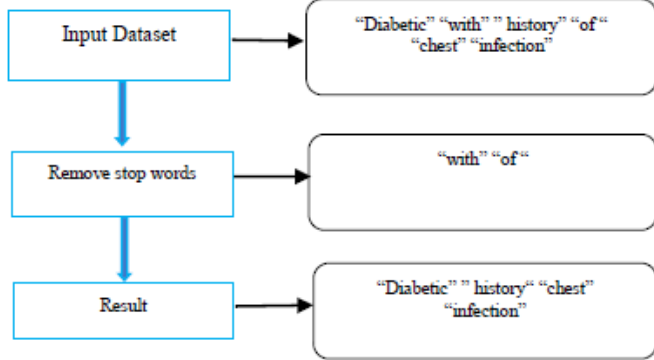


Fig. 2. Stop word removal.

- Extracting all needed features from the unstructured dataset, which includes 3 types of Chest X-Ray (CXR) images (front CXR, back CXR, and side CXR) using MATLAB. Haar wavelet and local binary pattern (LBP) were applied to over 150 CXR images. Haar was used since it is the fastest technique that can be used to calculate the feature vector [25]. This was performed based on applying the Haar wavelet 4 times to divide the input image into 16 sub-images, illustrated in Fig. 3. 64 features that include Energy, Entropy, Homogenous were found, Fig. 4 illustrate the resulted images after first level of Haar. Each CXR image represents certain features in the image of heart values of Energy_Entropy_Homo and wavelet features. A total of 16 for Energy_Entropy_Homo in each level since we have 4 levels in Haar wavelet so 64 features are extracted in total. On the other hand, LBP has been found to be a powerful and simple feature yet very efficient texture operator which labels the pixels of an image by calculating each pixels' neighborhoods' thresholding then considers the result as a binary number. We applied LBP to all CXR images by first, labeling all the pixels, absent the borders, using the LBP operator, then dividing the image into 60 segments. A feature vector is created by obtaining the histogram of each region, and finally concatenating all the histograms into one vector which result in finding 60 features. A typical LBP application to a CXR is shown in Fig. 5.

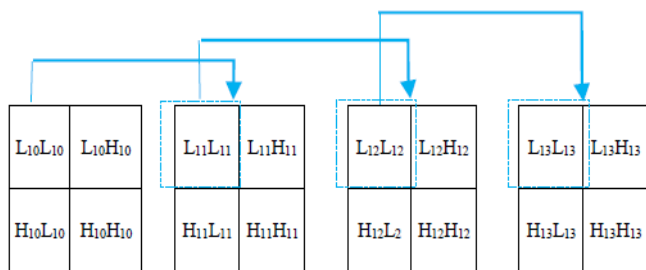


Fig. 3. Applying haar wavelet four times.

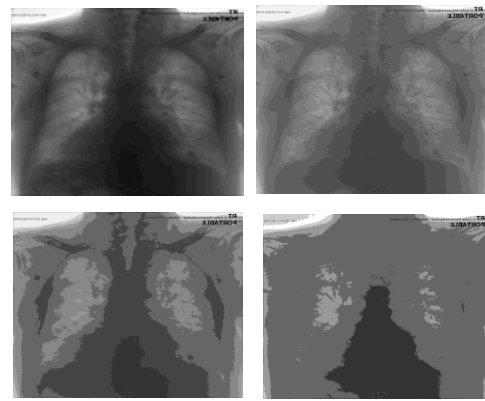
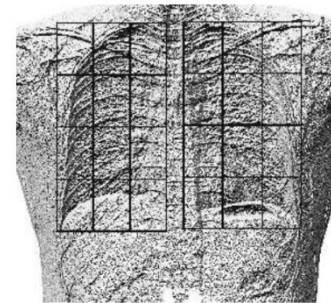
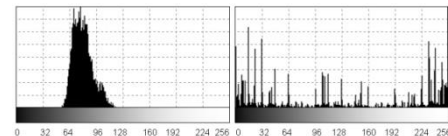


Fig. 4. The result from wavelet.



(a)



(b)

(c)

Fig. 5. Histogram of image obtained by applied LBP algorithm on CXR: (a) LBP applied image and (c) histogram before LBP (b) histogram of a.

- Principle component analysis (PCA) was applied to properly rank and compute the weights of the features to find the most promising attributes to predict HF from the features found. The selected attributes were used to train the classifiers to get a better accuracy. Also, to circumvent the imbalanced problem, we applied resampling method. This method alters the class distribution of the training data so that both the classes are well represented. It works by resampling the rare class records so that the resulting training set has an equal number of records for each.

C. Data Storage and Fusion

After pre-processing the data and extracting all the needed attributes, the statistics feature from CXR scan images with other attributes will be integrated using data fusion techniques to generate the needed data that will be used for training and testing and finally produce the predictive model. Complementary data fusion classification technique was used as each dataset represents part of the scene and was used to build a reliable information. We leverage the power of Hadoop as a framework for distributed data processing and storage. Hadoop is not a database, so it lacks functionality that is

necessary for many analytics scenarios. Fortunately, there are many options available for extending Hadoop to support complex analytics, including real-time predictive models such as Weka (Waikato Environment for Knowledge Analysis), which we used in our study. We added distributed WekaSpark to Weka' which works as a Hadoop wrapper for Weka.

D. Data Classification

In this study, each set of the data (Structured, Semi-structured, and Unstructured) trained and tested using data mining algorithms in Weka. Knowledge flow was used in Weka which presents, a workflow inspired interface, see Fig. 6. Data was trained using two state-of-the-art classification algorithms including, Random Forest (RF) and Logistic Regression (LR) as they both have been known to result in high accuracy in binary class prediction.

The ability to handle and analyze various types of data (structured, semi-structured or unstructured) is one of the most important characteristics of Big Data analytic techniques. We will perform the classifications in two phases to show that using the proposed integrated learning analytics technique is more efficient than a traditional single predictive model, especially if the data is multi-structured and has unique characteristics. In the end, model quality was assessed through common model quality measures such as accuracy, precision, recall and, Area under the Curve (AUC). Depending on the final goal of the HF prediction, the different evaluation measures are less or more appropriate. Recall is relevant as the detection of patients that belong to HF class is the main goal. The precision is considered less important as cost related to falsely predicting patients to belong to the class HF is low. The accuracy is the traditional evaluation measure that gives a global insight in the performance of the model. The AUC

measure is typically interesting in our study because the problem is imbalanced. It is observed that the number of instances with HF label significantly outnumbers the number of instances with class label Non-HF.

V. RESULTS AND INSIGHTS

It is clear from Tables V and VI that integrated dataset has the highest accuracy and AUC: ~92% and ~90%, respectively, then using each dataset by its own for HF patient's prediction. Also, using LBP features extraction methods achieved better performance results than Haar with 93% compared to 91% for Haar. We can also note that logistic regression did great in the integrated models compared to its poor performance in the single dataset models with over 90% recall, which can be resulted from the nature of the algorithm as it predicts better for problems with many attributes. Based on the experiment, we can provide evidence of the importance of the integration of unstructured, semi-structured, and structured data. This indicates that there are some indicators within textual patient report and images that can be extracted and used as important predictors of Heart failure. Also, the discovery of feature selection as a suite of methods that can increase model accuracy, decrease model training time and reduce overfitting.

Our proposed approach is also very important because it provides a knowledge discovery and intelligent model to the cardiologists and researchers such as: (1) the dataset contains 56% male patients, which mean male patients have more probability to get an HF diagnose than females. (2) 73% patient's age over 65 which indicate that aged people have more chance to get HF. (3) 70% of patients coming for HF complain having hypertension, diabetes, and SOB which means this disease has the main impact of HF.

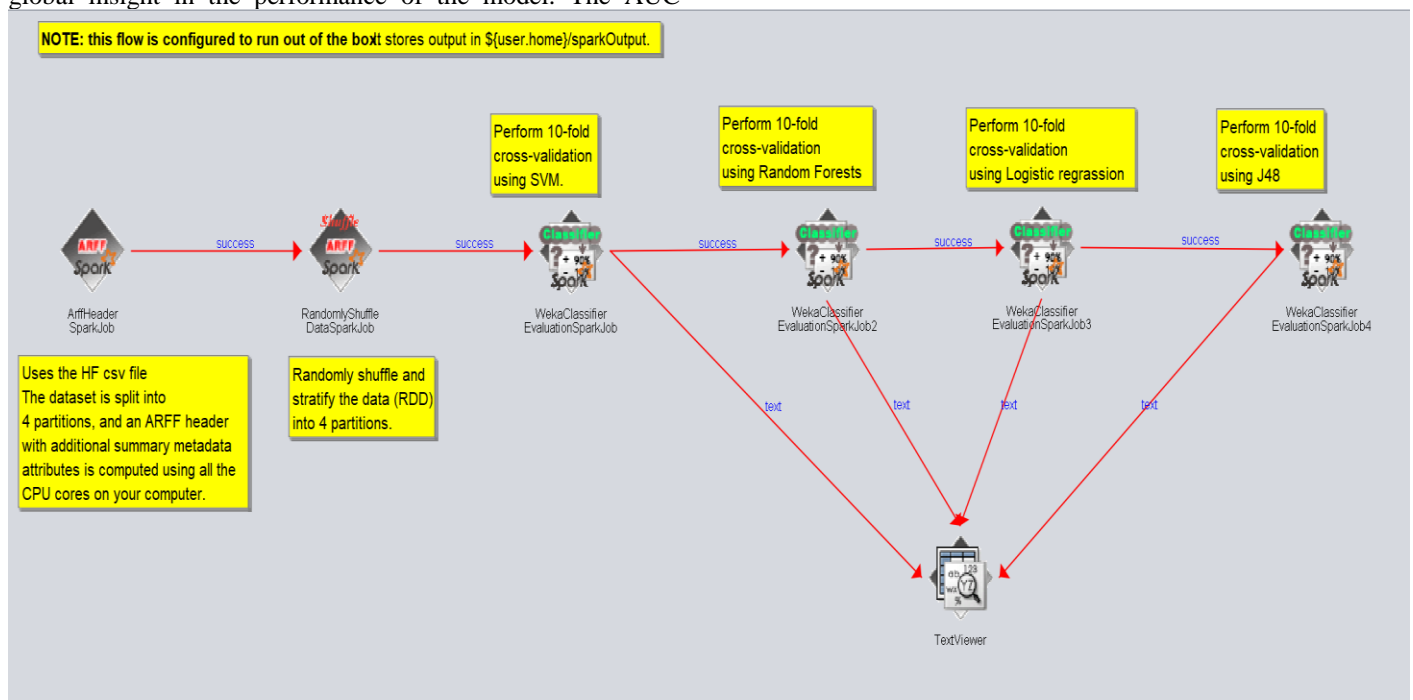


Fig. 6. The Proposed Knowledge flow using distributed Weka.

TABLE V. A PERFORMANCE MEASURE BASED ON RANDOM FOREST CLASSIFICATION ALGORITHM

	Precision %	Recall %	AUC %	Accuracy %
Structured	84.8	76.7	84.2	76.6
Semi-Structured	85.5	79.3	97.6	79.3
Un-Structured (Haar)	80.5	78.3	93.9	78.2
Un-Structured (LBP)	85.5	80	88.4	80
Integrated Dataset (Haar)	88.9	87.5	90	87.5
Integrated Dataset (LBP)	94.3	93.3	94.2	93.3

TABLE VI. A PERFORMANCE MEASURE BASED ON LOGISTIC REGRESSION CLASSIFICATION ALGORITHM

	Precision %	Recall %	AUC %	Accuracy %
Structured	79.2	60	78.3	60
Semi-Structured	40.1	41.3	70.5	41.3
Un-Structured (Haar)	80.5	78.3	93.9	78.2
Un-Structured (LBP)	76.1	56.7	66.5	56.6
Integrated Dataset (Haar)	91.7	91.7	80.3	91.6
Integrated Dataset (LBP)	93.3	93.3	94.3	93.3

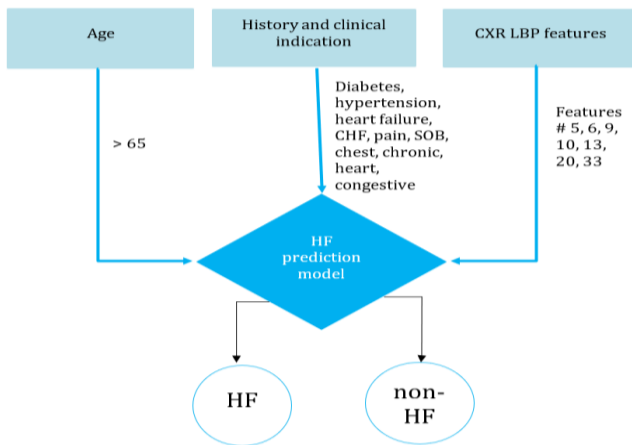


Fig. 7. HF integrated model.

The results of this study found 21 predictors, where the most powerful ones were older age, diabetes, hypertension and LBP features. It explains that if a patient came to the hospital having age > 65, suffering from Diabetes, hypertension, SOB and has some key features in the CXR, then there will be a high chance of having HF. Fig. 7 illustrated the proposed integrated model with the most promising attributes in all dataset. LBP features were used as LBP based model achieved better recall than Haar as mentioned in the previous sections. This also applied to the age range and words selected from the semi-structured data.

VI. CONCLUSION AND FUTURE WORK

Big Data Analytics provides a systematic way for achieving better outcomes of healthcare service. Non-Communicable Diseases like Heart Failure is one of a major health problems internationally. By transforming various health records of HF patients to useful analyzed result, this analysis will make the patient understand the complications to occur. The literature shows a gap in multi-structured predictors for HF prediction and data fusion which is our main task. It is easy to observe that our effort is orthogonal to this related work but, unlike us, none of these works deal with the problem semi-structured or unstructured HF predictor variable. Combining several characteristics from each patient demographical information, patient clinical information, and patient's Chest X-Ray is a very hard task. In this research, data fusion played a vital role in combining multi-structure dataset. We extracted different important factors of heart failure from King Saud Medical City (KSUMC) system. The extracted data were in the form of structured (patients demographics), semi-structured (patient history and clinical indication), and un-structured (patient chest X-Ray) data. Then we applied some preprocessing techniques to enhance the parameters of each dataset. After that, data was stored in HDFS to be trained and tested using different modeling algorithms on two phases to compare the performance measures of the resulted models before and after integrating them in the first phase we train each dataset as a traditional single predictive. Then, we integrated the most promising attributes from all dataset in the second phase and build 2 models based on Haar and LBP feature extraction. The results showed that the performances of the classifiers were better using the fused data (~93 % accuracy). For further improving, other intelligent algorithms need to be prospectively analyzed as well and more subjects should be investigated to keep upgrading the classifier. We will also incorporate more medical data into the model, better simulating how a cardiologist makes a decision.

REFERENCES

- [1] J. Gantz, and D. Reinsel, "The digital universe decade – are you ready?" External Publication of IDC (Analyse the Future) information and data, pp. 1- 16, 2010.
- [2] Hadoop Apache. Available at <http://hadoop.apache.org/>, Last accessed March 2017.
- [3] R project. Available at <https://www.r-project.org/about.html>, Last accessed March 2017.
- [4] P. Navas, Y. Parra, and J. Molano, "Big Data Tools: Hadoop, MongoDB and Weka". International Conference on Data Mining and Big Data, pp 449-456, 2017.
- [5] McKinsey and Company, McKinsey Global Institute, Big Data: The next frontier for innovation, competition, and productivity. Available at http://lateralpraxis.com/download/The_big_data_revolution_in_healthcare.pdf, Last accessed March 2017.
- [6] Meritalk, The Big Data cure. Available at <http://www.meritalk.com/bigdatacure>, Last accessed March 2017.
- [7] National Heart, Lung, and Blood Institute, What is heart failure. Available at <http://www.nhlbi.nih.gov/health/health-topics/topics/hf>, Last accessed April 2017.
- [8] World Health Organization (WHO) (2015). Cardiovascular diseases (CVDs). Available at <http://www.who.int/mediacentre/factsheets/fs317/en/>, Last accessed March 2017.
- [9] American Heart Association. Heart Disease and Stroke Statistics – At-a-Glance. Available at <http://www.heart.org/idc/groups/ahamah>

- public/@wcm/@sop/@smd/documents/downloadable/ucm_470704.pdf ,
Last accessed March 2017.
- [10] Mistry Of Health (MOH), “Cardiovascular Diseases Cause 42% of Non-Communicable Diseases Deaths in the Kingdom”. Available at <https://www.moh.gov.sa/en/Ministry/MediaCenter/News/Pages/News-2013-10-30-002.aspx>, Last accessed March 2018.
- [11] Ishwarappa, and J. Anuradha, “A Brief Introduction On Big Data 5Vs Characteristics And Hadoop Technology”. *Procedia Computer Science* 48, pp. 319-324, 2015.
- [12] J. A. Hartigan and M. A. Wong, “Algorithm AS 136: A k-means clustering algorithm,”, *Series C (Applied Statistics)*, vol. 28, pp. 100-108, 1979.
- [13] E. AbuKhoua, and P. Campbell, “Predictive data mining to support clinical decisions: An overview of heart disease prediction systems,” *Proc. IEEE, Innovations Information Technology (IIT)*, pp. 267-272, March 2012.
- [14] K. Zolfaghar, et al., “Big Data Solutions for Predicting Risk-of-Readmission for Congestive Heart Failure Patients”, *IEEE Inter. Conf. on Big Data*, 2013.
- [15] M. A. Vedomske, D. E. Brown, and J. H. Harrison, “Random forests on ubiquitous data for heart failure 30-day readmissions prediction”, *Proceedings of the 12th international conference on machine learning and applications*, vol. 2, pp. 415-421, 2013.
- [16] S. J. Shah, et al. “Phenomapping for novel classification of heart failure with preserved ejection fraction”. *Circulation*. Vol 131, pp. 269–279, 2015.
- [17] S. B. Roy, A. Teredesai, Zolfaghar K., Liu R., and Hazel D., “Dynamic hierarchical classification for patient risk-of-readmission”. *Proceedings of the 21st ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1691–1700, 2015.
- [18] G. Koulaouzidis, D. K. Iakovidis, and A. L. Clark, “Telemonitoring predicts in advance heart failure admissions”. *Int J Cardiol*, vol. 216, pp. 78–84, 2016.
- [19] L. Turgeman, J. H. May, “A mixed-ensemble model for hospital readmission.”, *Artif Intell Med*, vol. 72, pp. 72–82, 2016.
- [20] Y. Kang, M.D. McHugh, J. Chittams, K.H. Bowles, “Utilizing home healthcare electronic health records for telehomecare patients with heart failure. A decision tree approach to detect associations with rehospitalizations”. *Comput Inform Nurs*, vol. 34 no. 4, pp.175–182, 2016.
- [21] M. Panahiazar, V. Taslimitehrani, N. Pereira, J. Pathak, “Using EHRs and machine learning for heart failure survival analysis.” *Stud Health Technol Inform*, vol. 216, pp. 40–44, 2015.
- [22] M. Saqlain, W. Hussain, N. Saqib, A. Muazzam Khan, “Identification of Heart Failure by Using Unstructured Data of Cardiac Patients.”, *45th International Conference on Parallel Processing Workshops*, 2016.
- [23] G. Yang, et al., “A heart failure diagnosis model based on support vector machine”. *3rd International Conference on Biomedical Engineering and Informatics*. 2015;
- [24] C. Moral, A. Antonio, R. Imbert, J. Ramirez, “A survey of stemming algorithms in information retrieval.” *Information Research*, vol. 19 no. 1, March 2014.
- [25] S. Arora, Y. Brar, S. Kumar, “HAAR wavelet transform for solution of image retrieval.” *International Journal of Advanced Computer and Mathematical Sciences ISSN*, vol. 5, no. 2, pp 27-3, 2014.

Towards Privacy Preserving Commutative Encryption-Based Matchmaking in Mobile Social Network

Fizza Abbas¹, Ubaidullah Rajput¹, Adnan Manzoor², Imtiaz Ali Halepoto¹, Ayaz Hussain³

¹Department of Computer Systems Engineering, Quaid e Awam UEST, Nawabshah, Pakistan

²Department of Information Technology, Quaid e Awam UEST Nawabshah, Pakistan

³Department of Electrical Engineering, Balochistan University of Engineering and Technology, Khuzdar, Pakistan

Abstract—The last decade or so has witnessed a sharp rise in the growth of mobile devices. These mobile devices and wireless communication technologies enable people around the globe to instantaneously communicate with each other. This leads to the emergence of a new type of social networking known as Mobile Social Network (MSN). MSN offers a wide range of useful applications, such as group text services, social gaming, location-based services (to name a few). One of the popular applications of MSN is matchmaking where people match their interests/hobbies to find the like-minded people for a possible friendship. However, revealing personal hobbies can pose significant threats on a user's privacy. Therefore, a privacy preserving evaluation method is needed to find the similarity between users' interests. There are various techniques to achieve privacy preserving matchmaking, such as commutative encryption, oblivious transfer and homomorphic encryption. This paper discusses the feasibility of commutative encryption by evaluating recently proposed schemes. The paper attempts to identify various shortcomings in the present work and discusses future directions.

Keywords—Privacy; security; matchmaking; interests; mobile social network

I. INTRODUCTION

A Mobile Social Network (MSN) enables its users to make social ties between them using mobile devices and communication technologies [1]. MSN offers many useful applications such as locations-based services where nearby people share their experiences about restaurants, shopping malls, and social gaming that allows friends to play online games with each other (to name a few). One of the most popular applications of MSN is matchmaking where people find the similarity between their profiles to establish a possible friendship. Peoples' profiles consist of personal information such as political affiliations, sexual orientation and health status etc. Disclosure of such information to a stalker may seriously jeopardize the privacy of a user. In recent past, many researchers have proposed privacy preserving matchmaking schemes to privately evaluate the interest-wise similarity between their profiles. We can classify these techniques as a private set intersection (PSI) or private cardinality set intersection (PCSI) problem [2], [9], [10]. These techniques take their notion from the set theory where intersection operation is used to find the common elements in the sets. Here, private set intersection refers to the oblivious evaluation of intersection operation. There are other techniques to blindly

calculate the similarity such as cosine similarity can be used with the help of homomorphic encryption which incurs significant communication and computation costs. Moreover, it does not find interest to interest matching rather it calculates a similarity score [3], [4].

The remaining paper is organized as follows. The succeeding section discusses commutative encryption-based matchmaking protocols and their limitations. Section 3 provides the discussion. Section 4 concludes the paper along with future work.

II. MATCHMAKING PROTOCOLS

In this section, we discuss various commutative encryption-based protocols. Notations used in this paper are shown in Table I.

A. Agrawal et al. Protocol

Agrawal et al. presented the pioneer work regarding the commutative encryption. Originally, their work was intended to information sharing in between private databases [5]. They formulated their matching problem as a PSI problem. In case, the evaluation only finds the number of matches, the problem becomes PCSI. PSI and PCSI find the similar objects blindly [9], [10].

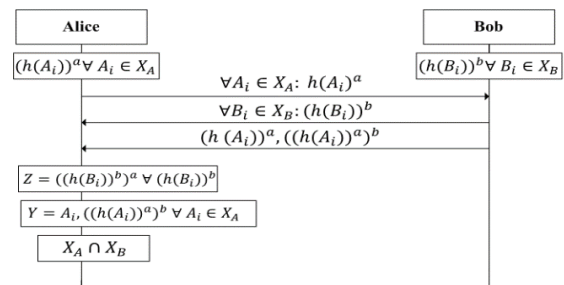


Fig. 1. Working of Agrawal et al. protocol [5].

The protocol proposed by Agrawal et al. uses the power function $f_e(x) = x^e \bmod p$ that has commutative properties i.e. the order of encryption is independent. Therefore, its security is based on Decisional Diffie-Hellman hypothesis (DDH). Suppose that a is the secret key of Alice, b is the secret key of Bob and m is the message then $((msg1)^a)^b = ((msg2)^b)^a$, iff $msg1 == msg2$.

TABLE I. NOTATIONS

Notations	Explanation
Alice	The protocol initiator
Bob	Responder
i	Index
X_A	Set containing Alice interests
X_B	Set containing Bob interests
A_i	i th interest Alice
B_i	i th interest Bob
a, b	Secret key of Alice and Bob respectively
h	Hash
PK_A, PK_B	Public key of Alice and Bob respectively
SK_A, SK_B	Secret key of Alice and Bob respectively
SK_{PIS}	Private key of third party server
SK_{VS}, SK_{IDV}	

A user's profile, consists of i number of interests require i number of modular exponentiations. The working of Agrawal's protocol is shown in Fig. 1. In this figure, we have two users Alice and Bob, each having i number of interests, who want to securely compute the intersection of their interests. First, Alice and Bob exponentiate (encrypt) their interests with their respective keys and exchange the exponentiated interests. After that, Bob commutatively exponentiates Alice's interests with his key and then makes pairs of each of these commutatively exponentiated interests and corresponding Alice's exponentiated value. Bob sends all these pairs to Alice. Similarly, Alice commutatively exponentiates Bob's exponentiated values with her key. If her commutative encryption matches with those sent by Bob, then Alice identifies it with the first element of the pair. However, there are many possible attacks on this scheme. Firstly, an attacker can freely choose his/her interests during various runs of the protocol against the same user and eventually finds all the interests of the victim. Secondly, there is no limit on the number of interests. Therefore, an attacker can form a very large set of interests that include nearly every possible interest.

There is a strong chance that victims set will become a subset of attacker set and the attacker will know all elements of victim's interests.

Another drawback is that the initiator only learns the result of the evaluation. This allows an adversary to learn the results and then run away without running the protocol as a responder. Finally, Bob can reorder the pairs $(h(Ai))^a, ((h(Ai))^a)^b$. Therefore, Alice incorrectly identifies the matched interests.

B. Xie et al. Protocol

Xie et al. [6] identify the attacks on Agrawal's protocol and propose their protocol to overcome the shortcomings of [5]. In their protocol, they utilize two trusted servers. One is used to certify a user and the other is used to certify the interests of a user. The protocol in [6] uses commitments to ensure that any of the user should not be able to maliciously reorder the encrypted interests' pairs in step # 5. Once the intersection has been computed and mutual interests are identified, both Alice and Bob exchange the matched interests through a shared secret computed with the help of Diffie-Hellman exchange to ensure each other that both have computed the same result. Fig. 2 shows the complete working of [6]. Although this protocol offers improvements over Agrawal's protocol, but it introduces new attacks and fails to prevent some attacks that were also present in [5]. First, the protocol of Xie et al. uses two servers and therefore, assumes that both the servers are fully trusted. These servers contain critical user identity and interests' information and in case of a compromise, the privacy of the participants may severely be jeopardized. Another major drawback of [6] is that it does not prevent the attack where any of Alice and Bob reorders the pair $(h(Ai))^a, ((h(Ai))^a)^b$. The protocol assumes that such attack can be detected in the end where Alice and Bob exchange the interests. However, once Alice receives maliciously reordered interests from dishonest Bob, then she will send those same presumed interests to Bob in the shared key. Bob will decrypt the message and will simply send those interests back to Alice and trick her to believe that the matching was successful.

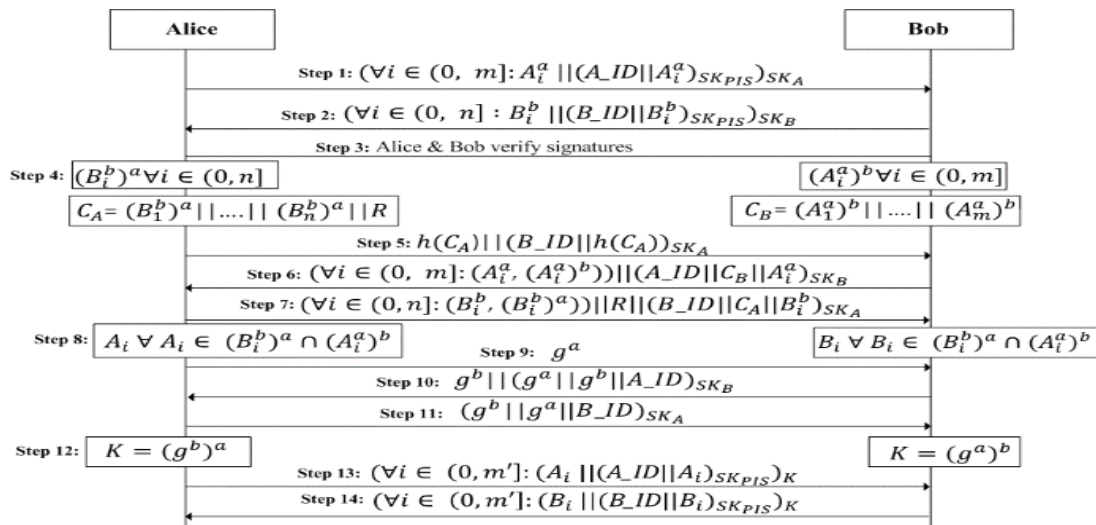


Fig. 2. Working of Xie et al. protocol [6].

Moreover, the protocol does not prevent from the attack where a malicious Alice will send gibberish values to Bob in step # 5 and then actual gibberish values in step # 7. Due to exponentiation, Bob will be unable to know the values are gibberish and will take it as an unsuccessful match. Therefore, Alice will learn the number of matched interests and Bob will know nothing.

C. Wang et al. Protocol

Wang et al. [7] proposed another protocol that attempt to overcome the shortcomings of both [5] and [6]. First, the

protocol in [7] combines the two trusted servers into a single server. Second, their protocol allows an initiator to run the protocol with several candidates in the first stage and finds the one candidate, described as the best match, with the most number of matches. Once the best candidate is found, the protocol proceeds almost in the same way as of [6]. However, in the end, instead of exchanging matched interests in a shared key, both Alice and Bob send the result to the trusted server that verifies the result and sends one's result to other. The Wang's protocol is given in detail in Fig. 3.

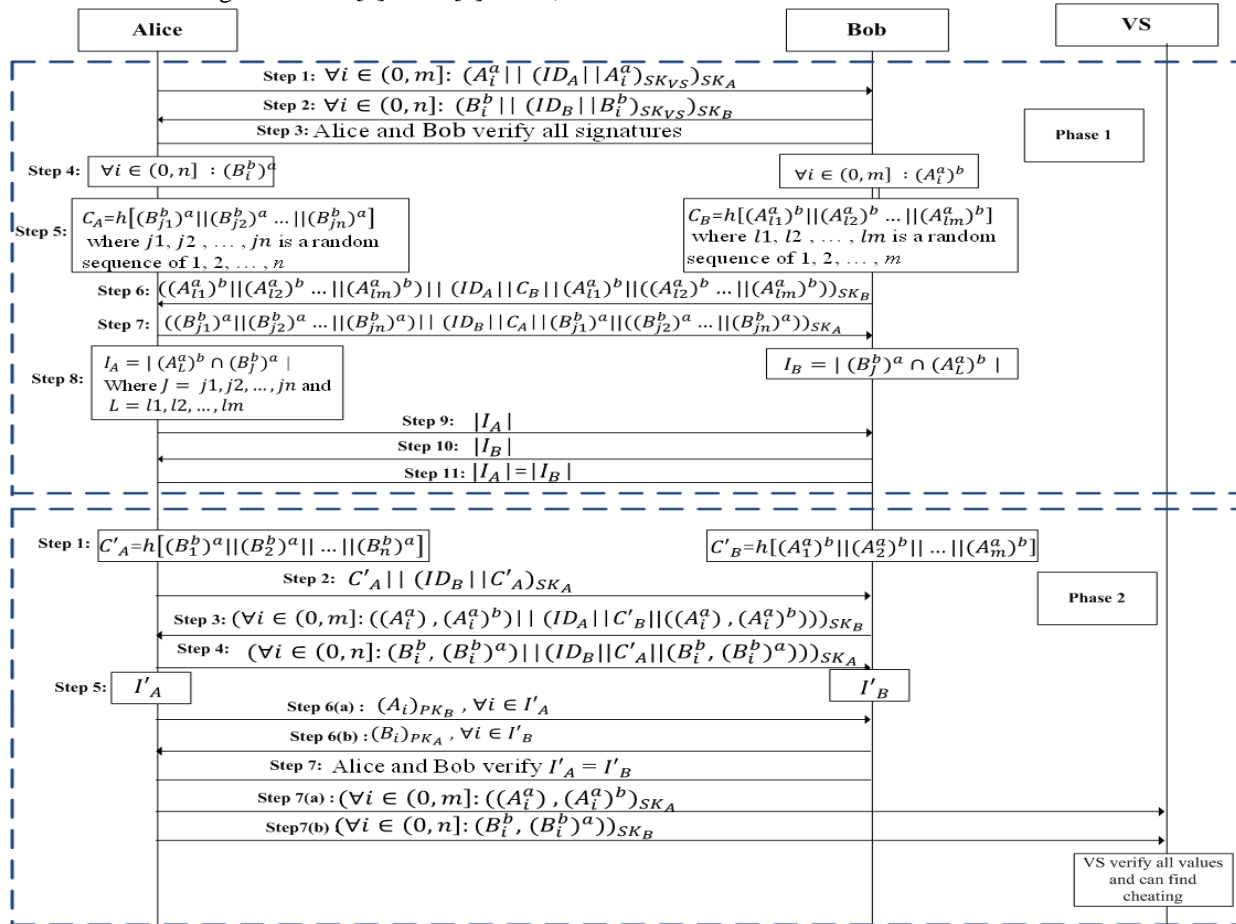


Fig. 3. Working of Wang et al. [7].

The protocol attempts to provide improvements over Xie's protocol. However, the main contribution seems to be the idea of finding the best match among a number of candidates. Many of the attacks on [6] are also possible in [7]. The unified servers till requires full trust of users, indeed, in the end of the protocol, the trusted server knows the result as well. Alice still can send the gibberish values to Bob and remains undetected. In the result, Alice knows the number of matched results. Similarly, the detection of cheating is not possible when both users exchange the actual interests in each other's public key. The malicious user will simply receive the actual values and send them back to other.

D. Fizza et al. Protocol

The authors of [8] propose a protocol to improve the work of [6] and [7] as shown in Fig. 4. They reduce the trust on

server by restricting the role of server in only verifying the number of interests of user. The server does not know the actual interests' values. Author in [8] uses the idea of introduces dummy interests in the interest set of both users. These dummy interests are known to both Alice and Bob but their relevant position in the set is only known to the set holder. Therefore, the gibberish values attack is nullified as the malicious user must correctly guess the position of dummy interests in the set which is very hard to guess. Moreover, author in [8] introduces a hash-based advantage less mechanism for interests exchange that ensures to find any mismatch in the exchanged results. However, one the drawback of [8] is the extra cost of exponentiating the dummy interests and the extra exchanges of commitments during the exchange of actual interests.

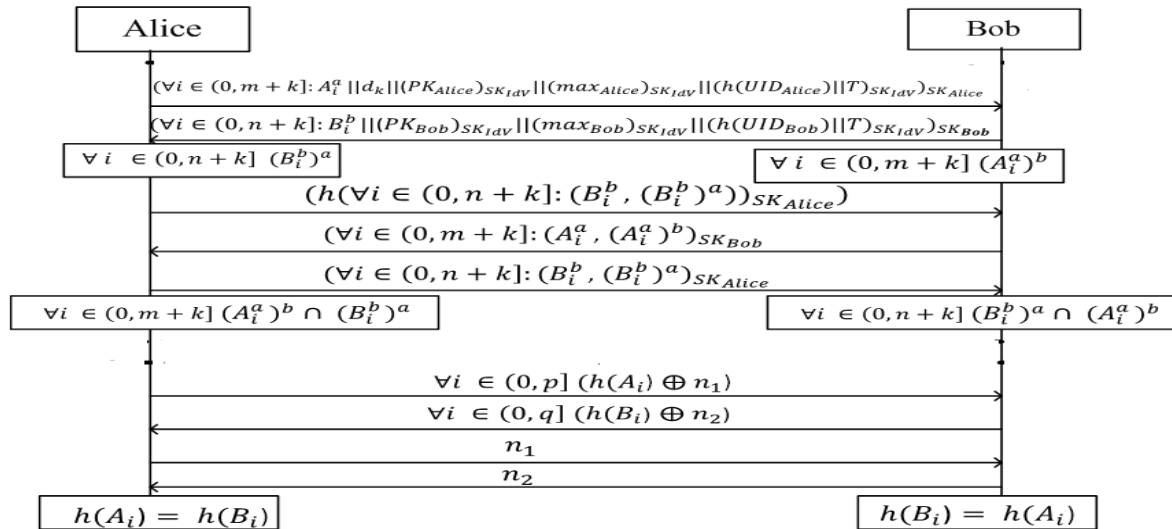


Fig. 4. Working of Fizza et al. protocol [8].

III. DISCUSSION

By looking at the literature survey it is evident that the commutative based encryption can provide significantly fast and reasonably secure PSI and PCSI based intersections. The protocols proposed for privacy preserving matchmaking based on commutative encryption have been increasingly secure. However, we a tradeoff where increased security and privacy requires increased computational and communication cost. The protocols in [5] paved the way for secure commutative based matchmaking. The protocols [6] and [7] improved the functionality and shortcomings of [5]; however, they introduced trust related issues as well as allowed malicious users to go unnoticed after cheating. The authors of [8] successfully eliminated above mentioned flaws but their protocol is slightly costlier in terms of communications and computation. Nonetheless, by keeping in mind the ever-increasing computing and communication capabilities of devices and telecommunication networks, these costs can be neglected by comparing the benefits being offered.

In future, there is need to improve the state of the art protocols by keeping the cost as low as possible. This becomes more significant because mostly matchmaking applications are designed for the people on the move carrying handheld devices and therefore, it required the matchmaking protocols to be light weighted both in terms of computation and communication.

IV. CONCLUSION

Matchmaking is one the famous application of mobile social network. Users share their private information with each other. To preserve users' privacy during matchmaking is not a trivial task. Therefore, many matchmaking protocols are proposed. This paper provides working of Agrawal et al., Xie et al., Wang et al. and Fizza et al. along with their limitations. In the end, a comparison is presented that compares the state of the art along with the benefit they offer over each other.

Finally, the paper signifies the need of light weight matchmaking protocols as possible future research directions.

REFERENCES

- [1] Y. Najafloo, B. Jedari, F. Xia, L. T. Yang, and M. S. Obaidat :safety challenges and solutions in mobile social networks," IEEE System Journal, vol. 9, no.3, pp. 834-854, 2015.
- [2] N. Kayastha, D. Niyato, P. Wang, and E. Hossain, " Applications, architecture, and protocol design issues for mobile social networks: A survey," Proceedings of the IEEE, vol. 99, no. 12, pp. 2130-2158, 2011.
- [3] L. P. Cox, A. Dalton, and V. Marupadi, "Smokescreen: flexible privacy controls for presence-sharing," in Proceedings of the 5th international conference on Mobile systems, applications and services. ACM, pp. 233-245, 2007.
- [4] X. Hu, T. H. Chu, V. C. Leung, E. C. H. Ngai, P. Kruchten, and H. C. Chan, "A survey on mobile social networks: Applications, platforms, system architectures, and future research directions," IEEE Communications Surveys & Tutorials, vol. 17, no. 3, pp. 1557-1581, 2015.
- [5] R. Agrawal, A. Evfimievski, and R. Srikant, "Information sharing across private databases," in Proceedings of the 2003 ACM SIGMOD international conference on Management of data. ACM, pp. 86-97, 2003.
- [6] Q. Xie and U. Hengartner, "Privacy-preserving matchmaking for mobile social networking secure against malicious users," in Privacy, Security and Trust (PST), 2011 Ninth Annual International Conference. IEEE, pp. 252-259, 2011.
- [7] Y. Wang, J. Hou, Y. Xia, and H. Li, "Efficient privacy preserving matchmaking for mobile social networking," Concurrency and Computation: Practice and Experience, vol. 27, no. 12, pp. 2924-2937, 2015.
- [8] F. Abbas, U. Rajput, and H. Oh, "Prism: Privacy-aware interest sharing and matching in mobile social networks," vol. 4, pp. 2594-2603, 2016.
- [9] M. J. Freedman, K. Nissim, and B. Pinkas, "Efficient private matching and set intersection," in International Conference on the Theory and Applications of Cryptographic Techniques. Springer, pp. 1-19, 2004.
- [10] C. Hazay and Y. Lindell, "Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries," in Proc. 5th Conf. Theory Cryptogr. (TCC), pp. 155-175, 2008.

Towards Security as a Service to Protect the Critical Resources of Mobile Computing Devices

Abdulrahman Alreshidi

College of Computer Science and Engineering
University of Ha'il
Ha'il, Saudi Arabia

Abstract—Mobile computing is fast replacing the traditional computing paradigms by offering its users to exploit portable computations and context-aware communications. Despite the benefits of mobile computing, such as portability and context-sensitivity, there are some critical challenges, such as resource poverty of mobile devices and security of mobile user's data that must be addressed. Implementing the security mechanisms to execute on mobile devices can be challenging as mobile devices lack the required processor, memory and battery resources to support continuous and long-term execution of computation intensive tasks. Cloud computing model can provide virtually unlimited hardware, software, and service resources to compensate for the resource poverty of mobile devices. In recent years, there is a lot of research and development of solutions and frameworks that preserve the security and privacy of mobile devices and their data. However, there has been little effort to secure mobile devices while also supporting an efficient utilization of the limited resources available on mobile devices. In this paper, we propose Security as a Service for mobile devices (SeaaS for mobile) that integrates mobile computing and cloud computing technologies to secure the critical resources of mobile devices. The proposed solution aims to support 1) security for the data critical resources of mobile devices, and 2) security as a service by cloud servers for an efficient utilization of the mobile device resources. We demonstrate the security as a service based on a practical scenario for the security of mobile devices. The evaluation results show that the proposed solution is 1) accurate to detect the potential security threats, and is 2) computationally efficient for mobile devices. The proposed solution as part of ongoing research provides the foundations to develop a framework to address SeaaS for mobile. The proposed solution aims to advance the research state-of-the-art on software engineering, mobile cloud computing, while it specifically focuses exploiting cloud-based services to secure mobile devices.

Keywords—Software engineering; mobile computing; cloud computing; computer security; mobile cloud computing; security as a service

I. INTRODUCTION

Mobile computing has fast emerged as a pervasive technology that empowers its users to exploit portability and context-awareness to perform a variety of tasks on the go [1]. Specifically, mobile computing can utilize the embedded sensors (i.e., GPS, Accelerometer as hardware resources) that can be combined with freely available apps (i.e., location services, maps as software resources) of a mobile device to support context-aware and portable computing [2]. For example, the mobile users can enable GPS based location

sensing to get live updates about traffic conditions or recommendations about the places/events of interest based on their geographical proximity. Despite these benefits, mobile computing in general and mobile devices in particular face two primary challenges [8], [9]. The first challenge relates to the resource poverty, i.e., the availability of limited processing, memory and battery resources to a mobile device. The second challenge is to protect the integrity of the mobile device that is prone to the threats of data security and privacy in a context-aware environment. The security threats relate to the security critical resources that are hardware (e.g., Microphone, GPS sensor) and software (e.g., Contacts List, Photos) resources. For example, a third part game installed on a mobile device can try to maliciously access the Microphone to spy on user's voice conversation or look into user's contact list for information [3]. If such private information can be compromised or exploited by entities with malicious access, it can put user's information and device's data on security risk [14].

A. Research Challenges

There is a need for a rigorous security mechanism(s) that protects the data critical resources of a mobile device to support secure mobile computing. However, any rigorous security solution(s) that continuously execute on a mobile device to protect itself may be impractical mainly due to the computation, memory and battery specific resource poverty of the mobile devices. There is a need for solutions that must ensure a rigorous security as well as efficient resource utilization of a mobile device [4]. Cloud computing represents an opportunistic computing model that relies on the 'pay-per-use' hardware and software services that can be used and released as required. Cloud computing model offers three main types of services referred to as Software as a Service (SaaS), Platform as a Service (PaaS), and Infrastructure as a Service (IaaS) [5], [15], [16]. In recent years, mobile and cloud computing technologies have been unified to enable Mobile Cloud Computing (MCC) as state-of-the-art mobile computing technology. Specifically, in MCC a mobile device represents a portable and context-aware user interface (as front-end technology) that relies on the resource sufficient cloud-based servers to perform the complex computations (as back-end technology) - off-loaded by mobile devices to the cloud-based servers. For example, the solutions of MCC in [6] allow mobile devices to off-load their computation intensive tasks to cloud-based servers to prolong the battery life and enhance the processor and memory performance of the mobile devices.

This means that existing solutions of MCC have been successful to support mobile devices that are portable, context-sensitive, as well as resource sufficient by relying on cloud computing.

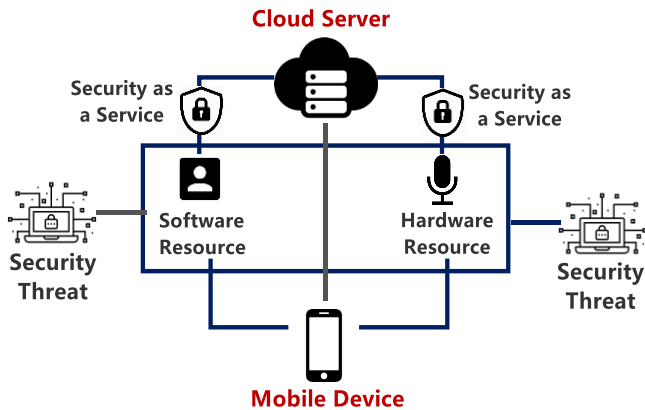


Fig. 1. Overview of the Solution to Support SeaaS for Mobile Devices.

B. Proposed Solution, Contributions, and Assumptions

We propose to address the security related issues that are relevant to the critical resources of mobile devices by integrating the mobile computing and cloud computing technologies. In the proposed solution, we aim to develop a framework that provides Security as a Service for mobile devices (SeaaS for mobile) offered as a cloud-based service to secure mobile computing. We define SeaaS for mobile as ‘security policies that are executed on cloud-based servers to continuously monitor and protect the critical (hardware and software) resources of a mobile device’. Critical resources of a mobile device are any hardware and software resources that produce or consume the information that can be prone to security threats and malicious access. We provide an overview of the proposed solution in Fig. 1. The figure highlights that the mobile device and cloud servers are connected. The cloud server runs the SeaaS to continuously monitor the critical hardware (e.g., Microphone) as well as software (e.g., Contacts) resource. Security specific issues of the mobile device resources are offloaded to the cloud-based server. Any potentially malicious access is detected and is eliminated or minimized by the cloud-based SeaaS. The primary contributions of research are summarized below.

- *Unification of Mobile and Cloud Computing Technologies* to exploit security as a service to protect mobile devices that operates in ad-hoc (unsecured) network environment.
- *Off-loading the Execution of Security Mechanism and policies* on cloud-based servers to support secure and efficient mobile computing.

Based on the proposed contributions, we have the following assumptions for the proposed solution.

- *A Continuous Network Connectivity* is required between the mobile device and cloud-server for the monitoring and protection of the device resources by cloud-based services.

- *Integration of Mobile and Cloud Computing* where a mobile device represents a client, whereas the cloud-based server(s) act as security provision resources.

The rest of the paper is organized as follows. Background details about Mobile Computing Environment is provided in Section 2. Related Work is presented in Section 3. The Proposed Solution is presented in Section 4. Solution Demonstration and Evaluation Results are presented in Section 5. Conclusions and Future Research are detailed in Section 6.

II. BACKGROUND ABOUT MOBILE COMPUTING ENVIRONMENT

We now present an overview of the mobile computing environment as in Fig. 2. We briefly discuss the elements of the mobile computing environment in the context of mobile security. The concepts and terms introduced here will be used in the remainder of this paper. As highlighted in Fig. 2, the elements of a mobile computing environment are introduced below. In Fig. 2, we have adopted this general model of mobile computing environment from [14].

A. Mobile Device

It represents any (handheld) equipment or machine that allows its user to perform computation, information sharing and other activities in a mobility-driven environment. A mobile device is a combination of hardware (e.g. GPS sensors, Camera) that is manipulated by means of software apps (e.g. Location Tracker, Image Editor) and both are vulnerable to the security threats.

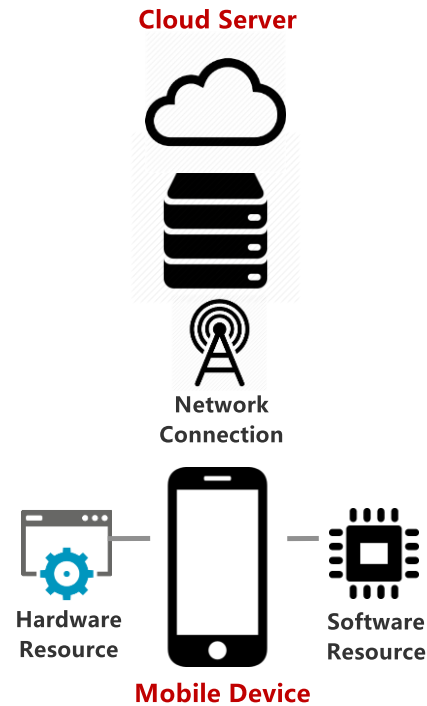


Fig. 2. Overview of the mobile computing environment.

- **Hardware Resources** are physical parts of a device such as sensors, processors that are sources of private

information. For example, to uniquely identify a device's hardware, the device-specific information of a mobile device is represented with Unique Device Identifier (UDID) or International Mobile Equipment Identity (IMEI) that are subject to security threats. In terms of software, a mobile device equipped with Global Position System (GPS) represents a scenario where user's location or context can be revealed or leaked to location sensing services.

- **Software Resources** are the applications or data that contains the useful information or private data of the device users. For example, the Contacts List or Photos represent the software resources of a mobile device and these resources are at a risk of comprising the private information.

B. Mobile Server

It provides an infrastructure that allows a mobile device (client) to store or retrieve data, request the desired services and to off-load the data for computation on the server. In Fig. 1, the server can be of many types (e.g., Communication, Proxy, and Database) that can provide a lot of new functionality such as location services, look-ups in directory services and enabling distributed data storage and processing. With an emergence of the mobile-cloud computing a mobile device as a resource-constrained computer can exploit virtually unlimited computation and storage via resource sufficient cloud server. Specifically, cloud servers have been proven successful solutions to compensate the limited computation, storage and power resources of mobile devices. In a mobile-cloud computing environment, the mobile acts as a context-aware and portable client to capture and display data. In comparison, the cloud represents a backend server that supports the computation and storage of all the data off-loaded to it by a mobile device. A mobile server is prone to invasion that can compromise the security of mobile data residing on the server.

C. Network Connectivity

It allows a mobile device to communicate with a server through network connections such as Wi-Fi or Bluetooth signals. This means that in addition to the data in a mobile device or the one residing on the server, the communication channel can be attacked to compromise the data that travels between a mobile and its corresponding server. A typical example is the attack on location queries that travel between a mobile device and location-providing server [3]. In addition, a direct communication between mobile devices is also subject to security threats on the communication channel. Based on the illustration in Fig. 1, we conclude that security in mobile computing environment allows protection and preservation of (user and device) data or information deemed as private from acts of malice [5], [6].

III. RELATED RESEARCH ON SECURE MOBILE COMPUTING

In this section, first we highlight the existing research (Section 3.A) and then discuss some proposed solutions as tools and frameworks (Section 3.B) that enable or enhance mobile computing security. By presenting the most relevant

related work, we justify the scope and proposed contributions of our research.

It is vital to mention that, according to the GSMA real-time tracker, the world is currently home to more than 7.2 billion mobile device connections, where 0 to 7 billion connections have been achieved in just three decades [7]. This implies that mobile connections are currently growing about five times faster than the human population. In a recent report, the Homeland Security has highlighted that security specific threats to mobile computing such as user location tracking, banking and transactions fraud, ransom-ware, identity theft puts at risk not just mobile device users, but the mobile carriers as well as infrastructure providers [3].

A. Summary of Existing Research on Security for Mobile Computing Environments

A survey-based study highlights the potential threats and proposed solutions for devices that operate in mobile and ad-hoc networks [8]. In recent years, there is a lot of focus on solutions to enable or enhance the security of the mobile devices. Specifically, in [9], the authors have highlighted the potential and a huge market for android applications, however; there are concerns relating to the security and privacy issues of these apps. Unless there is a strong mechanism for a device to protect itself, the widely available apps are subject to potential security threats. Moreover, a mobile device has limited resources, i.e.; memory and processing power which means that it becomes challenging to maintain security of its data.

The study [10], suggests a balance between IT infrastructure overhead and system security. The study has considered four application level security systems and evaluated them against a pre-defined scheme i.e., systems support for critical security related services such as authentication, authorization, maintainability, re-usability, productivity etc. On the basis of the evaluation results, the study concludes that that none of the selected systems fulfilled the evaluation schemes. In the mobile computing context, the challenge lies with providing a robust security mechanism while also supporting the efficiency of the devices memory, computation and energy efficiency.

TABLE I. A COMPARISON SUMMARY OF THE EXISTING FRAMEWORKS FOR MOBILE COMPUTING SECURITY.

Proposed Solution	Mobile Computing	Cloud Computing	Proposed Contributions
Xposed Framework	✓		Mobile Apps Monitoring
Android Monkey UI Exerciser		✓	Testing Android Apps
Resource Description Framework (RDF)	✓		Encode Data for Exchange
Proposed Solution	✓	✓	Mobile Device Security

B. Summary of Frameworks and Tools for Mobile Computing Security

After presenting the research challenges, we also highlight some existing frameworks and tool support that aim to automate and enhance adaptive security solutions. These tools and frameworks are mostly proof of the research concepts.

- *Xposed* framework is an open-source framework that allows monitoring of the installed apps can change system settings i.e., behaviour of the system or the installed apps as per the security needs [11]. The users can change the settings of the framework at runtime and can also bring the changes based on those dynamic settings.
- *Android Monkey UI Exerciser* framework is basically a third-party tool which helps in testing android applications. It is basically a command line tool which works with adb tool (Android Debug Bridge) [12]. It is basically used to perform stress testing on the android applications and to report back the errors, if occurred.
- *Resource Description Framework (RDF)* is the framework for data interchanges over the internet. This framework is capable of encoding, exchanging and reusing the structured metadata [13]. It is an application of XML that provides a description to specify the resources and the associated threats.

1) Comparative Summary – Existing vs Proposed Research

A comparative summary between the existing solutions and the proposed solution is provided in Table I. Specifically, Table I highlights that in comparison to the existing Solutions,

the proposed solution exploits mobile-cloud computing technologies to enable the security of mobile devices. Based on the discussion above, there is a lack of research that supports cloud-based security as a service to protect resource constrained mobile devices. The proposed research aims to off-load security mechanisms to cloud-based servers to support secure and efficient mobile computing.

IV. A LAYERED SOLUTION FOR SECURITY AS A SERVICE FOR MOBILE COMPUTING

We now present an overview of the framework in Fig. 3. Fig. 3 highlights that the proposed solution is based on a layered architecture system. Specifically, the proposed solution has the following two layers namely Mobile Computing Layer and Cloud Computing Layer.

A. Layer I - Security Critical Mobile Computing Layer

It is the front-end layers that allow a user to perform computing and communication in a portable and context-aware fashion. As highlighted in Fig. 3, the mobile computing layer has some hardware and software resources that are critical from a security point of view. For example, the *Microphone* that is a hardware resource must be protected from any malicious access to comprise the privacy of *user's voice* and *audio messaging*. In a similar scenario from Fig. 3, *Contacts* that represents a software resource containing the critical information such as contacts' *names, numbers, email, photos* that can be compromised. The intent for any unauthorized access to the hardware and software resources can be due to spying on user's private information or selling user's private information to the third part advertisement providers. In such circumstances there is a need to secure mobile computing resources.

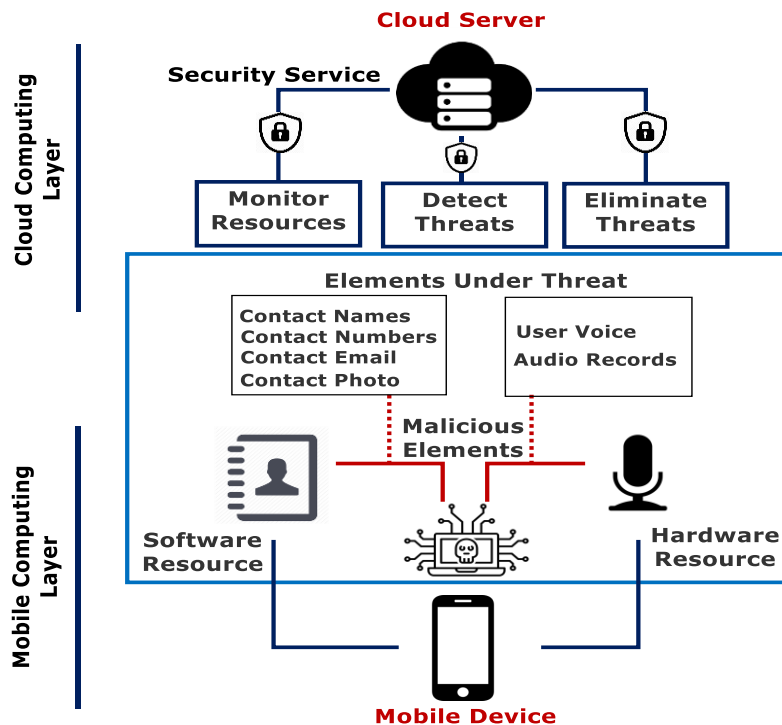


Fig. 3. Overview of the mobile cloud based security as a service.

However, the main challenge relates to the resource poverty of mobile devices that cannot support an effective execution of the rigorous and computational intensive security protocols.

B. Layer I - Security Critical Mobile Computing Layer

To compensate for the resource constrained mobile devices, the backend cloud computing layer can provide the security as a service for the critical resources of a mobile device [16]. As highlighted in Fig. 3, the cloud-based Security as a service can relieve a mobile device from securing its resources in a way that cloud-based server continuously monitors the critical resources of a device and secures them by eliminating and or minimizing any security threats. As highlighted in Fig. 3, the cloud based Security as a Service follows a three steps process that supports: 1) Monitoring of the device resources, 2) Detection of any unwanted and

potentially malicious access as a security threat, and 3) Eliminating or minimizing the security threat to secure resources of a mobile device.

V. DEMONSTRATION AND EVALUATION OF THE FRAMEWORK

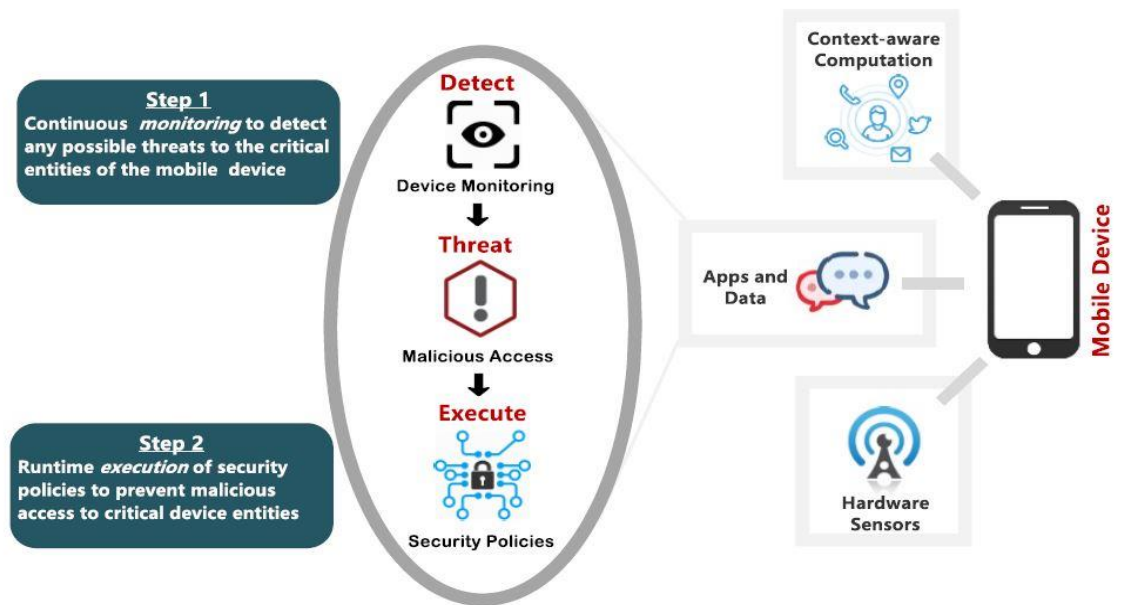
After presenting the overview of the solution, we now demonstrate the usability of the framework based on a case study (Section 5.A). We then present the results of the evaluation of the framework (Section 5.B).

A. Demonstration of the Security as a Service Framework

The scenarios presented here are taken from [14]. Scenario-based demonstration highlights the applicability of the proposed framework based on scenario-driven approach as in Table II. Table II presents the scenarios for Mobile Computing security.

TABLE II. OVERVIEW OF SCENARIOS FOR SECURITY AS A SERVICE IN THE CONTEXT OF MOBILE COMPUTING

<p>Scenario I - Malicious Access to Hardware Resources</p> <p>Security Challenge: how to protect a mobile device’s hardware resources (e.g.; accelerometer, gyroscope, wireless or GPS sensors) from any act of malice that compromises the de-vice’s data and security?</p> <p>Proposed Solution: The proposed solutions offer Perceptual Monitoring as a mechanism to monitor, control and customise access to a devices sensor to safeguard user’s private data (e.g.; locations, actions, movement). Such perceptual monitors enable or enhance a device’s security by working as:</p> <ul style="list-style-type: none"> • Monitor the access of the device’s sensors. • Customise the sensor usage policy of third-party apps (e.g.; grant, deny, or selective permission). • Runtime modification of the sensor access permission as per user’s needs and requirements.
<p>Scenario I - Malicious Access to Mobile Device Data and Resources</p> <p>Security Challenge how to enable the security of mobile devices that prevents or minimizes any malicious access to mobile apps and leak-age of private data?</p> <p>Proposed Solution provide Reconfigurable Security Policies and device data monitors that are executed on a mobile device. These policies and runtime monitors can dynamically configure their behaviour -depending on the context of data/app being accessed -to prevent or minimise any malicious access to device data and apps. Reconfigurable security policies work as follows:</p> <ul style="list-style-type: none"> • Define access policies for device’s data and apps. • Execute policies as backend processes to monitor malicious access to app or misused data • Reconfiguration policies as per the context of app or data access and taking into account user’s customisation of the policies.



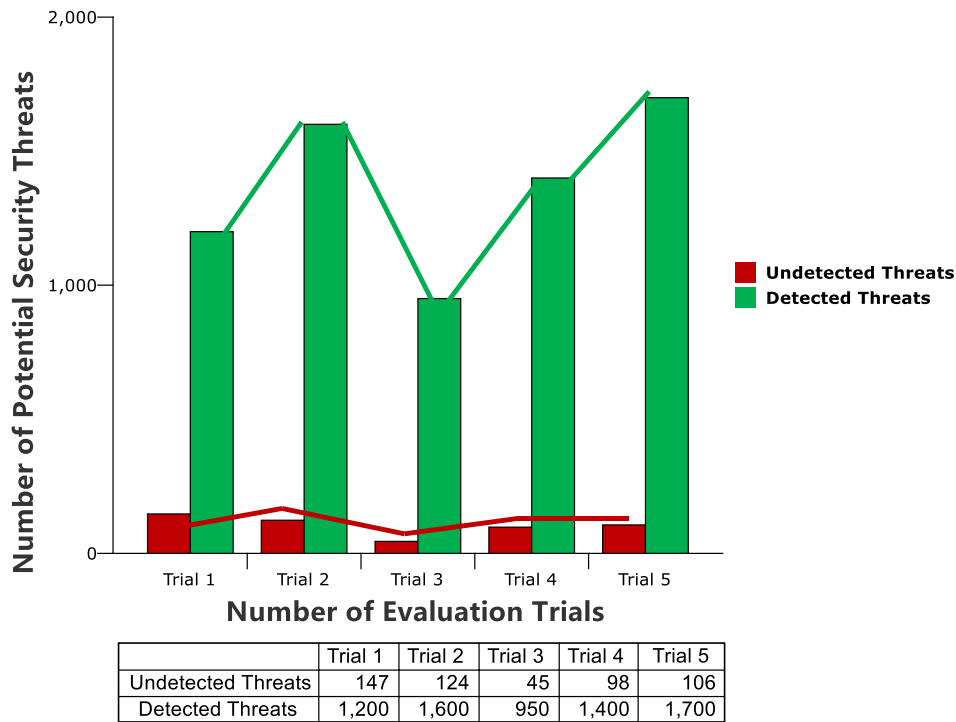


Fig. 4. Overview of the mobile cloud based security as a service.

A. Results of Preliminary Evaluations

After presenting the demonstration of the framework, we now discuss the results of the evaluation to validate the accuracy and efficiency of the framework. Measurement of both the accuracy and efficiency represent the qualitative evaluation of the framework based on the ISO-IEC-9126 model for software quality¹. To conduct the preliminary evaluations of the proposed solution, we have used:

- **Mobile Computing Layer:** We have used HTML5 technologies for mobile front-end development to target multiple mobile platforms. Another reason of using hybrid mobile application instead of native mobile application development is that we carry out all performance intensive tasks over cloud layer. Considering the mobile services development via HTML5 technologies, we exploit ionic framework to target android and iOS platform.
- **Cloud Computing Layer:** We exploit Amazon cloud services for storage and computing efficiency. Pricing model adopted by AWS is pay-as-you-go. We launch a virtual server on Amazon cloud called Amazon EC2 instance and set up Red Hat Linux operating system over the instance. We develop server-side application using Node.js and set up Node.js web server on Amazon EC2 Instance. For the sake of efficient data retrieval, we use MongoDB. We install MongoDB on

Amazon EC2 Instance. To use files and media we utilize Amazon S3 storage services.

- **Accuracy to Detect the Potential Threats:** We now present a summary of the framework's accuracy to detect the potential security threats. An overview of the results of the preliminary evaluations is presented in Fig. 4. Specifically, Fig. 4 shows the total number of trials (X-axis) along with the number of detected/undetected threats. We conducted a total of 5 trials where each trial engaged 10 users on average to evaluate the accuracy of potential security threat detection. As per the ISO/IEC-9126 model, accuracy refers to the system's ability to correctly compute the results. As illustrated in Fig. 4, based on 5 trials, the total number of detected threats is 6850, while the detected threats are 520. Based on the average, the ratio of detected to undetected threats (Detected/Undetected) is 13:1 that demonstrates a high-level of accuracy for threat detection.
- **Computational Efficiency of the Proposed Solution:** In addition to the accuracy, as highlighted earlier, we also need a solution that is computationally efficient. By computationally efficient we mean that the proposed solution has an efficient utilisation of the mobile device processor. The efficiency can be enabled by off-loading the computational intensive tasks to the cloud-based servers [6], [14]. Fig. 5 highlights the CPU utilization for the mobile device. Fig. 5 highlights two pieces of information, (i) the ratio of processor consumption and (ii) total number of trials.

¹It is noteworthy that, ISO/IEC 9126-1 was first published in 1991; and later on from the year 2001 to year 2004 ISO published an international standard (ISO/IEC 9126-1) as well three technical reports (ISO/IEC 9126-2 to ISO/IEC 9126-4)

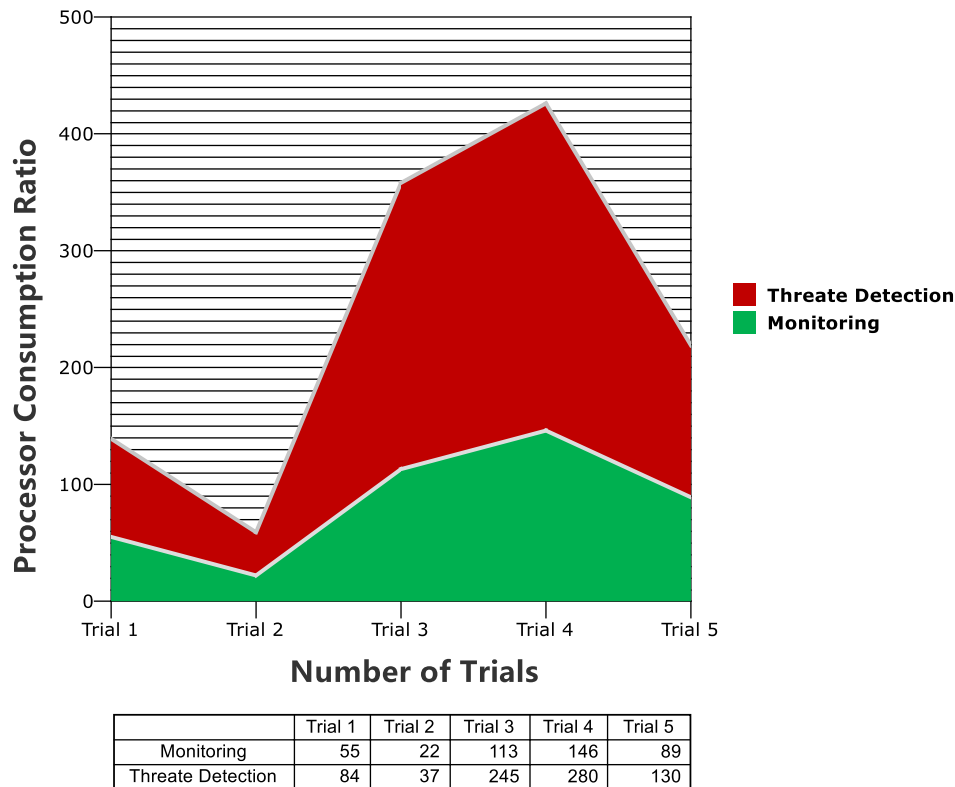


Fig. 5. Overview of the mobile cloud based security as a service.

The results in Fig. 5 show that while monitoring the potential security threats, how much of the device's processor has been utilised. We conducted a total of 5 trials where each trial engaged 10 users on average to evaluate the accuracy of potential security threat detection. Fig. 5 also demonstrates that only a small percentage of the device's CPU has been used while monitoring for the security threats. Such efficient utilisation of the processor is only possible due to the fact that all computation intensive tasks have been performed on the cloud-based server.

We conclude that in the scope of the proposed solution that is part of the ongoing project, we have only conducted the preliminary evaluation of the proposed solution in terms of system's accuracy and efficiency. The preliminary results show that the proposed solution is efficient and accurate. However, for more concrete and objective evaluation, we need more trials and further development of the system that is part of the future work.

VI. CONCLUSIONS AND FUTURE WORK

Mobile computing supports portability, context-sensitivity, and enhanced user interaction to replace the traditional computing paradigms. Despite these benefits, mobile computing faces a number of challenges such as resource poverty of a mobile device and threats to the security and privacy of users' information and device's data. Specifically, to address the issues of mobile device security in an efficient way we have proposed a novel solution that relies on the integration of mobile devices and cloud servers to enable or

enhance a mobile device's security. The proposed solution aims to address two of the most prominent challenges for mobile computing namely security and efficiency of mobile computing. The proposed solutions exploit a layered approach to address these challenges by offloading the security mechanisms of mobile device to cloud-based servers. Specifically, the front-end layer (i.e., mobile device) represents a portable and context-aware computer that relies on the back-end layer (i.e., cloud server) to monitor and protect the critical resources of the mobile device. The integration of mobile and cloud computing technologies as state of the art mobile computing technology aims to support secure mobile computing.

We have presented the solution architecture and demonstrated its application to enable mobile device security. The results of the preliminary evaluation suggest that proposed solution is (i) accurate for the detection of potential security threats, and (ii) it off-loads computation intensive tasks from mobile devices to cloud-based servers to enable efficient mobile computing. We conclude that the proposed solution can be helpful for

- Advancing the state-of-the-art on mobile computing technology to support mobile-cloud driven security framework.
- Enabling the solution for secure and efficient mobile computing by means of cloud-based security.

A. Possible Future Research

The proposed solution provides a framework and the foundations to develop a comprehensive tool support as a proof-of-the-concept to enable and automate the concept of Security as a Service for the critical resources of mobile devices. Therefore, as part of the future work we mainly focus on the development of the framework that provides an executable solution for further evaluation. Moreover, the proposed solutions need concrete scenarios of security threats that can be executed and analysed to demonstrate the applicability and validation of the proposed solution. We are particularly interested in exploiting the existing algorithms and solutions that can leverage the cloud computing resources to secure mobile devices.

REFERENCES

- [1] Pejovic, and M. Musolesi. Anticipatory Mobile Computing: A Survey of the State of the Art and Research Challenges. In ACM Computing Surveys (CSUR), vol 47, no 3, pp. 47, 2015.
- [2] C. Perera, A. Zaslavsky, P. Christen, D. Georgakopoulos. Context aware computing for the internet of things: A survey." IEEE Communications Surveys & Tutorials 16.1 (2014): 414-454.
- [3] Study on Mobile Device Security - Homeland Security, April 2017. [Online:] <https://www.dhs.gov/sites/default/files/publications/>, Accessed 23-12-2017.
- [4] A. N. Khan, M. L. Mat Kiah, S. U.Khan, S. A. Madani. Towards Secure Mobile Cloud Computing. In Future Generation Computer Systems, vol. 16, no. 1, pp: 1278 – 1299, Science Direct, 2013.
- [5] P. Mell and T. Grance. The NIST Definition of Cloud Computing. Special Publication 800-145 (Draft), National Institute of Standards and Technology, Gaithersburg, Maryland, 2011. [Online:] <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>.
- [6] G. A. Lewis, P. Lago, G. Procaccianti. Architecture Strategies for Cyber-Foraging: Preliminary Results from a Systematic Literature Review. In 8th European conference on Software Architecture (ECSA'14), 2014.
- [7] GSMA Intelligence - Definitive Data and Analysis for the Mobile Industry: [Online:] <https://www.gsmainelligence.com/>, Accessed 24-12-2017.
- [8] [D. Djenouri, L. Khelladi, and N. Badache. A Survey of Security Issues in Mobile Ad hoc Networks. In IEEE communications surveys, vol 7, no 4, pp: 2-28, 2005.
- [9] M. Garcia, D. Llewellyn-Jones, F. Ortin, and M. Merabti, Applying Dynamic Separation of Aspects to Distributed Systems Security: A Case Study, Software, IET, vol. 6, no. 3, pp. 231–248, June 2012.
- [10] S. Alampalayam and A. Kumar, "An adaptive security model for mobile agents in wireless networks," in Global Telecommunications Conference, 2003. GLOBECOM '03. IEEE, vol. 3, Dec 2003, pp. 1516–1521 vol.3.
- [11] Xposed - general info, versions & changelog, xda developers, 2016. [Online]. Available: <https://forum.xda-developers.com/xposed/xposedinstaller-versions-changelog-t2714053>
- [12] Android Testing. www.tutorialspoint.com, 2017. [Online]. Available: https://www.tutorialspoint.com/android/android_testing.htm
- [13] E. Mille. (1998) An introduction to the resource description framework. [Online]. Available: <http://www.dlib.org/dlib/may98/miller/05miller.html>
- [15] M. Sajjad, A. Ahmad, A. Malik, A. B. Altamimi, I. M. Alseadoon, I. M. Classification and Mapping of Adaptive Security for Mobile Computing. IEEE Transactions on Emerging Topics in Computing, 2018.
- [16] Jamshidi, Pooyan, Aakash Ahmad, and Claus Pahl. "Cloud migration research: a systematic review." IEEE Transactions on Cloud Computing 1, no. 2 (2013): 142-157.

Comparison of Task Scheduling Algorithms in Cloud Environment

Babur Hayat Malik, Mehwashma Amir, Bilal Mazhar, Shehzad Ali, Rabiya Jalil, Javaria Khalid

Department of CS and IT
The University of Lahore, Gujrat Campus, Pakistan

Abstract—The enhanced form of client-server, cluster and grid computing is termed as Cloud Computing. The cloud users can virtually access the resources over the internet. Task submitted by cloud users are responsible for efficiency and performance of cloud computing services. One of the most essential factors which increase the efficiency and performance of cloud environment by maximizing the resource utilization is termed as Task Scheduling. This paper deals with the survey of different scheduling algorithms used in cloud providers. Different scheduling algorithms are available to achieve the quality of service, performance and minimize execution time. Task scheduling is an essential downside within the cloud computing that has to be optimized by combining different parameter. This paper explains the comparison of several job scheduling techniques with respect to several parameters, like response time, load balance, execution time and makespan of job to find the best and efficient task scheduling algorithm under these parameters. The comparison of scheduling algorithms is also discussed in tabular form in this paper which helps in finding the best algorithms.

Keywords—Task scheduling; algorithms; cloud computing; min-max; genetic algorithm; load balancing; resource utilization

I. INTRODUCTION

In scientific community, Cloud Computing has gained a vast amount of attention. Cloud Computing provides an environment which is more flexible rather than its counterparts. Cloud Computing provides the facility to access the data anywhere from your cloud [1]. Organizations are shifting their businesses toward cloud computing because cloud computing providing resources in large quantity and user/organizations are using resources freely.

Cloud Computing is a model which provide easy access to available resources to cloud users on their demand [2]. Cloud provides a variety of services on the demand of to its user, e.g. dynamically network access, rapid elasticity. The popularity of cloud depends on its performance, manage resource and optimally job scheduling.

This paper is mainly focusing on different task scheduling approaches. Task scheduling can be defined as choosing the most appropriate and suitable resources for the execution of the task. The task can also be defined as user's queries send to the different server, and these queries also accomplished within required time period [8]. Task scheduling works on principle of distributed the task on available resources.

The main objective of scheduling algorithms in

decentralized environment is to extend different task on servers to balance the load, this maximize the utilization of processors and minimize the execution time of user task. The central objective is to schedule available resource according to the available time for its execution. The task may include entering a query, a process that query, accessing the required software and memory [1]. Then data centre classifies user's queries on the requirements made on the service requested and agreement of services.

The user task is appoint to one of the available servers, and result or response of task is sent back to the user. Task by the cloud users are dispatched to available resources for their timely execution is task scheduling [26]. Several task scheduling Algorithms are used to increase the performance of cloud and enhanced throughput of servers. The different parameters of scheduling are used to increase the overall cloud performance [2]. There are several limitations while scheduling a task such as a cost, throughput, time, resource utilization and make span [26]. The main contribution in task scheduling is to minimize cost and time to produce an optimal result which causing to increases the performance of the cloud.

The coming part of the paper will explain classification of scheduling and scheduling process. Section III is comprised of the literature review and Section IV explained the working of several task scheduling algorithms helping to make a brief comparison of different algorithms and discusses the results of this research. After the comparison of scheduling algorithm at last we conclude the best and efficient scheduling algorithm.

II. CLASSIFICATION OF SCHEDULING

Scheduling methods are classified into three main groups: task scheduling, resource scheduling and workflow scheduling. The distribution of virtual resources among Servers (physical machine) is done by Resource Scheduling. To scheduling workflow comprised by an entire job in an efficient order. Task Scheduling is to assigned the task to available resource for its execution. Task Scheduling method is for centralized as well as decentralized structure and also for the homogenous and heterogeneous environment [25]. The paper mainly focuses on the task scheduling algorithms and their comparison. Several Task Scheduling Algorithms are diagrammatically shown in Fig. 1. Allocating resources to any task is considered as task scheduling, and it is the main component of cloud computing. The most significant factor in task scheduling is time and cost is required for its completion.

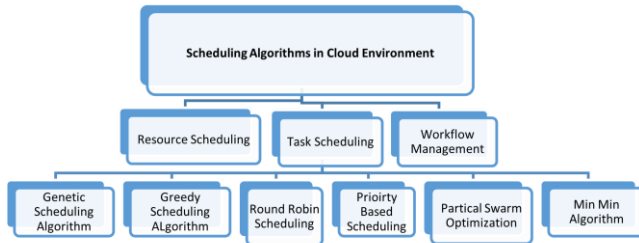


Fig. 1. Classification of scheduling techniques [2].

III. SCHEDULING

Distributed computing environment comprised of several scheduling techniques. These techniques is implemented to cloud environment with applicable parameters. The central purpose of these techniques is to improve the response time of system and the performance of the cloud [3]. Conventional Scheduling techniques are not able to increase the performance, it maximize the cost and execution time as well. Scheduling algorithms which are discussed in this paper are min-min, First Come First Serve, Round Robbin, Genetic Algorithm, Greedy Algorithm, Partical Swarm Optimization, Priority based scheduling, etc. [4], [24].

A. Scheduling Process

Cloud Task scheduling process is generally classified into three stages [7] are shown in Fig. 2:

Resource Dicorvey and Filtering: The cloud service provider discover list of available resource in a secific network and also collect and check their working status.

Selction of Resources: It is the most important stage in task scheduling, is also known as deciding stage. Required resources are selected on specific parameter and according to the requiremt of task execution.

Submission of Task: After slecting the required resource the task is submit to the resource for execution [34].

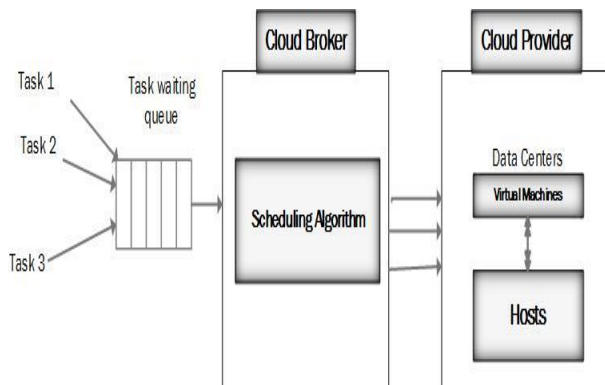


Fig. 2. Task scheduling process [4].

IV. LITERATURE REVIEW

The major issue in task scheduling is the allocation of efficient and available resources to the new task enter by the user. If several tasks arrive at the same time, then dynamically

resource allocation process is become more complex. Therefore, S. Ravichandran and D. E. Naganathan [7] had proposed a new system to get rid from this problem, this system works when a new task is arrived it is sent into the queue for waiting and the job scheduler will easily order each task and allocate resources for their execution. Genetic Algorithm is considered as a best practice in this regard, all the tasks are sent to the queue, scheduler pick the task from the queue allocate resources and execute the task. The central purpose of this system is to minimize the execution time of the task and optimize the resource utilization.

In this paper authors, V. V. Kumar and S. Palaniswami mainly focused on enhancing the efficiency of job scheduling techniques for cloud computing service. They have also purposed an algorithm which optimally utilized the turnaround time by giving high priority to different job for its completion and less priority for termination issues of real-time task [9]

Moreover, a new task scheduling algorithm was purposed by Z. Zheng which is based on genetic algorithm which is termed as Parallel Genetic Algorithm. The main objective of this algorithm is to optimize the cloud scheduling problem mathematically.

Siad Bin Alla and Hicham Bin Alla in [22] have explained a novel based dynamic task scheduling technique which is based on improved genetic algorithm. The working of proposed algorithm is mainly focused to reduce the execution time, effectively improve the throughput and the scalability of the cloud in task scheduling.

In [33], author proposed a novelistic approach for task scheduling algorithm M Quality of Service with Genetic Algorithm and Ant Colony "QOS-GAAC" with multi-QOS constraint, in this algorithm author mainly focused on expenditures, security, time-consuming and reliability in the process of task scheduling. This algorithm is the combination of genetic Algorithm and ant colony optimization algorithm. The result represent that this algorithm has great performance in both guaranteeing QOS and resource balancing in task scheduling [27].

Author proposed an optimized algorithm to minimize the cost and bi-objective makespan used by heuristic search techniques for independent tasks scheduling [32]. Integer PSO is a new variant is proposed; the main objective of this variant is to solve the task scheduling problem in cloud. Integral-PSO is an improved and continuous form of Particle swarm optimization.

V. TASK ANALYSIS ALGORITHMS

A. Genetic Algorithm

The most transformative algorithms are the genetic algorithms which are dependent on the concept of natural transformation [14]. This genetic algorithm is promptly emerging in the field of artificial intelligence [6], [3]. It works on the processing of every task as shown in Fig. 3, in which the quality of each task is being processed according to the user requirements unless the user is being satisfied [9]. Darwin's theory introduced the idea of "Survival of the fittest" which basically processes the tasks according to their

allocation to the resources on the base of their fitness value functions [12]. It doesn't process that task as whole rather it evaluates each parameter of that task on basis of fitness value [4] [20]. The generic terms of this algorithm are as follow:

a) Initial Population

In this algorithm there are several number of individuals which operates the tasks in an iterative way and so several number of solutions are being fixed up, such solutions are termed as populations, in every specific iteration. In that population every solution is termed as chromosome. Ten chromosomes are being selected from that population [5]. From this an initial population, ten chromosomes are selected unsystematically [6].

b) Fitness Function

The basic purpose of this function is quality evaluation of each individual in population while depending upon approach of optimization. It is more often dependent on deadline but in few cases it is dependent on the budget constraints [7].

c) Selection

In this process an operator is used known as proportional selection which is used for evaluating the probability and fitness between two algorithms. It identifies that either selected probability or next groups are proportional to fitness of each individual [10].

d) Crossover

Purpose of this process is the selection of best fitted pair of individuals for crossing over and this is not done without the usage of an operator known as single-point crossover operator. The benefit of crossing over is that both sides' portions can be exchanged [6], [10].

e) Mutation

New individuals are not generated in easy way for that purpose; some of gene locus is being substituted by other gene locus values and it is done in the coding series. A very small value (0.05) is chosen as mutation probability [11].

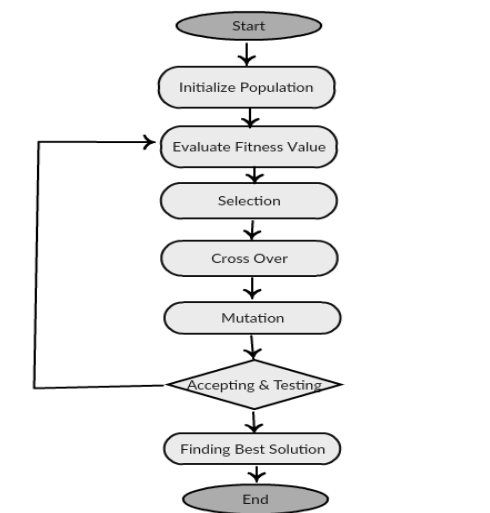


Fig. 3. Working of Genetic Algorithms [5].

B. Greedy Algorithm

Greedy Algorithm is used for solving the problems by making decisions considered best in that particular situation. Working of Greedy Algorithm is explained in Fig. 4. Optimization issues can be easily solved with the help of this algorithm. Though some problems do not seem easy enough with efficient solutions but they are being solved with the help of greedy algorithm with the finest solutions [11]. There are some deviations to the greedy algorithm:

- Pure greedy algorithms
- Orthogonal greedy algorithms
- Relaxed greedy algorithms

When some agitations occur such as bad weather and so on, few constraints in above model are being effected due to which the entire schedule become totally unworkable. The basic purpose of this algorithm was to overcome such problems in each and every step and makes the finest decisions. Its main aim was to get the finest solution and keep on working that schedule unless all the problems are being solved [13]. Due to this optimization of the large problem was divided into small size problems and this helped in identification of solutions in less time [12]. Basic working of Greedy Algorithm is as follows:

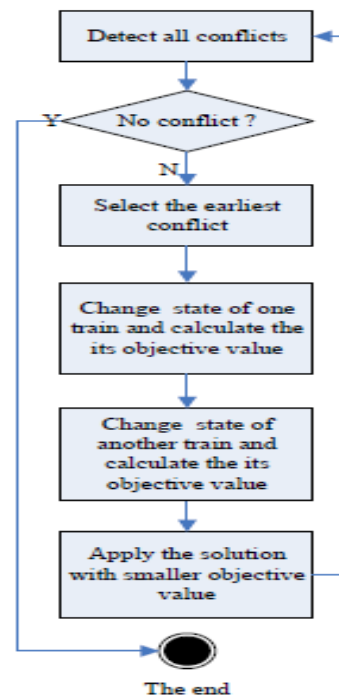


Fig. 4. Working of Greedy Algorithms [12].

Some standards of Greedy Algorithm are as follows [12]:

a) Kruskal's Minimum Spanning Tree (MST)

In this MST is created by selecting edges not collectively but individually. The greedy choice is always selecting the edge of lightest weight because it wouldn't create a cycle in MST.

b) Prim's Minimum Spanning Tree

MST is being created again in it but we manage two sets: set of vertices which are already being added up in MST and the set of vertices which is not added yet. Those edges are selected which are less in weight [11].

c) Dijkstra's Shortest Path

It is very similar to Prim's algorithm. In this the shortest path tree is being built up by every single edge. We manage two sets: set of vertices which are already being added up in MST and the set of vertices which is not added yet [18]. Greedy choice is selection of the smallest weight path.

d) Huffman Coding

Loss-less compression technique is considered as the base of this algorithm. It allocates variable length bit codes to different characters. The Greedy Choice is to assign least bit length code to the most frequent character [11], [12].

C. Priority-based Job Scheduling Algorithm

In Cloud computing, an innovative approach to deal with programming work is presented by Shamsollah Ghanbari and Mohamed Othman by using mathematical measurements [19]. The significance of the job for programming is considered by this algorithm and is called the Algorithm for priority based job scheduling algorithm "PJSC". It is centred as the multiplicative standards decision-making model. There are three levels of priorities in this algorithm that are programming level, resource level and work level which is shown in Fig. 5.

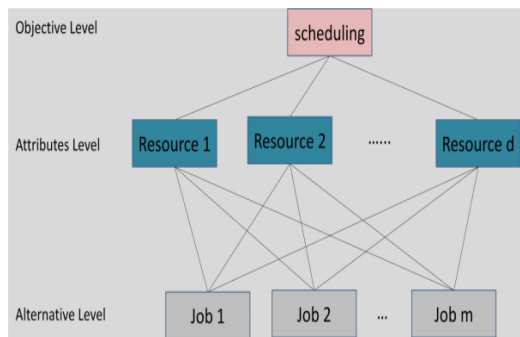


Fig. 5. Three level of priority based scheduling [20].

In this algorithm jobs set are taken as $J = \{J_1, J_2, J_3, J, \dots, J_m\}$ that demands assets in a cloud atmosphere and resources set are taken as $C = \{C_1, C_2, C_3, C_4, \dots, C_d\}$ that is presented in cloud atmosphere as input where $(d \ll m)$. Each job set demands a resource with the required priority. Priority of different jobs set is compared independently [28]. Each job is allocated a resource with the specified priority. Hence, the correlation networks of each activity/job set are computed according to the prospects of retrieving the resources, and the matrix of comparison of the resources is also computed. Now the normal matrix of all jobs with the name Δ is calculated by comparing the each of the job set matrices and priority paths are also calculated [21]. Then the normal resources matrix is calculated, and the name of the matrix is given as γ .

Now the PVS (priority vector of S) is calculated in this algorithm and S is stated as a set of jobs. The matrix Δ is multiplied with the matrix γ which is resulted as PVS. Now the highest ranked job is chosen, and resource is allocated to that job. Job's list is upgraded with the time, and the programming procedure proceeds until the point that all jobs are planned in a suitable resource [33]. The trial/experiments come about show that the calculation of the algorithm has the rational complexity. There are additionally a few issues identified with this calculation, for example, completion time, consistency and complexity [19], [21].

D. Round Robin

The round robin is a simple example of load balancing technique. A round robin technique was designed to divide scheduling time among all scheduled task equally, in which all tasks get in queue list, and each task gets an equally small unit of time as clearly explained in Fig. 6. The major concern of RR is to focus on dividing load to all resources equally [14]. A cyclic approach is applied in round robin. The scheduler picked a task and assigned to the controller and after time expires of the first task then move to next task [17]. This is the cyclic approach in which all task assigned to the controller at least once and then scheduler again pick up the first task again.

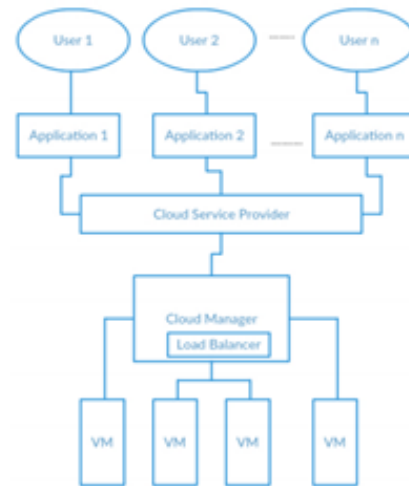


Fig. 6. Load Balancer Round Robin [10].

The load balancing and response time are much better compared to other algorithms. The working of Round Robin in cloud environment is same as round robin in process scheduling [15]. In round robin, each task has an equal opportunity to be chosen [28].

E. Min-Min Algorithm

This scheduling technique works on strategy in which task has minimum execution time is selected for all task. This algorithm starts when a set of all jobs are not assigned and continue to execute until the whole set of the job is empty [16]. In Min-Min jobs having the greater time or long task may not be considered first and the task having greater time will always follow the short job. In this algorithm completion

time of all tasks is computed then job having smallest completion time is scheduled on resource [16].

The formula for calculating completion is as follows:

$$T_{finish} = T_{exe} + T_{start} \quad (1)$$

T_{finish} is finish time, T_{exe} is expected execution time, and T_{start} is starting time of the task

$$T_{comp} = \max(T_{finish} = T_{exe} + T_{start}) \quad (2)$$

The Job which is mapped to resource first after its completion is deleted and the process repeated until all task is mapped. In-Min causes all set of tasks executed get a longer time and unbalanced load even in some cases long task cannot be considered [23].

F. Particle Swarm Optimization (PSO)

This is meta-heuristic population-based algorithm exhilarated by social manners of fish schooling and bird flocking [29]. The algorithm contains set of particles, and each particle depicts a solution for the problem in given search space which is then used to approach convenient solutions [19]. This algorithm is initialized by a set of random particles and then finding a best solution in problem space. In PSO we use iteration to find out each particles position which is referred as Personal best P_b achieved by particle i and global best P_g . Position found by neighbour particle i . Equations to update particles velocity and position after finding both global personal values are [22].

Equation (1):

$$v_i^{t+1} = \omega v_i^t + c1r1(p_i^t - x_i^t) + c2r2(p_g^t - x_i^t) \quad (3)$$

$$x_i^{t+1} = x_i^t + v_i^{t+1} \quad (4)$$

Where particles velocity and position in dimension d of the i^{th} particle are represented by v_i^t and x_i^t respectively. The PSO parameter ω , $c1$ and $c2$ should be considered properly to increase the efficiency of algorithms. This helps in finding out best solution in short computing time [30].

Once all the tasks are queued in a cloud environment, the optimization algorithm is then used to calculate minimum waiting time values of all jobs. These minimum values used for correct order of task which in return minimize the overall waiting time [31]. After getting bested optimal order of task, Queue generated algorithm is applied to find the threshold then dispatches a task to this queue. The scheduler then schedules a task to a suitable resource (server) [35].

The main objective of PSO is to allocate a user request to a suitable resource [33], [35]. To schedule a task on cloud environment efficiently, the task scheduling process requires such optimal algorithm that takes a task and resources into consideration. The PSO algorithm considers both resource and task and helps in keeping the resource as busy as possible and minimizing the processing time of the task [31], [33].

G. First Come First Serve Algorithm

The most fundamental and simplest techniques which uses the task arrival time to schedule the task on cloud environment. The task will be schedule and executed depending on which task has arrived first in queue.it totally depends on arrival time and doesn't consider any other parameter. The tasks will be scheduled by selecting correct order of jobs. The task or user request which comes first to data centre will be assigned to VM first for execution. The data centre controller checks for free Virtual machine and then assign task to that VM then remove that task from queue.

If four task arrives on cloud environment having three virtual machines then FCFS scheduler will schedule three task on VM parallel leaving one task until one VM becomes free for first schedule. For second schedule if task 4 has Childs then childs can't be executed until their parent executed. When task 4 is executing on VM then two other VM remains idle which cause the less utilization of resources.

VI. COMPARISON OF EXISTING TASK SCHEDULING ALGORITHMS

Task scheduling algorithms that we discussed earlier are compared in by their methodology.

TABLE. I. COMPARISON OF VARIOUS TASK SCHEDULING ALGORITHMS

Algorithms	Methodology	Parameters	Merits	Demerits
Genetic Algorithm	Genetic algorithm wants a depiction of the solution domain and suitable function to estimate the solution domain.	1.population size 2. Crossover probability 3.mutation probability	1. It can solve the mathematical problems and financial problems more accurately. 2. Easy to understand the concepts. 3. Some applications required less time for processing.	1. This algorithm work very slowly. 2. This algorithm cannot find the exact solutions. 3. Method of selection should be appropriate.
Greedy Algorithm	This algorithm tries to find the global optimum by following the problem-solving heuristic approach for choosing every step.	1. Parameter μ 2. Domain D 3. Population n	1. Easy to implement. 2. This algorithm needs fewer resources. 3. Execution is very fast. 4. Scheduling is very fast.	1. Global optimization solution is not fulfilled by this algorithm. 2. Very difficult to make changes in parameters.
Priority-Based Job Scheduling Algorithm	Dependency mode	1. Priority to each queue	1. Priority of the process increases with the increases in the time. 2. Easy to use and user-friendly. 3. Best for the applications which require time and	1. Jobs having lowest priority will be lost when the system crashes. 2. Starvation for resources they need.

			resources. 4. Less finish time	
Round Robin	The algorithm works on cyclic approach in which each task has equal chance to be chosen and has an equally small unit of time for execution	1. Arrival time 2. Time slice	1. Response time is good. 2. Load is balanced. 3. less complex	1. Pre-emption causes the process out once time slice expires
Particle Swarm Optimization	The algorithms use population to find the optimal minimum values that help in creating a correct order of tasks and schedule task to a suitable resource	1. Inertia, 2. C1, C2 constants	1. high utilization of resources, finding the optimal solution, minimizing processing time	1. Slow convergence speed if search space is large
Min-Min Algorithm	This algorithm works on strategy in which task having minimum execution time is selected for all task	1. Makespan	1. Better makespan	1. load imbalance 2. poor QoS
First Come First Serve	This algorithm manages the task scheduling with FIFO queue. Task which comes first will be executed first on VM	1. Arrival time	1. Simple and fast execution	1. Task scheduling is based on arrival time, doesn't consider any other criteria Less utilization of VM

VII. DISCUSSION

Task Scheduling is one of the biggest challenges in Cloud Computing. The principle motive of task scheduling is to distribute the incoming tasks from users to the available virtual machines keeping in mind the different parameters Load Balancing, execution time, load balance, Quality of service, performance, response time and fairness resource allocation in which task can be executed. Some algorithms consider only load balance while some consider response time. As most algorithms works with one or two parameters, due to which good result can't be achieved effectively. Better results can be produced by coupling more scheduling metrics to generate one efficient algorithm as an enhancement but this can be little bit complex.

VIII. CONCLUSION

Efficient scheduling algorithm can yield more desirable services to users and increase the performance provided by cloud environment. The main objective of task scheduling in cloud environment is to reduce the execution time of tasks and to maximize the resource utilization. In this paper, a study related to different existing task scheduling algorithms in a cloud environment has been presented. A short description of each algorithm methodology has been presented and most algorithms consider on one or two parameters. More satisfactory results can be achieved by adding more metrics to existing algorithms. Table I is based on different scheduling parameters such as execution time, load balance, Quality of service, performance, response time and makespan. The major problem in task scheduling is load balancing, response time, resource utilization and memory storage. Efficient scheduling algorithm can be achieved by combining different parameters to existing algorithms which will improve their overall performance of cloud environment.

REFERENCES

- [1] S. Sonia, "Task Scheduling in Cloud Computing," International Journal of Advanced Research in Computer Engineering & Technology (IJARECT), vol. 4, no. 6, June 2015.
- [2] D. a. S. N. R. Nallakumar, "A Survey of Task Scheduling Methods in Cloud Computing," International Journal of Computer Science and Engineering (JCSE), vol. 2, no. 10, 31 Oct 2014.
- [3] P. H. Srimathi, "Survey and Analysis of Task Scheduling in Cloud Environment," Indian Journal of Science and Technology, vol. 9(37), October 2016.
- [4] M. B. Yougita Chawla, "A study on Scheduling Methods in Cloud Computing," International Journal of Emerging Trends & Technology in Computer Science (IJETTCS), vol. 1, no. 3, October 2012.
- [5] F. A. O. Safwar A. Hamad, "Genetic-Based Task Scheduling Algorithm in Cloud Computing Environment," International Journal of Advanced Computer Science and Applications (IJACSA), vol. 7, no. 4, 2016.
- [6] T. Y. K. J. K. K. a. J. S. L. S. H. Jang, "The study of genetic algorithm-based task scheduling for cloud computing," International Journal of Control and Automation, vol. 5, pp. 157-162, 2012.
- [7] T. G. a. A. Agrawal, "Host Scheduling Algorithm Using Genetic Algorithm In Cloud Computing Environment," International Journal of Research in Engineering & Technology (IJRET), vol. 1, 2013.
- [8] R. R. a. R. N. C. R. Buyya, "Modeling and simulation of scalable Cloud computing environments and the CloudSim toolkit: Challenges and opportunities," in High Performance Computing & Simulation, 2009. HPCS'09. International Conference on, 2009.
- [9] B. M. P. D. . J. K. M. a. S. D. K. Dasguptaa, "A Genetic Algorithm (GA) based Load Balancing Strategy for Cloud Computing," International Conference on Computational Intelligence: Modeling Techniques and Applications (CIMTA), 2013.
- [10] S. A. H. a. A. Omara, "Genetic-Based Task Scheduling Algorithm in Cloud Computing Environment," (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 7, no. 4, 2016.
- [11] A. S. a. M. V. S. A. Malik, "Greedy Algorithms," International Journal of Scientific and Research Publications, vol. Volume 3, no. Issue 8, August 2013.
- [12] Z. HE, "Research on Improved Greedy Algorithm for Train Rescheduling," Seventh International Conference on Computational Intelligence and Security, 2011.
- [13] V. K. Z. B. S. B. N. Krivsha, "Greedy algorithms for granular computing problems in spatial granulation technique," XIIth International Symposium «Intelligent Systems», INTELS'16, 5-7 October 2016, Moscow, Russia.
- [14] M. Ann Arbor, "J.H. Adaptation in Natural and Artificial Systems," in University of Michigan Press, Holland.
- [15] J. P. a. X. F. Gaochao Xu, "A load Balancing Model Based on Cloud Partitioning for the Public Cloud," TSINGHUA SCIENCE AND TECHNOLOGY, vol. 18, no. 1, pp. 34-39, February 2013.
- [16] J. L. a. J. X. Gang Lui, "An Improved Min-Min Algorithm in Cloud Computing".
- [17] E. a. GibetTaniHicham, "Smarter Round Robin Scheduling Algorithm for Cloud Computing and Big Data," Journal of Data Mining and Humanities, 23 January 2017.
- [18] M. Ann Arbor, "J.H. Adaptation in Natural and Artificial Systems;," University of Michigan Press.
- [19] S. a. M. Othamn, "A Priority Based Job Scheduling Algorithm in Cloud Computing," International Conference on Advances Science and Contemporary Engineering, 2012.
- [20] W. H. Z. J. Yin H., "An Improved Genetic Algorithm with Limited Iteration for Grid Scheduling," Yin H., Wu H., Zhou J., "An Improved

- Genetic Algorithm with Limited Iteration for Grid Scheduling", IEEE Sixth International Conference on Grid and Cooperative Computing, 2007. GCC 2007, Los Alamitos, CA.
- [21] S. P. a. U.Bhoi, "Priority Based Job Scheduling Techniques in Cloud Computing: A Systematic Review," INTERNATIONAL JOURNAL OF SCIENTIFIC & TECHNOLOGY RESEARCH , vol. 2, no. 11, NOVEMBER 2013 .
- [22] S. B. A. A. E. a. A. M. Hicham Bin Alla, "A Novel Architecture with Dynamic Queues Based on Fuzzy Logic and Partical Swarm Optimization Algorithm for Task Scheduling in Cloud Computing," Advances in Ubiquitous Networking 2, Lecture Notes in Electrical Engineering, 2017.
- [23] A. I. E.-D. M. F. A.-r. El-Sayed T. El-kenawy, "Extended Max-Min Scheduling Using Petri Net and Load Balancing," International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231-2307, vol. 2, no. 4, September 2012.
- [24] D. B. J. K. J. B. Shalmali Ambike, "An Optimistic Differentiated Job Scheduling System for Cloud Computing," International Journal of Engineering Research and Applications (IJERA) , vol. 2, no. 2, pp. 2248-9622.
- [25] M. J. ., M. B. S. a. M. K. T. Arash Ghorbannia Delavar, "'RSDC (RELIABLE SCHEDULING DISTRIBUTED IN CLOUD COMPUTING)";" International Journal of Computer Science, Engineering and Applications (IJCEA) , vol. 2, no. 3, June 2012.
- [26] ". Saeed Parsa and Reza Entezari-Maleki, " RASA: A New Task Scheduling Algorithm in Grid Environment," World Applied Sciences Journal 7 (Special Issue of Computer & IT), pp. 152-160, 2009.
- [27] Mrs.S.Selvarani1 and D. Sadhasivam, "improved cost-based algorithm for task scheduling in Cloud computing," IEEE, 2010.
- [28] A. I. E.-D. M. F. A.-r. El-Sayed T. El-kenawy, "'A Priority based Job Scheduling Algorithm in Cloud Computing," International Conference on Advances Science and Contemporary Engineering 2012 (ICASCE 2012).
- [29] M. S. S. Kalra, "A review of metaheuristic scheduling techniques in cloud computing Egypt. Inf. J. 16(3), 275–295 (2015). Elsevier".
- [30] J. J. e. a. Durillo, "Multi-objective particle swarm optimizers: An experimental comparison," International Conference on Evolutionary Multi-Criterion Optimization. Springer Berlin Heidelberg, 2009.
- [31] A. a. S. K. Verma, "Bi-criteria priority based particle swarm optimization workflow scheduling algorithm for cloud," Engineering and Computational Sciences (RAECS), 2014 Recent Advances in. IEEE, , 2014.
- [32] Y. L. Y. L. X. Dai, "A task scheduling algorithm based on genetic algorithm and ant colony optimization algorithm with multi-QoS constraints in cloud computing," In: 7th International Conference on Intelligent Human-Machine Systems and Cybernetics, pp. 428-431, IEEE (2015).
- [33] A. R. M. Beegom, "Beegom, A., Rajasree, M.: A particle swarm optimization based pareto optimal task scheduling in cloud computing," Lecture Notes in Computer Science, p. 79–86, Springer (2014).
- [34] P. Salot, "A survey of various scheduling algorithm in cloud computing environment," International Journal of Research in Engineering and Technology, vol. 2, no. 2, pp. 131-135., 2013.
- [35] "Parallel Workloads Archive: NASA Ames iPSC/860. http://www.cs.huji.ac.il/labs/parallel/workload/l_nasa_ipsc/."

Technical and Perceived Usability Issues in Arabic Educational Websites

Mohamed Benaida, Abdallah Namoun
Faculty of Computer and Information Systems
Islamic University of Madinah
Medina, Saudi Arabia

Abstract—Educational websites are often used as effective communication mediums to provide useful information for students and course instructors. The current study explores the perceived usability of three top-ranked Arabic educational websites across seven key usability components: effectiveness, efficiency, learnability, memorability, errors, satisfaction, and content. Moreover, the study also identifies the key technical and usability issues that currently exist within Arabic educational websites. A two-phase process encompassing automated tools and user testing was adopted to evaluate the technical performance and student acceptance of Arabic educational websites. In the automatic evaluation, two tools, namely, Web Page Analyser and GTMetrix, assessed the websites against a number of well-known performance guidelines and criteria. The student evaluation entailed 150 students completing three interaction tasks and evaluating the sites using the CSUQ questionnaire. The findings indicate that Arabic educational websites suffered from various technical issues, such as a high number of HTML objects and their large size and, consequently, slow loading speed. Moreover, the websites failed to satisfy all usability components, and students rated them negatively. Relevant guidelines for the effective design of Arabic educational websites are also discussed in this paper.

Keywords—Arabic educational websites; perceived usability; automatic evaluation; student perception

I. INTRODUCTION

Three major North African countries include Libya, Algeria, and Morocco. The majority of these countries' residents are young people, and the number of Internet users in these North African countries has risen dramatically from 300,000 in year 2000 to 47,721,000 users in the year 2017 [44].

Internet users in general and students in North Africa in particular face many barriers in gaining high-speed Internet access. The main problems that discourage students, for example in Algeria, is a slow internet speed due to the limited reach of AT's fixed-line network and the inflation in the cost of Internet usage. In addition, the third generation (3G) connectivity has been in operation only recently, starting in December 2013 [1].

Research efforts investigating the design of Arabic websites are still few and so are the design guidelines for creating effective Arabic websites. However, in 2017 there were more than 185 million Arab users connected to the Internet, and Arabic is the third-used language on the Internet after English and Chinese [44]. This necessitates the need to carry out

research studies exploring the usability and design of Arabic websites.

This research makes three distinct contributions. Firstly, it adds to the existing body of knowledge related to the usability of Arabic web user interfaces, which is still undeveloped. Secondly, it gauges students' overall perception of the current usability of Arabic educational websites and identifies the core technical and usability challenges from which these websites suffer. Thirdly, it proposes tentative recommendations to guide the effective design of Arabic educational websites.

The websites of three prominent universities in North Africa are examined in this paper, including Technology Houari Boumediene (USTHB) from Algeria, the University of Mohammed Premier in Oujda from Morocco, and the University of Benghazi from Libya. The remainder of this paper is divided into five sections. Section II reviews related studies, Section III reports on the experimental procedures and instruments, Section IV summarizes the main results, and Section V discusses the implications of the findings in respect to the design of Arabic websites.

II. RELATED WORK

A. Usability and Usability Components

Website usability has become more important than ever before to the success of a web design. Moreover, enabling users to find what they are looking for effortlessly and rapidly is considered to be a key ingredient in websites' success. Historically, the term 'usability' replaced the term 'user-friendly' in the 1980s [2]. Essentially, usability can be defined as the ease of use and acceptability of a system for a particular category of users carrying out specific tasks [3]. Usability, therefore, is concerned with achieving a specific goal and reducing the frustration that occurs when a situation hinders or stops someone from reaching their goal [4]. Research studies have demonstrated that various design factors, such as images, content, color, and logos, may directly influence users' satisfaction [5]-[7]. Indeed, the main usability components, including efficiency, effectiveness, learnability, and satisfaction, are intertwined [8].

Numerous studies have examined the usability of websites in different application areas, such as education, banking, healthcare, and mobile settings [8]-[12], [5]. Sometimes, websites fail to retain their users as a result of poor application of usability principles, which could be an indicator of poor knowledge and use of web design theories [14]. For instance,

Juristo, Moreno, and Sanchez-Segura [13] emphasized that web designers should possess sufficient knowledge about psychology, ergonomics, and linguistics theories and principles that constitute the cornerstone of effective website creation [13]. Masood and Thigambaram [20] confirmed this finding in their observation of how children, aged 4-5 years, are effected by the usability of mobile applications. Their data analysis, collected through eye tracking technology, showed that there is a mental model gap between website designers and preschoolers, and they advised developers to adopt a user-centered design methodology when creating educational applications.

B. Evaluation of Educational Websites

There exists various usability methods to evaluate educational websites, including questionnaires, usability inspection, heuristic evaluation, field study, analysis of site usage logs, formal usability testing, and focus groups [5]. Tan et al. [25] compared the effectiveness of heuristic analysis and user testing during the evaluation of four commercial websites. Their results revealed that those two methods complement each other, as they tackle different usability problems. Heuristic analysis is generally associated with the initial stage of the design process whilst user testing is related to the latter stages of process design. However, [24] argued for the need to create a compound index to measure the perceived usability by combining major metrics such as task time, task completion rate, error rates, and satisfaction ratings.

In [22], the authors employed Shackel’s usability model to develop a survey to measure the usability experience of final-year students on e-learning websites. Indeed, prior experience was shown not to have any influence on the overall perceived usability. Instead, other usability attributes such as effectiveness, learnability, flexibility, and attitude were considered critical for achieving a pleasurable experience with e-learning webpages. Likewise, [26] analyzed the usability of Namik Kemal University’s (NKU) website and found that four factors, namely attractiveness, helpfulness, efficiency, and learnability, are positively correlated with website usability and overall user perceptions. Both gender and web experience had moderating effects on users’ usability perceptions. However, Kiget [18] performed a case study to assess various factors, including learnability, user-friendliness, culture, technological infrastructure, gender, and policy, that may drive the usability of Kenyan e-learning websites. Only learnability was identified as a key contributor to the usability of e-learning systems.

Sengel [9] discussed the level of usability pertaining to a particular university website in Turkey. Results showed that the majority of the users found the website easy to use, and the website proved to be a very useful source of information related to the university. In contrast, Sengel and Oncu [19] showed that the content provided by the Uludag University website needed more attention and should be reviewed and updated regularly. Moreover, gender was found to influence the perceived usability of the website, with females holding a more positive view about the website. Furthermore, [23] conducted a statistical study to assess the perceived usability and accessibility levels of three educational websites. Task completion times was found to correlate with the overall level of the user’s satisfaction. In [38], the researchers discussed the

participation of students in e-learning systems by focusing on various engagement methods such as participation of students in forums, blogs, and wikis. The results showed that the forums and wikis were more beneficial to the students’ learning process than were the blogs.

In [17], the researchers examined multilingual websites and argued that these websites must be designed in a way that satisfies all users at the local and international levels. The authors proposed a framework and set of guidelines, including use of appropriate color and layouts, to be followed while designing Arabic-English websites. Moreover, [21] emphasized the need for designers to consider the cultural context when creating multi-lingual websites. In the context of our research, timid work efforts, such as [15], [16], [33], have been conducted in respect to establishing design recommendations and models serving the design of Arabic websites.

III. EXPERIMENTAL SETUP

A. Selected University Websites

In this study, the top ranked university in the year 2018 in three North African countries was considered, namely, Algeria, Morocco and Libya, for this experiment. This selection relied on Webometrics ranking [45], which is carried out by an independent research group in Europe that specializes in the analysis of web presence of universities.

Webometrics scores university websites on four weighted key indicators mainly, web presence (5%), visibility (50%), openness (10%) and excellence (35%) [46].

Web presence refers to the number of pages and rich files. Visibility refers to the number of backlinks from external networks to the university website. Transparency refers to the number of citations about the university in Google Scholar. Excellence refers to the number of top 10% research papers cited in Scimago Journals within the last 5 years.

TABLE I. WEB RANKING OF SELECTED UNIVERSITY WEBSITES ACCORDING TO WEBOMETRICS

Name of University (Country)	University of Sciences and Technology Houari Boumediene (USTHB - Algeria)	Mohamed Premier University Oujda (MPOU - Morocco)	University of Benghazi (UB – Libya)
World Rank	2250	2345	4030
Arab Rank	37	47	134
Country Rank	1	2	1
Indicator one:	1522	3646	9105
Indicator two:	6448	10703	12072
Indicator three:	2414	2385	4270
Indicator four:	1783	1106	3168

The results of Webometrics analysis, as depicted in Table I, showed disparate ranking of the three websites in the Arab world, with the University of Sciences and Technology Houari Boumediene (USTHB) Website placed first (Fig. 1), followed by Mohamed Premier University Oujda (MPUO) Website (Fig. 2), and finally University of Benghazi (UB) Website (Fig. 3).

Table II compares the three websites across a range of design features and media elements including background color, languages used, font type and size used, and menus. Generally, all websites employ black, white and green text on a whitish background.

TABLE II. THE MAIN DESIGN FEATURES OF UNIVERSITY WEBSITES, KEY CROSS-DIFFERENCES ARE HIGHLIGHTED IN GGEY

	University of Sciences and Technology Houari Boumediene	Mohamed Premier University Oujda	University of Benghazi
	(USTHB Website)	(UMPO Website)	(UB Website)
Web Ranking using Webometrics in the Arab World	37	47	134
Background colour	White and light green	White and light grey	White and light grey
Languages available	Arabic, English, French	Arabic, English, French, Amazigh	Arabic, English
Search engine	Not available	Top right	Top left
Text colours used	White, Black, Green	White, Black	White, Black and Grey
Font type and size (Title)	Arial, 21 PX	Titillium Web, 23PX	GE SS TWO, 48 PX Bold
Font type and size (Text)	Arial, 12 PX	Titillium Web, 13PX	GE SS TWO, 12 PX Bold
Menus	One at the top and another on the right	Two at the top	One at the top
Pictures on the home page	2 small pictures	5 big and 25 small pictures	6 big and 11 small pictures
Logo	Top left	Top right	Top Middle
Languages mixed on same page	No	Yes (Arabic and French mixed)	No
Layout	Traditional, three column structure	Traditional, grid design	Traditional, two column structure
Social network links	Not available	Facebook, Twitter, YouTube	Facebook, Twitter, YouTube, Instagram
Animated pictures	Not available	Top of the page	Top of the page

The text font employed on the websites varies between Arial 12, Titillium Web 13, and GE SS TWO 12. Each website offers its content in at least two languages, Arabic and English. Remarkably, the UMPO website mix two languages, particularly Arabic with French, on the same page. Moreover, all websites use images on their home pages and provide a search engine for searching the content, apart from USTHB website. All three websites organize their content and media features within a column or grid-structured design. In respect to social exposure, the UMPO and UB websites enable their students to connect to social networking sites from their pages.

The following screenshots represent the home page of the three educational websites examined in this research study.



Fig. 1. Website One: University of Sciences and Technology Houari Boumediene Website [47].

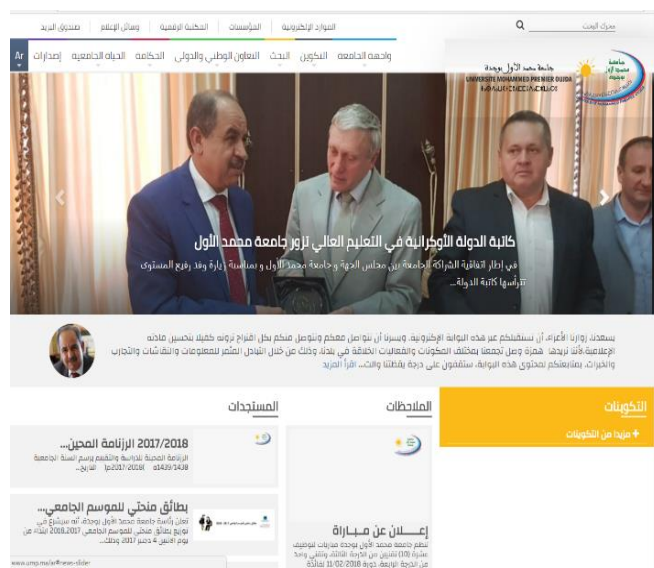


Fig. 2. Website Two: Mohamed Premier University Oujda Website [48].



Fig. 3. Website Three: University of Benghazi website [49].

B. Research Methodology

This research employed a two-step evaluation procedure consisting of two differing yet complimentary approaches, a technical performance evaluation and user testing as depicted in Fig. 4. The technical assessment employs freemium website performance tools to check compliance with web technical guidelines and best practices that focus on optimizing various page aspects such as the page size, items number, and load time. Previous research showed that browsing frustration is strongly affected by slow downloads and connection drops [27]. However, the subjective evaluation engages actual students into the evaluation process and considers their opinions and feelings about the perceived usability of educational websites. Overall, this procedure empowers the creation of highly dedicated research recommendations for designing effective Arabic educational websites.

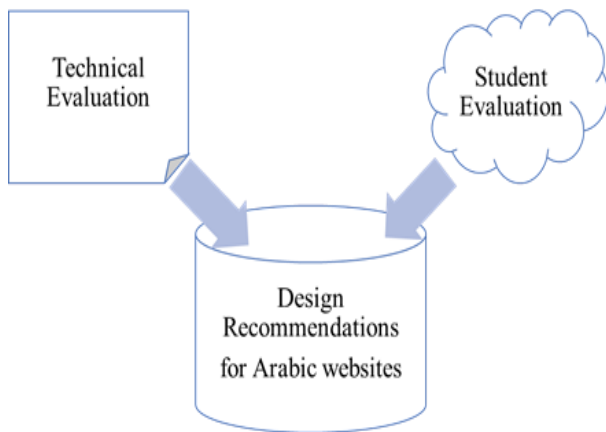


Fig. 4. Evaluation procedure of university websites.

C. Automated Technical Evaluation

The technical evaluation aims to assess web pages' performance by investigating a number of characteristics such as load speed of websites, size of pages and build-in features such as, number of objects and size of these objects. To this end, two automated evaluation tools were used to analyse the

performance of our three educational websites. These tools are GTmetrix and Web Page Analyzer, respectively.

Two tables were produced using two website analysis tools, GTmetrix and Web Page Analyzer respectively. GTmetrix produced details about five main performance indicators namely PageSpeed Grade, YSow Grade, Fully loaded time, Total page size, and Total number of requests as listed in Table III. PageSpeed Grade is a measure from Google tools which indicates the degree of overall compliance to best practices, whereas Yahoo YSow grade is calculated using 23 rules, such as minimal use of DOM elements and small scale of images [50] believed to affect web page overall performance. Page speed scores from Google tools showed that none of the university websites surpass the average result from other websites (i.e. less than 71%). Similarly, YSow results showed that only University of Benghazi website exceeded the average score from other websites (more than 69%). When high Internet speed was used, all educational websites loaded fully in less than 6.7 seconds. However, in a more realistic scenario when low Internet speed was used, only the UB website fared quite well (4.9 seconds). Notably, UMPO webpage size totaled approximately 5.05MB (2.86MB is the average recommendation) and has received substantially more HTTP requests (i.e. 108) during a month than the remaining websites.

TABLE. III. WEB PAGES PERFORMANCE RESULTS FROM GTMETRIX [51], GREY CELLS REPRESENT THE ACCECTABLE RESULTS USING CHROME (DESKOP) WITH HIGH INTERNET SPEED (MAIN 38.4 MBPS) AND LOW INTERNET SPEED (MEAN=2.69 MBPS) [52]

	Average Results from Other Websites	USTHB Website	UMPO Website	UB Website
PageSpeed Grade (Higher is better)	71%	D (54%)	F (23%)	D(63%)
YSow Grade (Higher is better)	69%	D (61%)	D (60%)	C (75%)
Fully loaded time (in Second) using high Internet speed (Lower is better)	6.7	4.4	4.9	3.7
Fully loaded time (in Second) using low Internet speed (Lower is better)	6.7	8.7	6.0	4.9
Total page size (MB) (Lower is better)	2.86	1.59	5.05	2.63
Total number of requests (Lower is better)	89	56	108	103

Although UB scored better than all websites, it received a high number of HTTP requests (103) which negatively influences the overall page loading time.

Table IV reports performance results from Web Page Analyzer [53] which given a URL of a website provides an objective performance analysis of the total number of objects, images, CSS, and scripts and size of HTML, images, CSS, and scripts on a web page. Overall results show that all three educational websites do not conform to most of the Web Page Analyzer guidelines with regards to the number and size of website items. However, only two criteria were adhered to namely HTML size and CSS size (except UMPO website). The evident issues, on the other hand, were the use of many objects, images, and scripts that exceed the recommended number and size for websites. Overall, UB website outperformed the other websites in respect to total and items' size, followed by USTHB website and finally UMPO website. This confirms the results of GTmetrix analysis.

TABLE IV. WEB PAGE PERFORMANCE RESULTS FROM WEB PAGE ANALYSER (NUMBERS REPRESENT AVERAGE ON A SINGLE PAGE), GREY CELLS PRESENT THE ACCEPTABLE RESULTS [53].

	WPA Guideline (Should be equal to or less than)	(USTHB Website)	(UMPO Website)	(UB Website)
Total HTML files	1	1	2	3
Total Objects	20	84	55	66
Total Images	reasonable number	65	35	53
Total CSS	1	3	4	1
Total Script	1	15	14	9
Total Items		168	110	132
HTML Size (bytes)	50000	7205	25843	45862
Image Size (bytes)	100000	1572344	1909354	694568
Script Size (bytes)	20000	154195	570609	542042
CSS Size (bytes)	8000	2807	10291	2478
Total Size (bytes)	100000	1736551	2516097	1284950

D. Student Evaluation

User testing refers to evaluating the educational websites with actual users, allowing to collect direct input and feedback from the respective students of these universities. In such kind of testing, the students are requested to carry out a set of realistic tasks and performance is measured using various usability metrics such as completion time, number of clicks and satisfaction scores.

E. Perceived Usability Instrument

Usability is defined by the international standard, ISO9241-11 as: 'the extent to which a product can be used by specified

users to achieve specified goals with effectiveness, efficiency and satisfaction in a specified context of use'. Moreover, Nielsen [54] specified five important components that constitute usability: efficiency, learnability, memorability, errors, and satisfaction.

Accordingly, we have chosen to measure seven components believed to contribute to the perceived usability of educational websites. These components include effectiveness, efficiency, learnability, memorability, errors, satisfaction, and content. Below are their definitions.

- Effectiveness: refers to the rate of successfully completing tasks on the educational websites. This is also referred to as completion rate which is considered a critical performance metric.
- Efficiency: refers to the time taken, in seconds, to successfully complete designated tasks. Therefore, efficiency measures mainly the speed with which a certain task can be accomplished. It is referred to as task time.
- Learnability: refers to the efforts needed to learn to use the functionalities of the educational websites from the first time.
- Memorability: refers to the ability to recall how to use the educational websites on revisits after a period of discontinuity.
- Errors: refers to the number of mistakes or error rate when completing the intended tasks and actions.
- Satisfaction: refers to the liking of the system and pleasurable experience when using the educational websites. This is usually captured at the end of the test using a post study questionnaire.
- Content: refers to the quality and clearness of information, such as text, graphics or videos, provided on the educational websites.

Indeed, questionnaires have been used as a primary method to gather user views and satisfaction about web design [5]. In this study, we have opted to use the IBM Computer System Usability Questionnaire (CSUQ) to quantify the above usability components [28]. The CSUQ usability questionnaire has been shown to yield accurate results and requires only between 12 and 14 users to run [29]. Moreover, the CSUQ questionnaire is well known for its high reliability [28]. Overall, the CSUQ encompasses a total of 19 questions that need to be rated on a 7-point Likert scale. The students therefore had to rate the questions from 1= strongly disagree to 7= strongly agree. The rating scale included 'a not applicable' option as well. All questions of the CSUQ were translated to Arabic to enable the students to respond accurately in their first language.

We relied on the above definitions of the seven usability components in order to categorize the CSUQ questions as follows.

TABLE. V. PERCEIVED USABILITY COMPONENTS AND RELEVANT ITEMS FROM CSUQ QUESTIONNAIRE

Perceived Usability Component	Question Item
Effectiveness	<ul style="list-style-type: none"> It was simple to use this educational website I can effectively complete my work using this educational website The information is effective in helping me complete the tasks and scenarios This educational website has all the functions and capabilities I expect it to have
Efficiency	<ul style="list-style-type: none"> I am able to complete my work quickly using this educational website I am able to efficiently complete my work using this educational website I believe I became productive quickly using this educational website
Learnability	<ul style="list-style-type: none"> It was easy to learn to use this educational website
Memorability	<ul style="list-style-type: none"> Whenever I make a mistake using the educational website, I recover easily and quickly
Errors	<ul style="list-style-type: none"> This educational website gives error messages that clearly tell me how to fix problems
Content / Information	<ul style="list-style-type: none"> The information (such as online help, on-screen messages, and other documentation) provided with this educational website is clear The information provided for the educational website is easy to understand The organization of information on this educational website is clear It is easy to find the information I needed
Satisfaction	<ul style="list-style-type: none"> I feel comfortable using this educational website The interface of this educational website is pleasant I like using the interface of this educational website I am satisfied with how easy it is to use this educational website I am satisfied with this educational website

F. Procedure and Tasks

The designated tasks and CQUS Arabic questionnaire were pilot tested with a total of 10 students from the three above universities to ensure smooth execution of the actual research study. However, no major concerns were raised by the participants. The research study was conducted during off peak times to ensure Internet highest speed.

After providing a short introduction about the study and its procedure, students were instructed to carry out the following activities:

- Read the study information sheet and sign a participation consent form.
- Fill out a short background information form about students' gender and experience using their university website.

- Explore freely the different sections of their university website for 5 minutes.
- Carry out three interactive search tasks using their university website.
- Upon completion of all tasks, evaluate seven usability components, listed in Table V, of their university website using the CSUQ Arabic questionnaire.

The information search tasks completed by the students were varied in complexity and included three tasks. The direct URL link to the Arabic version of the university website was provided to the students.

- Task One: Find the year the University was established.
- Task Two: Find names of three research centers established in the University.
- Task Three: Find the different undergraduate degrees offered by the Faculty of Science.

These tasks were similar and were achievable through all three educational websites. These search tasks enabled the students to explore different sections and aspects of their university website. Two main usability metrics were recorded as the interaction unfolded. These metrics comprised of completion rate and task time.

IV. RESULTS

In total 150 undergraduate students took part in this experiment, with 50 students coming from each university. The students were undertaking degrees from different faculties of their University. Overall, 51% of the students were male and 49% were female. Approximately 48% of the students reported using their university website on a regular basis (i.e. a few times a week or more) as shown in Fig. 5.

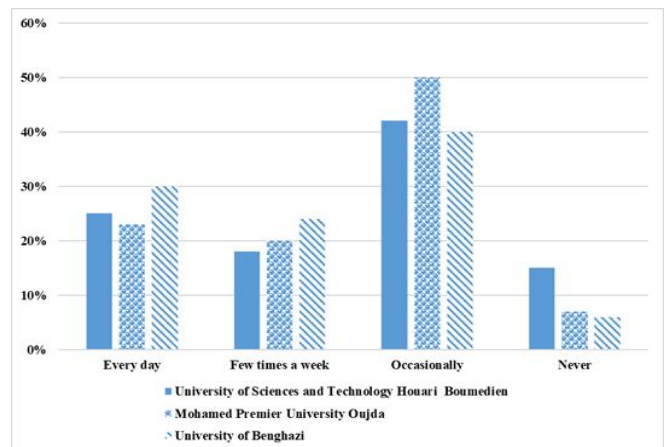


Fig. 5. Frequency of website use.

On average, the results indicated that the students spent the longest time to complete the third task (99.67 seconds) in comparison to the remaining tasks ($p < .001$), as shown in Fig. 6. Similarly, the third task was deemed as the most challenging and resulted in the lowest completion rate in the UMPO website. In the third task, the students were instructed to find the Bachelor programs offered by the Faculty of

Science. In the UMPO website, these Bachelor programs were available in an external page (i.e. not belonging to the same domain of the Arabic version of the University website) and they were written in French which might have confused the students. Similarly, task one resulted in 62% completion rate in the UMPO website due to the unclear navigation structure employed within the website, as shown in Fig. 7.

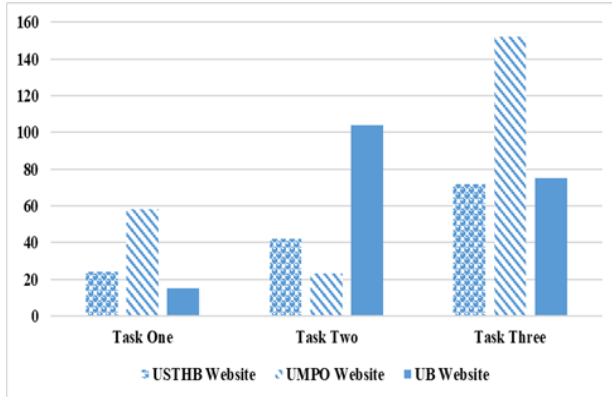


Fig. 6. Task time (seconds).

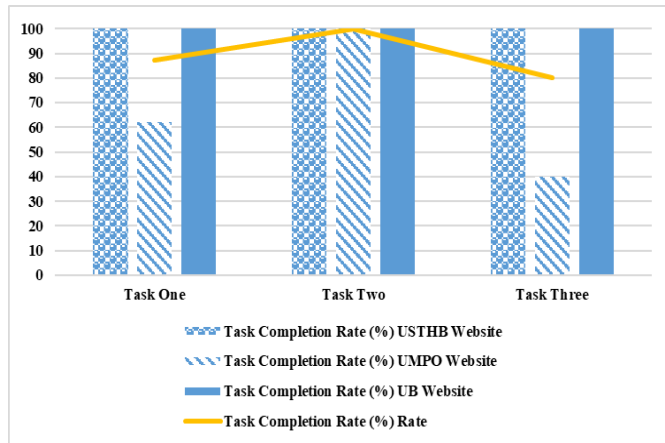


Fig. 7. Task completion rate (%).

A contrast between the three educational websites using their average score of the two performance metrics reveals that UMPO website performed the worst as it took longer time to complete the tasks (77.66 seconds per task) and its students failed to complete all information search tasks (67.33% average completion rate) ($p < .001$), as depicted in Fig. 8. However, the USTHB website outperformed the remaining websites as its students answered all questions correctly within an average speed of 46 seconds per task. This indicates that a good navigational structure is employed in the USTHB website. Finally, UB website trailer behind the USTHB website as it scored the same completion rate but its students took slightly longer to complete the search tasks (65 seconds per task).

Following the completion of search tasks, the students rated various aspects of the perceived usability of their university websites. Cronbach's alpha test showed good reliability of all usability components proposed in Table V ($\alpha = 0.86$). All questions of the CSUQ were rated on a 7-point Likert scale. Table VI summarizes the results of the students' rating of the three educational websites. It shows that the average score of almost all usability components were scored below 4, the acceptable threshold, across the three websites. This indicates low perceived usability and satisfaction by the students towards their university websites.

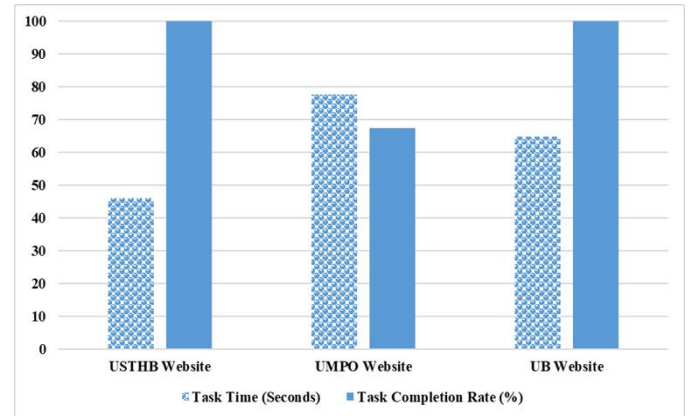


Fig. 8. Performance comparison across the educational websites.

Overall, UB website received a better evaluation than the other websites across most of the usability components ($p < .001$), except efficiency and satisfaction. On the other hand, UMPO website received the worst rating score in respect to four components mainly effectiveness, learnability, content, and errors. The effectiveness mean for all educational websites was 3.01, with UB website receiving the highest average and UMPO receiving the lowest average. This result agrees with the performance results. The efficiency mean did not differ across the three websites (average= 3.13).

On average, the learnability component received the highest rating score (mean=3.75) across all educational websites, with the UB website as the favorite (mean=4.22). This component had the highest score amongst all other usability factors. However, the memorability mean for all websites was quite low (mean=2.52), with the UB website receiving the highest score and USBTH website receiving the lowest score. On average, the quality of content received the second best rating score (mean=3.15) among all components, with the UMPO website scoring the worst (mean=2.98). On the contrary, the websites were rated very poorly for providing error messages to help resolve problems (mean=2.19), with the UMPO website perceived as the worst on this component. Finally, the students did not like the design of their educational websites and were general unsatisfied (mean=2.97).

TABLE. VI. STUDENTS RATING OF PERCEIVED USABILITY COMPONENTS OF THE THREE EDUCATIONAL WEBSITES (LIGHT GREY REPRESENTS WORST RATING; DARK GREY REPRESENTS BEST RATING)

Usability Component	USBTH Website		UMPO Website		UB Website		Mean Score (all websites)	Sig value
	Mean	STD	Mean	STD	Mean	STD		
Effectiveness	2.88	0.82	2.63	0.62	3.51	0.76	3.01	.000
Efficiency	3.13	0.95	2.96	0.56	3.29	0.59	3.13	NS
Learnability	3.76	1.49	3.26	1.01	4.22	1.25	3.75	.000
Memorability	2.1	1.34	2.5	0.94	2.96	1.03	2.52	.001
Content / Information	3.21	0.79	2.89	0.45	3.36	0.58	3.15	.000
Errors	2.1	1.04	1.8	0.77	2.68	1.3	2.19	.001
Satisfaction	2.97	1.02	2.99	0.96	2.95	0.82	2.97	NS
Average perceived usability score	2.88		2.72		3.28			

V. DISCUSSION

To the best of our knowledge, this research study is the first to examine the perceived usability of North African educational websites, with a focus on the top web ranked universities in Morocco, Algeria, and Libya. These websites used Arabic as the main language in addition to at least one secondary language. Typically, educational websites provide useful information for students and instructors and are used as a communication medium with these stakeholders [30]. It is therefore important to ensure high usability of these websites and the use of effective design elements to maximize educational gains. This will, in turn, help achieve the educational institution's goals.

Our study comprised two methods of evaluation, automatic evaluation and empirical evaluation, of the educational websites. Barnes et al. [31] showed that using data from multiple methods of testing yields better results and gives more insights about users' perceptions. In automatic evaluation, a tool checks for whether or not the educational websites adhere to web standards and best practices [32]. In empirical evaluation, actual students are recruited to complete search tasks and provide their views about the perceived usability of these websites [32]. The results of both evaluations indicated that our three Arabic educational websites failed to meet the performance and usability expectations of their students.

The web performance results from Web Page Analyzer and GTMetrix showed weak scores for all educational websites. The major technical problems involved the large size of images and JavaScript within the sites, lack of scaled images, absence of browser caching, the parsing of JavaScript during initial page load, the large number of objects, the lack of compression to optimize transfer size, and serving the same web resources from multiple URLs. These problems negatively impact the page load times. Even using high Internet speeds, page loading results were still below the recommendations. The page loading speed of all educational websites was less than 6.9 seconds. However, Kissmetrics has suggested that around 40% of site users will abandon the website if the loading time

exceeds more than 3 seconds [43]. It is well-documented in the literature that a slow website increases user frustration, affects the overall judgment of the website, and discourage users from returning to the website [34].

The student evaluation encompassed the completion of three information search tasks and the rating of seven usability components, namely effectiveness, efficiency, learnability, memorability, errors, content, and satisfaction, which were derived from the CSUQ questionnaire. The overall rating showed weak usability perceptions by the students across all usability components. The UMPO website obtained the worst rating for its content, possibly as a result of mixing two languages on the same pages. Similarly, it received a poor rating for error handling and prevention, as students were often directed to the French version of the site and experienced difficulty returning back to the Arabic version. Students were also unable to complete all tasks on the UMPO website due to the poorly designed navigation menu. Previous research has emphasized the importance of content and navigation within educational websites for overall student preference of websites [35]. Further, Nielsen suggested paying particular attention to response time, content, and navigation mode [36], and Naidu [37] claimed that the terminology, length of pages, and organization of links affects the search performance. In our view, the poor performance results might have also influenced the perceived usability of the Arabic educational websites.

VI. CONCLUSION AND FUTURE WORK

This paper proposes, based on the findings of our evaluation, design recommendations to guide the development of educational Arabic websites. These recommendations aim at enhancing the performance of web pages and improving the perceived usability of these pages for students.

A. Website Performance Implications

Top-ranked educational websites in three Arab countries were overburdened with the use of excessive web elements, including images, objects, and scripts. Both the number and size of these elements exceeded the existing recommendations. Optimization was absent, which led to weak performance and

prolonged loading times in the educational websites. It is advised to minimize the number of website objects, which in turn reduces the number of HTTP requests. This will expedite the web response and loading time. Moreover, the size of website objects such as images should be maintained at minimum sizes to accommodate slow Internet connections. This design recommendation supports previous findings where users were found to be mainly frustrated by long download times [27]. Lazar et al. [39] confirmed that the time lost during the interaction increases student frustration levels.

B. Perceived Usability Implications

Students' ratings of seven usability components include effectiveness, efficiency, learnability, memorability, error, content, and satisfaction, revealed poor student experience. This student dissatisfaction could result in website abandonment. For instance, SC Chang and FC Tung [40] confirmed that perceived ease of use is major indicator of students' intention to use educational websites. In addition to including relevant and useful content, it is important to design educational websites that are easy to use. This agrees with the suggestions provided in [41].

The current study has initiated research into the perceived usability of Arabic educational websites and raised several concerns regarding the performance and satisfaction of students. Although our study provides initial recommendations, further research is encouraged to better understand the weak student satisfaction towards their Arabic educational websites. Quantitative results indicated some key findings, however these findings would need to be complemented by a qualitative investigation to derive conclusive guidelines and a deeper understanding of students' overall experiences. Moreover, the link between educational websites' performance and student satisfaction needs to be explored in future research. Previous research has shown that longer loading times correlate with user frustration, which could consequently affect overall acceptance and use of websites [27], [39]. Finally, aesthetics might have played a role in framing the perceived usability, as demonstrated in previous research studies [42].

REFERENCES

- [1] S. Chaabna and H. Wang, "Analysis of the State of E-commerce in Algeria," *International Journal of Marketing Studies*, vol. 7, no. 2, 2015.
- [2] N. Bevan, J. Kirakowski and J. Maissel, "What is usability?. Human Aspects in Computing Design and Use of Interactive Systems and Work with Terminals," 4th International Conference on Human Computer Interaction, 1991.
- [3] A. Holzinger, "Usability engineering methods for software developers," *Communications of the ACM*, vol. 48, no. 1, pp. 71-74, 2005.
- [4] V. Mendoza and D.G. Novick, "Usability over time," *Proceedings of the 23rd annual international conference on Design of communication: documenting and designing for pervasive information*, ACM, pp. 151-158, 2005.
- [5] V. Venkatesh, H. Hoehle and R. Aljafari, "A usability evaluation of the Obamacare website," *Government information quarterly*, vol. 31, no. 4, pp. 669-680, 2014.
- [6] I.M. Hanafy and R. Sanad, "Colour preferences according to educational background," *Procedia-Social and Behavioral Sciences*, vol. 205, pp. 437-444, 2015.
- [7] E.Ş. Ekici, C. Yener, N. Camgöz and E.Ahin, "Colour naming," *Optics and Laser Technology* vol. 38, no. 4-6, pp. 466-485, 2006.
- [8] J. Jeng, "Usability Assessment of Academic Digital Libraries: Effectiveness, Efficiency, Satisfaction, and Learnability," *Libri*, vol. 55, pp. 96-121, 2005.
- [9] E. Şengel, "Usability level of a university web site, " 4th International Conference on New Horizons in Education, *Procedia-Social and Behavioral Sciences*, vol. 106, pp. 3246-3252, 2013.
- [10] C. Buchmayer, M. Greil, A.L Hinkl, O.Kaiser-Dolidze and C. Miniberger, "Usability on the Edge: The Implementation of u: cris at the University of Vienna," *Procedia Computer Science*, vol. 33, pp. 103-109, 2014.
- [11] C.A. Gumussoy, "Usability guideline for banking software design," *Computers in Human Behavior*, vol. 62, pp. 277-285, 2016.
- [12] E. Çetin and S. Özdemir, "A Study on an Educational Website's Usability," *Procedia-Social and Behavioral Sciences*, vol. 83, pp. 683-688, 2013.
- [13] N. Juristo, A. M. Moreno and M. I Sanchez-Segura, "Analysing the impact of usability on software design," *Journal of Systems and Software*, vol. 80, no. 9, pp. 1506-1516, 2017.
- [14] S. Qayyum and S. Rafiq, "Website Design Usability Issues Faced by the User in Pakistan," *Computer Engineering and Intelligent Systems*, vol.7, no.9, 2016.
- [15] L. Hasan, "Heuristic evaluation of three Jordanian university websites," *Informatics in Education*, vol. 12, no. 2, pp. 231-251, 2013.
- [16] L. Hasan, "Evaluating the usability of educational websites based on students' preferences of design characteristics," *International Arab Journal of e-Technology*, vol. 3, no. 3, pp. 179-193, 2014.
- [17] M. A. Ababtain and A. R. Khan, "Towards a Framework for Usability of Arabic-English Websites," *Procedia Computer Science*, vol. 109, pp. 1010-1015, 2017.
- [18] N. K. Kiget, G. Wanyembi and A. I. Peters, "Evaluating usability of e-learning systems in universities," *International Journal of Advanced Computer Science and Applications*, vol. 5, no. 8, pp. 97-102, 2014.
- [19] E. Şengel, and S.Öncü, "Conducting preliminary steps to usability testing: investigating the website of Uludağ University," *Procedia-Social and Behavioral Sciences*, vol. 2, no. 2, pp. 890-894, 2010.
- [20] M. Masood and M. Thigambaram, "The usability of mobile applications for pre-schoolers," *Procedia-Social and Behavioral Sciences*, vol. 197, pp. 1818-1826, 2015.
- [21] M. Hillier, "The role of cultural context in multilingual website usability," *Electronic Commerce Research and Applications*, vol. 2, no. 1, 2003.
- [22] M. H. Thowfeek and M. N. A. Salam, "Students' Assessment on the Usability of E-learning Websites," *Procedia-Social and Behavioral Sciences*, vol. 141, pp. 916-922, 2014.
- [23] S. Roy, P. K. Pattnaik and R. Mall, "A quantitative approach to evaluate usability of academic websites based on human perception," *Egyptian Informatics Journal*, vol. 15, no. 3, pp.159-167, 2104.
- [24] G. J. Esmeria and R. R. Seva, "Web Usability: A Literature Review," Presented at the DLSU Research Congress 2017 De La Salle University, Manila, Philippines, 2017.
- [25] W. S. Tan, D. Liu and R. Bishu, "Web evaluation: Heuristic evaluation vs. user testing," *International Journal of Industrial Ergonomics*, vol. 39, no. 4, pp. 621-627, 2009.
- [26] S. A. Mentos and A. H. Turan, "Assessing the usability of university websites: an empirical study on Namik Kemal University, " *TOJET: The Turkish Online Journal of Educational Technology*, vol. 11, no. 3, 2012.
- [27] I. Ceaparu, J. Lazar, K. Bessiere, J. Robinson and B. Shneiderman, "Determining Causes and Severity of End-User Frustration," *International Journal of Human-Computer Interaction*, vol. 17, no. 3, pp. 333-356, 2004.
- [28] R. J. Lewis, "IBM computer usability satisfaction questionnaires: Psychometric evaluation and instructions for use," *International Journal of Human-Computer Interaction*, vol. 7, no. 1, pp. 57-78, 1995.
- [29] T. S. Tullis and J. N. Stetson, "A comparison of questionnaires for assessing website usability," *Usability professional association conference*, 2004.

- [30] H. Yu-chang, "Better educational website interface design: the implications from gender-specific preferences in graduate students," *British Journal of Educational Technology*, vol. 37, no. 2, pp. 233-242, 2006.
- [31] S. J. Barnes and R. T. Vidgen, "Data triangulation and web quality metrics: A case study in e-government," *Information & Management*, vol. 43, no. 6, pp. 767-777, 2006.
- [32] P. Helen and N. Bevan, "The Evaluation of Accessibility, Usability, and User Experience," *The Universal Access Handbook*, C Stepanidis (ed), CRC Press, vol. 1, no. 16, 2009.
- [33] L. Hasan, "The usefulness of user testing methods in identifying problems on university websites," *JISTEM-Journal of Information Systems and Technology Management*, vol. 11, no. 2, pp. 229-256, 2014.
- [34] F. H. Nah, "A study on tolerable waiting time: how long are web users willing to wait," *Behaviour & Information Technology*, vol. 23, no. 3, pp. 153-163, 2004.
- [35] S. Ozkan and R. Koseler, "Multi-dimensional students' evaluation of e-learning systems in the higher education context: An empirical investigation," *Computers & Education*, vol. 53, no. 4, pp. 1285-1296, 2009.
- [36] J. Nielsen, "Designing web usability," Indianapolis, IN: New Riders Publishing, 1999.
- [37] S. Naidu, "Evaluating the usability of educational websites for children," *Usability News*, vol. 7, no. 2, pp. 1-7, 2005.
- [38] T. Miyazoe and T. Anderson, "Learning outcomes and students' perceptions of online writing: Simultaneous implementation of a forum, blog, and wiki in an EFL blended learning setting," *System*, vol. 38, no. 2, pp. 185-199, 2010.
- [39] J. Lazar, A. Jones, M. Hackley and B. Shneiderman, "Severity and impact of computer user frustration: A comparison of student and workplace users," *Interacting with Computers*, vol. 18, no. 2, pp. 187-207, 2005.
- [40] S. C. Chang and F. C. Tung, "An empirical investigation of students' behavioural intentions to use the online learning course websites," *British Journal of Educational Technology*, vol. 39, no. 1, pp. 71-83, 2008.
- [41] S. M. Almahamid, A. F. Tweiqat and M. S. Almanaseer, "University website quality characteristics and success: lecturers' perspective," *International Journal of Business Information Systems*, vol. 22, no. 1, pp. 41-61, 2016.
- [42] A. Sonderegger and J. Sauer, "The influence of design aesthetics in usability testing: Effects on user performance and perceived usability," *Applied ergonomics*, vol. 41, no. 3, pp. 403-410, 2010.
- [43] <https://www.kissmetrics.com/>
- [44] <http://www.internetworldstats.com/stats19.htm>
- [45] <http://www.webometrics.info/en>
- [46] <http://www.webometrics.info/en/node/200>
- [47] <http://www.usthb.dz/ar/>
- [48] <http://www.ump.ma/ar>
- [49] <http://uob.edu.ly/>
- [50] <http://yslow.org/>
- [51] <https://gtmetrix.com/recommendations.html>
- [52] <https://speedof.me/>
- [53] <http://www.websiteoptimization.com/services/analyze/>
- [54] <https://www.nngroup.com/articles/usability-101-introduction-to-usability/>

Automatic Sign Language Recognition: Performance Comparison of Word based Approach with Spelling based Approach

Shazia Saqib

Department of Computer Science, GC University,
Lahore, Pakistan,
Works at Lahore Garrison University, Pakistan

Syed Asad Raza Kazmi

Department of Computer Science, GC University,
Lahore, Pakistan

Dr. Khalid Masood

Lahore Garrison University,
Pakistan

Saleh Alrashed

University of Dammam,
Dammam, KS

Abstract--Evolution of computer based interaction has been through a number of phases. From command line interface to menu driven environment to Graphics User Interface, the communication has evolved to a better user friendly environment. A new form of communication is on the rise and that is Gesture Based Communication, which is a touch free environment basically. Although its applications are mainly for deaf community but smart mobiles, laptops and other similar devices are encouraging this new kind of communication. Sign languages all over the world have a dictionary of signs of several thousand words. Mostly these signs are word based which means that these signs do not make use of basic alphabet signs, rather a new sign has to be designed for every new word added to the dictionary. This paper suggests use of spelling-based gestures especially while communicating with smart phones and laptops.

Keywords—Feature extraction; human computer interaction; image segmentation; object recognition

I. INTRODUCTION

A “sign” means a meaningful unit element of communication in a sign language. A sign may be a simple sign or it may represent a word. Signs may represent hand gestures or other signs. Whole sign recognition system works around to identify signs conveyed in a gesture. Sign languages pose the challenge that there is unfortunately no international sign language. All sign languages suffer from lack of rules and regulations [1], [2].

A gesture conveys an information. When we press a key, although it creates a piece of information but it is not a gesture, so we need a different class of input devices to understand gestures. Human hand and fingers have larger degree of freedom as compared to mouse or other pointing devices. The human hand has roughly 27 degrees of freedom. While dealing with hand gestures we need to learn two concepts:

- Hand posture: A hand posture is a simple hand sign that has no hand movement.

- Hand gesture: A hand gesture is collection of one or more continuous hand postures. When we hold our fist in a particular way it is a hand posture. A hand gesture is a dynamic movement, such as moving hand to wish farewell [3].

Many fiction writers like in “Star Trek” had foreseen this future development and used these ideas in the movies. They used the idea of shouting in space and let the computer process the command [3]. Another way has been Gaze Detection, use of Bolt added gaze detection to improve the interface. There were up to thirty different moving images at a time on wall with all sound tracks as one track. Gaze detection helped to resize moving images.

The gesture-based input gives us a natural and interactive interface. When we use touch free interface we can be more focused on our problem rather than the input [4]. The following are the application areas of hand gestures system:

A. Evolution of Sign Language

Gestures have gained popularity all over the world for communication with smart devices. Another major area is helping deaf community to make smooth interaction with mobile and computers. American Sign Language (ASL) and British Sign Language (BSL) are based on English language [5]-[7] whereas, German Sign Language (GSL) [8], Argentinian Sign Language (ArSL) [9], Indian Sign Language (ISL) [10], Persian sign language[11], Arabian Sign Language [12], [13] and Chinese Sign Language (CSL) [14] are also among the well-known sign languages.

B. Robot Control

Robot control is another area where gestures are highly useful. One interesting application of robotics, i.e. all finger gestures are assigned special meaning e.g. 1 may mean come forward. These systems used one hand fingers only.

C. Gestures in Graphic Editing

In this hand gesture movement is tracked and depending on the movement of hand a particular action is taken care of.

Min et al. has involved 12 dynamic gestures to make tools available for creating graphic systems [15]. It uses building blocks like triangle, rectangular, circle, arc, lines for creating the environment for various graphics functions.

D. Virtual Environments (VEs)

Another popular application of gestures is virtual environments VEs for 3D pointing gesture recognition.

E. Numbers Identification

Gestures are used to identify numbers. Elmezain et al. has proposed algorithm for recognition of Arabic numbers in real time using HMM [16].

F. 3D Modeling

Sign languages are more powerful in building 3D models, sometimes even hand shadows are used to build 2D and 3D objects [17].

The following are the reasons, why users love gesture-based interface:

Reason 1: It is the safest method to communicate

Gestures are most ancient way of communication. Even for a foreign language communication, where we fail to use words, we speak our mind by using our hands, and you can be sure that you'll have a good chance of being understood. Moreover regardless of nationality, No matter where you come from, you have the same methods of indicating that something is big, that the meal you've just eaten was tasty or that something smells awful.

Reason 2: Trend of adopting gesture-based design in applications

In last few years, the trend of adopting gesture-based design in applications (mobile applications) has become more common thanks to growing usage of devices with touchscreen interfaces. In the not-so-distant past, everyone mainly used mobile devices such as smartphones and tablets. This new segment of devices requires a completely new approach to the problem, because the interface that we will be working with is different. It's worth pointing out that we can use these new devices on our lap or on the table, so both hands are free. This wipes out a whole load of restrictions so designers are able to discover and develop new solutions in the area of interfaces and interactions.

Reason 3: Touch Free interface is fun and invokes Powerful Response

The popularity of applications such as Snap Chat or Tinder has influenced users' awareness and feel of interfaces based on gestures. This has opened up the option of honing app design down to simplicity by getting rid of some visual elements of navigation. A key trend to notice in this field is the gradual disappearance of differences between popular platforms (iOS, Android) and increasing unification in terms of possible interactions based on gestures.

Simpler design and unified solutions allow the user to concentrate fully on the content and absorb it in a more intuitive way. This type of application is way more fun for the

user and elicits a more positive and powerful emotional response. A more intuitive approach to usage that allows us to get into the app's content more directly and simply makes our apps more attractive and, therefore, more highly rated.

Reason 4: Gesture-based approach is constantly evolving and is a game changer

The gesture-based approach is constantly evolving and, in the future when combined with voice recognition, it may well turn out to be a game changer.

Usually gestures are a part of a continuous sequence of signs. This is very active area of research in gesture recognition as it is very hard to locate start and end of a gesture. Different signers have different gesture boundaries thus making it quite hard to enumerate all gestures. A signer gives a gesture starting from a pause state and ending in a pause state even when gesturing continuously [3].

The following performance parameters were kept in mind during the design of the research:

Recognition time: A good recognition time varies from 0.25 to 0.50 seconds.

Continuous and automatic recognition: System accuracy should be high enough to recognize many dynamic gestures in one go.

Recognition Accuracy: A good accuracy rate varies in between 80-90 percent [18], [19].

Gesture based systems are affected by image transition, scaling and illumination affects. Response time and the cost of interpretation of gesture are other important factors that measure the effectiveness of the system [20]. Fig. 1 shows few such factors. Our goal is to help this disabled population and side by side providing a touch free environment for our smart devices [21].

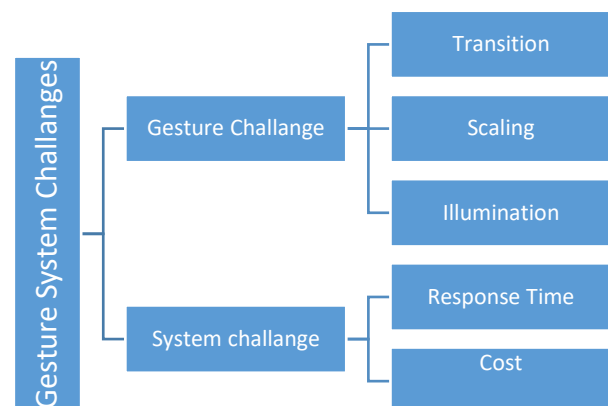


Fig. 1. Gesture system challenges.

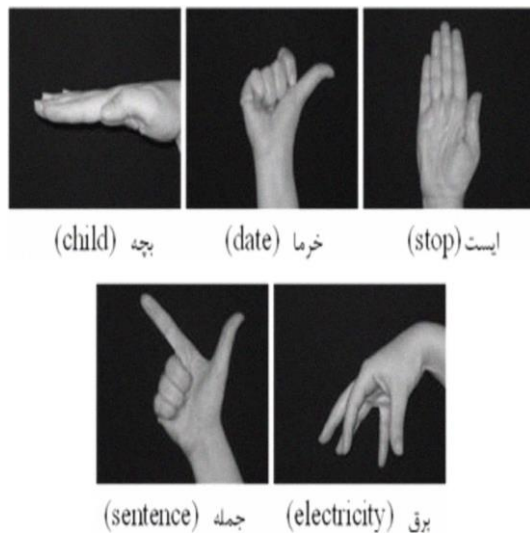


Fig. 2. Few words in Persian Sign language.

Different words have different gestures in different sign languages, e.g. the word “date” in Persian Sign Language has same symbol as that of Seen (س) in Pakistan Sign Language and the gesture for the word “sentence” in Persian Sign Language resembles the symbol laam (ل) in Urdu language. That means there are no universal signs [11], [22], [23]. Fig. 2 shows few word based signs in Persian sign language.

As compared to verbal languages, sign languages are regional languages. Even in Arab countries lots of efforts have been made to establish same standard sign language used in individual countries [7]. The significance of using hand gestures for communication becomes clearer when sign language is considered. Sign language uses gestures, postures, and facial expressions to identify the signer’s input [3], [24].

A sign language uses manual communication and body language to convey meaning rather than using sound to communicate [25]. Every gesture represents a letter or a word and it may hold different meaning in different languages for the word “What” the sign in different Sign languages is shown in Fig. 3.



Fig. 3. Sign of “What” in different sign languages.

We can see here some sign languages use two hand gestures while some use one hand gesture.

II. PROPOSED APPROACH

Most of the sign languages used are using word based gestures, which means suppose there are 5000 words in the sign language and we need to add one more word to our gesture database, we will need a new gesture which does not base on any algorithm or automata structure. So whenever a new word becomes part of sign language, the word gesture has to be learned.

The proposed **spelling based gestures** are based on the fact that alphabets make words. So static gestures can be used to form gestures for dynamic gestures. In this way whenever new words are added to dictionary of sign language, there will be no need to learn a new sign. New words will be actually sequence of static signs.

Moreover, if word gestures are spelling based, we can use any technique like regular expression, expression tree or finite automata for gesture recognition. The gesture based input are very powerful means for communication. Their advantages are: naturalness, more expressive, direct interaction and freedom of expression through signs. Fig. 4 shows word-based gestures for shoe while Fig. 5 shows same when we use spelling-based gestures.

Recognition and interpretation of a sign holds biggest challenge in gesture recognition. When a gesture starts or ends, this is another research question at the moment for image processing experts. However, when we deal with word-based gestures the intensity of problem is reduced considerably. In spelling based gestures a sequence of inputs is given however boundaries are blurred [4].

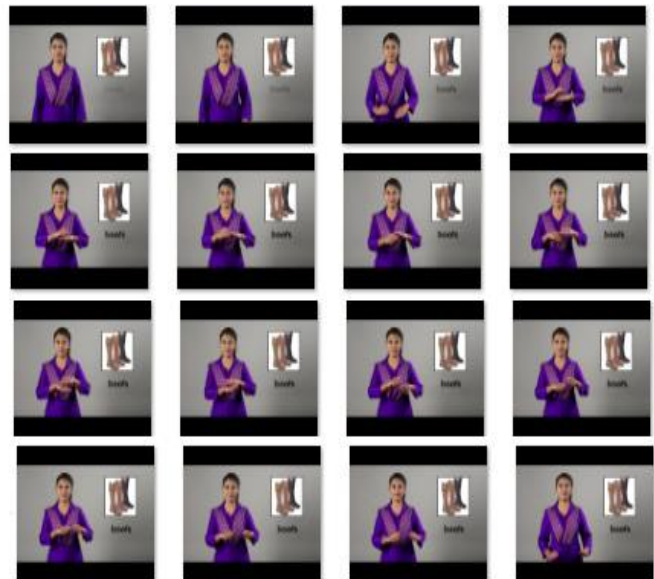


Fig. 4. Word based gesture for shoe.

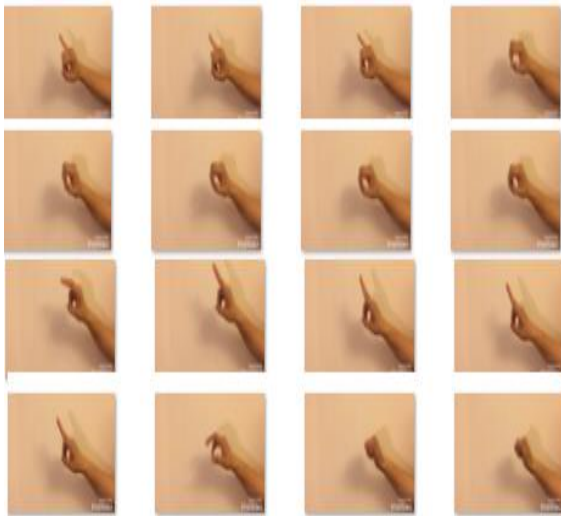


Fig. 5. Spelling based gesture for shoe.

III. EXPERIMENTAL SETUP

Spelling based gestures can be one hand or two hand, they are usually one hand. That is another added benefit, we are usually holding something in our one hand. So, one handed interaction will give us more freedom.

Videos for both have been converted to frames. It was found that the average time taken by word-based gesture shoe(joota) took 4 seconds while the spelling-based gesture for shoe(joota) took 3 seconds.

To compare different word based and spelling based gestures, we used 30 words with five different signers. The camera for capturing images can be any webcam, mobile camera or ordinary laptop camera. The images are captured at a resolution of 8 mega pixels. The images are taken from an average distance of 4 to 5 feet.



Fig. 6. Word-based gesture for banana.



Fig. 7. Spelling based gesture for banana.

Fig. 6 and 7 shows gestures for Banana in both word based and spelling-based gestures. Fig. 8 shows a chart of comparison of word based and spelling-based gestures of few words.

Using 5 signers, different videos for different words were made by ordinary mobile and laptop camera. The average distance between the signer and image capturing device was 5 feet approx. The experiment was repeated time and again for different signers. Table I gives measure of average time taken by each gesture.

TABLE. I. TIME COMPARISON OF SPELLING BASED AND WORD BASED GESTURES (TIME IN SECONDS)

	Spelling based gestures	Word based gestures
Apple	6	4.5
Banana	7	6.5
Raisins	7	6.5
Lychee	7	7
Mango	4	5.5
Shoe	3.5	5

IV. RESULTS AND DISCUSSION

Pakistan sign language has been taken as a case study for experimentation however any sign language can be chosen. By looking at the video comparison result of both types, we can draw the following conclusions:

- 1) Both types of gestures take almost the same time. Difference of time in both is almost negligible.
- 2) Spelling based gestures are very easy to learn and generate. A signer needs to learn only 26 alphabets and 10 numeric signs.
- 3) If we can devise an algorithm that can separate character boundaries, we can use a construct like finite automata or fuzzy logic to identify the sign.

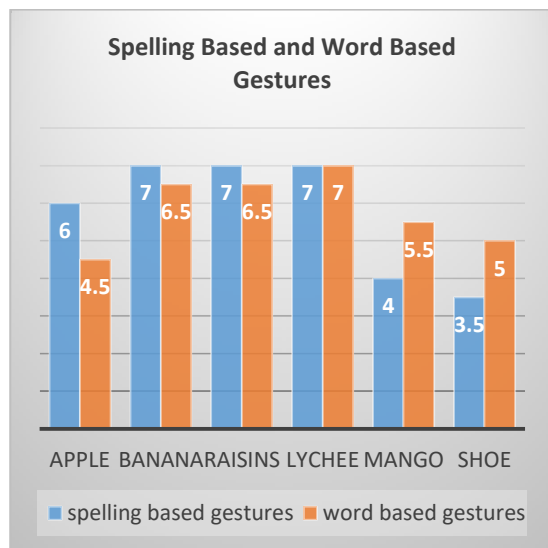


Fig. 8. Graph of word based gestures and spelling based gestures time measured in second.

4) If we can devise an algorithm that can identify first alphabet of the gesture, searching the gesture in data set becomes very easy.

V. CONCLUSION AND FUTURE WORK

Although word-based gestures are in use all over the world, adopting spelling-based approach can bring a revolution for deaf as well as for an effective touch free interface. Every region has its own sign language. Although the whole research in this area started with an idea to facilitate the deaf community but it actually formed the basis of touch free interface for smart devices. This leads to a very strong need for **universal sign language**. Moreover, estimating start of a gesture and end of a gesture can be a very attractive area of research. At the moment it is very hard to separate symbols however in next few years advancements in technology will be able to easily handle this issue. Present research can recognize static gestures with great accuracy. Very soon spelling based gestures will be adopted as this is most easy way to expand the sign language.

ACKNOWLEDGMENTS

The authors are indebted to The Deaf Reach Program Pakistan for their research on Pakistan Sign Language. They have provided huge vocabulary for 5000 plus words which has been very helpful for this research.

REFERENCES

[1] S. Tehsin et al., "Text Localization and Detection Method for Born-digital Images Text Localization and Detection Method for Born-digital Images," vol. 2063, no. February, 2016.
[2] M. Melnyk, V. Shadrova, and B. Karwatsky, "Towards Computer Assisted International Sign Language Recognition System: A Systematic Survey," Int. J. Comput. Appl., vol. 89, no. 17, pp. 44–51, 2014.

[3] Q. Chen, M. D. Cordea, E. M. Petriu, T. E. Whalen, I. J. Rudas, and A. Varkonyi-koczy, "Hand-Gesture and Facial-Expression Human-Computer Interfaces for Intelligent Space Applications," 2008.
[4] M. Van Beurden, W. Ijsselsteijn, and Y. De Kort, "User experience of gesture-based interfaces: A comparison with traditional interaction methods on pragmatic and hedonic qualities," pp. 121–124.
[5] R. Varga and Z. Prekopcsák, "Creating a Database for Objective Comparison of Gesture Recognition Systems," pp. 1–6, 2011.
[6] J. R. Pansare, S. H. Gawande, and M. Ingle, "Real-Time Static Hand Gesture Recognition for American Sign Language (ASL) in Complex Background," vol. 2012, no. August, pp. 364–367, 2012.
[7] P. Pandey and V. Jain, "Hand Gesture Recognition for Sign Language Recognition: A Review," vol. 4, no. 3, 2015.
[8] J. Forster, C. Oberdörfer, O. Koller, and H. Ney, "Modality Combination Techniques for Continuous Sign Language Recognition," Pattern Recognit. Image Anal., pp. 89–99, 2013.
[9] F. Ronchetti, F. Quiroga, and L. Lanzarini, "LSA64: An Argentinian Sign Language Dataset," pp. 794–803.
[10] J. Singha and K. Das, "Indian Sign Language Recognition Using Eigen Value Weighted Euclidean Distance Based Classification Technique," vol. 4, no. 2, pp. 188–195, 2013.
[11] A. K. Sarkalehl, F. Poorahangaryan, B. Zan, and A. Karami, "A neural network based system for persian sign language recognition," ICSIPA09 - 2009 IEEE Int. Conf. Signal Image Process. Appl. Conf. Proc., pp. 145–149, 2009.
[12] N. R. Albelwi, "Real-Time Arabic Sign Language (ArSL) Recognition," pp. 497–501, 2012.
[13] B. Jalilian and A. Chalechale, "Face and Hand Shape Segmentation Using Statistical Skin Detection for Sign Language Recognition," vol. 1, no. 3, pp. 196–201, 2013.
[14] C. Wang, X. Chen, and W. Gao, "Expanding Training Set for Chinese Sign Language Recognition," pp. 0–5, 2006.
[15] B. Min, H. Yoon, J. Soh, T. Ejimau, and I. Engineering, "Recognitron Using," pp. 4232–4235, 1997.
[16] M. Elmezain, A. Al-Hamadi, J. Appenrodt, and B. Michaelis, "A Hidden Markov Model-Based Isolated and Meaningful Hand Gesture Recognition," Proc. World Acad. Sci. Eng. Technol., vol. 43, no. July, p. pages 394–401, 2009.
[17] R. Zaman Khan and N. A. Ibraheem, "Comparative Study of Hand Gesture Recognition System," Comput. Sci. Inf. Technol. (CS IT), vol. 6, no. 4, pp. 203–213, 2012.
[18] A. K. Alvi, A. Muzaffar, and M. Usman, "Project Mentor: Team Members:,"
[19] A. K. Alvi et al., "Pakistan Sign Language Recognition Using Statistical Template Matching," vol. 1, no. 3, pp. 1–4, 2007.
[20] S. M. Darwish, M. M. Madbouly, and M. B. Khorsheed, "Hand Gesture Recognition for Sign Language: A New Higher Order Fuzzy HMM Approach," vol. 8, no. 3, 2016.
[21] M. S. Abdalla and E. E. Hemayed, "Dynamic Hand Gesture Recognition of Arabic," vol. 13, no. 5, 2013.
[22] M. Moghaddam, M. Nahvi, and R. H. Pak, "Static Persian Sign Language Recognition using Kernel-based Feature Extraction," 2011.
[23] S. N. Sawant, "Sign Language Recognition System to aid Deaf-dumb People Using PCA," vol. 5, no. 5, pp. 570–574, 2014.
[24] N. El-bendary, H. M. Zawbaa, M. S. Daoud, A. E. Hassanien, and K. Nakamatsu, "ArSLAT: Arabic Sign Language Alphabets Translator," vol. 3, pp. 498–506, 2011.
[25] N. S. Khan, A. Abid, K. Abid, U. Farooq, M. S. Farooq, and H. Jameel, "Speak Pakistan: Challenges in Developing Pakistan Sign," vol. 30, no. 2, pp. 367–379, 2015.

Geographical Distance and Communication Challenges in Global Software Development: A Review

Babur Hayat Malik, Saeed Faroom, Muhammad Nauman Ali, Nasir Shehzad, Sheraz Yousaf, Hammad Saleem
Department of CS&IT
University of Lahore, Gujrat Campus Pakistan

Abstract—Due to innumerable advantages the Global software engineering is trending now a days in software development industry. Basic drivers for this trend are flexibility, faster development and expected cost saving. Software development has moved from traditional development to the global software development (GSD). Global software development is very important and ordinary practice in the software industry. In GSD, the developers are distributed across different sites and different countries, and lots of problems arise due to the physical social and cultural barriers. Global Software development is facing a number of challenges including Geographical distance, Communication and collaboration, time, culture, trust, tasks distribution, requirements gathering and collaboration. In this paper, authors conducted a detailed study on geographical distances and communication challenges in GSD, their inter dependencies, and also the proposed solutions and guidelines to address these challenges that are very critical in the success of GSD projects. Also in this paper a detailed literature review is provided, combined results are summarized and on the basis of these studies, a comparative study is made. This research will be helpful for other researchers to draw new strategies to tackle these challenges.

Keywords—Global Software Development (GSD); distributed software development; geographical distance challenges; communication and collaboration

I. INTRODUCTION

Global software development (GSD) is a phenomenon that is receiving significant interest from all over the companies in the world. In GSD, stakeholders from different national and organizational cultures are involved in developing software. No doubt Global software development complicates the collaboration among the team members who are working on the same project but on the different sites. GSD can offer benefits such as improving time to market, improve quality, access to a larger and Better-skilled developer pool, reduced development costs, save time and shared knowledge [3]. Author's contribution in the paper is the discussion of all the challenges in GSD and dependencies between challenges. Also, will list the benefits of GSD but main focus is GSD challenges.

The number of organizations distributing their software development processes globally keeps increasing and this change is having a deep impact on the way products are

considered, designed, constructed, tested and supplied to customers GSD takes several forms. Distance (time and space) creates many challenges in communication, coordination, organization, project planning and follow up, and work allocation. Advances in communication technology and tools have carried GSD in focus [6].

In this paper, authors will discuss the Challenges of Geographical Communication in global software engineering. The purpose behind this study is to find the factors that badly affect the communication effectiveness and how they work. Section 2 contains the detailed definition of Global Software Development and *Identifying* the factors that introduce problem in global software development; Section 3 contains the detailed Literature Review and Section 4 contains the proposed methods and defining strategies to minimize GSD problems then Section 5 contains the motivation followed by acknowledgment and conclusion. Finally, the references are mentioned.

II. GLOBAL SOFTWARE DEVELOPMENT

After its first foundation in the conference which is sponsored by "NATO science committee at the end of 1960" Software Engineering industry is growing continuously [2]. Due to internal and external improvements in development method, its evolution is continuing [2]. Today one of the most and important change in software industry is Global software development which is also known as distributed software development [1]. In Globalization Technology is geographically distributed and this helps the organizations to change their operating and development models [1]. In Global Software development different teams work from different places on a same project [3]. This help companies to save cost by outsourcing developments work to low-cost countries [5] and also save time by using the strategies like follow the Sun [4]. In order to support collaborative work on projects, software engineers communicate directly and through meetings [5]. Communication, particularly informal communication plays an important role in the success of any GSD team [6]. Due to different cultures, communication and coordination among developers is a major challenge [3]. Lot of work is done to overcome on these issues like regular meetings with manger of the project can be reduced communication and coordination gap [3]. Ideal solutions of these problems still lack.

A. Identifying the Factors that Introduce Problem in Global Software Development

Most of the time the professionals working on Global Software Development projects mention that inadequate communication is the key problem in performing requirements engineering activities [7], [8]. This problem arises mainly due to the loss of communication richness due to lack of one to one interaction among the teams. Other factors which arises the challenge of communication in global software development are also the geographical distance among the teams as teams are distributed across different countries. The language proficiency is another factor that causes the problem of communication in Global software development projects.

- First of these problems are the Time difference among different countries [7].
- Second problem is the Time separation which is the additional to time difference as this includes the problems of breaks, holidays and timetable laps [8].
- Cultural diversity is another problem as development teams are distributed across the different countries so every country has it's their language, culture and religions [9], [10].
- Knowledge management is another problem in the global software development as huge amount of information is coming from many sources and need to share all the information with all the teams working on the same project [7].

III. LITERATURE REVIEW

In the section of Literature Review Authors consider the approaches used by many of researchers to discuss the geographical communication challenges in global software development [6]. They will discuss the work done by many authors and their research results on the geographical communication challenges in context of global software development.

A. Analyzing and Evaluating the Main Factors that Challenge Global Software Development

In this paper the author the author consider that Global software development is increasingly trending and adopted in development organizations due to its innumerable advantages like minimizing cost, quick delivery [10]. However, the culture diversity and the difference of time are challenge in the performance of teams especially in the activities of requirements engineering as it is very crucial to have all of the stakeholders on board.

According to author, the major problem in global software development is the inadequate communication, the time, language and the cultural difference. Also the communication got mired because of the non-availability of the knowledge management strategies.

B. Communication Effectiveness in Global Virtual Teams: A Case Study of Software Outsourcing Industry in China

As the “global virtual teams” (GVT) having staff members from across the countries working on the same project. The success of the projects is highly relay on the communication among the teams working on the project. The purpose behind this study is to find the factors that badly affect the communication effectiveness and how they work. From the literature review the author found the two aspects which are as follows “critical success factor” (CSFs) and the team characteristics [6], [11].

The author focus on importance of communication among the global working teams with specified references. Also to identify the factors that affects the communication among the teams including selection and use of ICT, GVT management, task characteristics and the demographic diversity (Fig. 1).

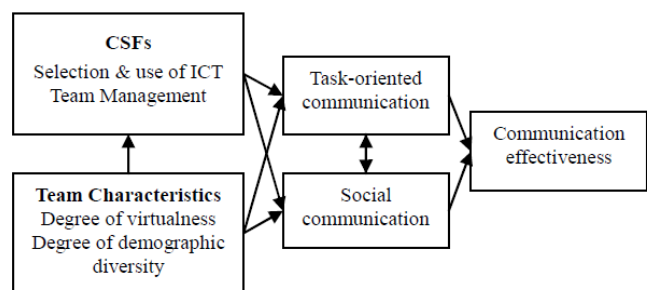


Fig. 1. GVT communication effectiveness model from the work of Min et al.

C. How do Distribution and Time Zones affect Software Development? A Case Study on Communication

According to the author the software projects now crossed the borders in search of talent and now consist of intra country teams that working on the same project. In this paper the author use a case study method to analyze the geographic communication difference in global software development. This case study is all about the three teams working on same project while they are student and the project continue for the two semester and the project teams are located at ten different countries [6]. This case study results that there is much difference in the communication size in the two location project and a three location project. The total amount of communication is much higher in the nearby locations or in two locations as compared to three location project. This case study also analyzes the effect caused by different time zones. On the basis of different time zones Authors can classify the project in to three time ranges which are:

- Large
- Medium
- Small

From this case study, authors found that in the small time zone range the amount of the communication is higher than the medium and the large [12]. Author also analyzes that in the small time range projects the reply to any e-mail comes faster than the projects with medium and large time range (Table I).

TABLE. I. STUDENT’S FEEDBACK FOR THE EFFECT OF TIME ONE AND CULTURAL DIFFERENCES IN DOSE 2010, VALUES RANGE FROM 1-5 , FROM THE WORK OF NORDIO ET AL.

Table # 1				
	Large	Medium	Small	Average
Times Zone affected quality	2.6	2.0	1.4	2.1
Times Zone affected productivity	3.1	2.5	1.4	2.5
Times Zone caused communication overhead	3.2	2.6	1.7	2.6
Cultural differences affected quality				
Cultural differences affected productivity	2.2	2.0	1.9	2.1
Cultural differences caused communication overhead	2.3	2.2	1.8	2.1
Local projects: the development would be easier	2.5	2.6	1.8	2.4
Local projects: the quality would be better	4.1	3.7	4.1	4.0
Local projects: the productivity would be higher	3.7	3.1	4.0	3.6
Local projects: the communication overhead would be lower	4.1	3.3	4.0	3.8
	4.0	3.7	4.2	3.9

D. Requirements Engineering During the Global Software Development: Some Impediments to the Requirements Engineering Process. A Case

The author presents in this paper that requirement Engineering is the most crucial task when teams are distributed across the countries or in case of global software development. There are two teams working on this projects that are situated in the UK and the other is software house working on the same project from New Zealand. The Phase of requirement is not easy for any software project. This paper present a case study on a project that contains distributed teams in two countries and the project was completed in the time of seven months.

The main drawback faced by the RE process during the Global software development team is communication. This issue may be further divided into the four categories [13]:

- Distribution of the clients and the development team
- Distribution of the development team
- Cultural Differences among the Clients and development team
- Cultural Differences among development team.

In Fig. 2, the author represent the intensity of communication and changes in requirement that occurs during the development of this project. As all the changes are embraced during the development of project as the Requirement engineering process is ongoing due to iterative in nature. The Curve line shows that the communication become more challenging as the requirements got changes during the development. At each stage the communication become more intense mainly due to [13]:

- Miscommunication/Misinterpretation
- Invalid Requirements

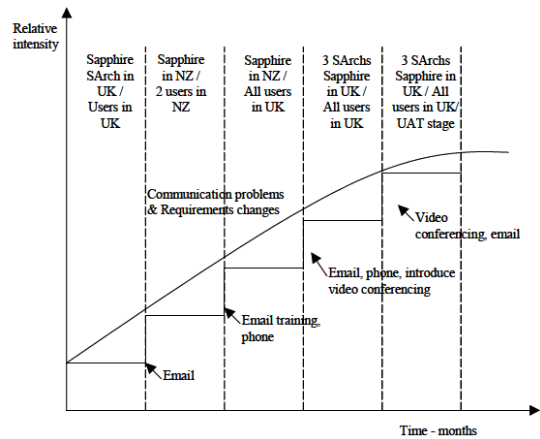


Fig. 2. Intensity of communication during the project from the work of Hanisch and Corbitt.

E. A Case Study of Customer Communication in Globally Distributed Software Product Development

In all of the cases communication of the customers was active and similar communication channels used to verify that different kinds of information are used. In development task and in developer position, coordination network is very important. Only difference in communication media is used of videoconferencing in case 2 but this is not available in case 3. Most interesting comparison which is discussed is between case 2 and case 3. IN case 3 they used agile approach and all members of 3 units were integrated through regular planning, meetings on daily basis despite time-zone difference. In all of these cases Indian and Irish development organization did not involve any user and customer. In US development organization customer group is slower in reacting due to transition from the traditional to agile approach. So rapid communication and regular agile meetings, involved customers, so they can be seen more successful communication between case 1 and in case 3 compared by US case 2 as shown in Fig. 3. Based on all case studies, Agile method in which customer involvement is consider best, due to involvement of the customers [14].

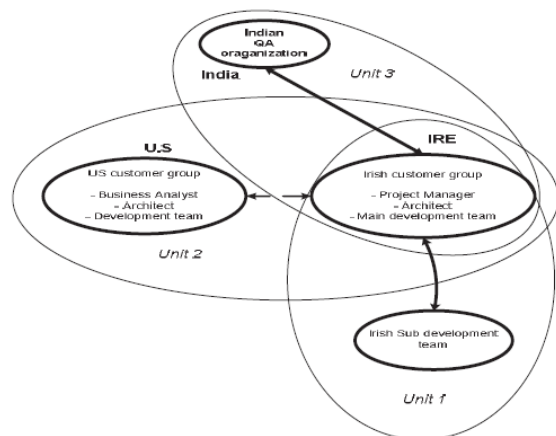


Fig. 3. The project organization and customer communication environment of the case project from the work of Korkala et al.

F. Communication Patterns in Geographically Distributed Software Development and Engineers' Contributions to the Development Effort

Unusual patterns of communication and coordination can be critical for the success of software project development. All roles of the customers in agile development is pre-defined. If all of the customers not present in the meetings then some of group of customers must be present in meetings. So base on different cases the result which they obtain is narrowed and limited view to the phenomenon observed. Even all of the information from different sources cannot be accessed, as many of the data utilized and analyzed by different researchers. So the Data triangulation which is introduced by stake used to check the validity of the results. This study shows that how communication patterns in the "Geographically distributed software development" (GDSD) evolve time to time [15].

G. Building Social Ties for Global Teamwork

The commitment of social ties and learning sharing to effective cooperation in distributed information system improvement groups has been investigated. Authors presume that in addition to technical solutions, human-related issues as social ties and learning sharing were revealed as keys to effective joint effort [16]. Specifically, the significance of compatibility and transitive memory was evident in the studied project. Besides, authoritative instruments that make and keep up social ties between scattered colleagues were reported for in detail. Authors recommend that future investigations ought to lead an overview over the information system industry in which the causal connections between these three primary ideas will be additionally examined as shown in Fig. 4.

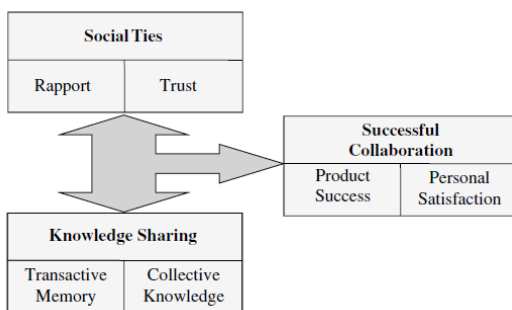


Fig. 4. Main concepts and their categories from the work of Oshri *et al.*

H. Exploring Collaboration Patterns among Global Software Development Teams

Figuring out how to work in worldwide software development student groups is challenging and in some cases even troublesome. Colleagues need to learn the most effective method to configuration, actualize, and approve software systems, as well as they should figure out how to function in socially assorted work groups, manage time, express thoughts, and speak with other peoples. Students should figure out how to utilize collective innovations for example, teleconferencing, video conferencing, email, voice mail, and groupware applications to speak with colleagues who might be found in different urban areas and even nations. The investigation

detailed in this paper looks at correspondence practices in worldwide programming improvement student groups. The creators of this paper describe the kinds of correspondence practices that happen when students groups are engaged with a software development project. Utilizing content and bunch investigations techniques, Authors recognized particular examples of cooperation and analyzed how these designs were related with task, culture, GPA, and performance of collaborative groups. Our outcomes propose that communication patterns among global software students might be identified with task, culture and GPA. It is hoped that these discoveries will prompt the advancement of new procedures for improving communications among global software teams [17].

I. Non-Optimized Temporal Structures as a Failure in Virtual

Global software development is trending day by day and making its worth in the market, but the management is worried about the failure of some projects and they intend to find the root causes for the failure of projects. This study consider the two virtually working teams and compare their effort and the time they take to complete their projects, to fine the success rate of both teams. It is found that the only reasons of the poor performance of one team are [18]:

- Entertainment of the temporal norms of the country
- Social situations of the members

This study defines the only reason behind the project failure is the ineffective communication and absence of meetings among the virtual teams working on the project [6]. "The core reason is that the teams remain in a limbo and cannot maintain momentum due to lack of discussions, feedback and supervision".

J. Culture in Global Software development - a Weakness or Strength

In these paper different cultures is discussed in global software development. As most of the complex issue in global software development is culture difference so Authors discuss how Authors can minimize these cultures issues in global software development. Global software development emphasis the need of knowledge of different culture issues to ensure the project success. Its Project manager and senior executive's responsibility to check existing culture difference and take steps to manage cultural diversity. Most of the strategies which Authors discussed already implemented in various globally distributed software development teams and companies. So mainly, Author discuss about Indian companies in this paper. Authors discussed different strategies which are suitable in Indian culture. So studying different cultures and strategies helps to manage global software development teams more efficiently [19].

K. The Impact of Intercultural Factors on Global Software Development

The main decision Authors can make from the previous foregoing is that there is more work to do. Some project managers have perceived the effect of intercultural factors on their global software development projects. A few analysts

have watched that intercultural factors influence the working relations of software engineers. The need remains for the improvement of tangible processes through which project managers can perceive the potential effect of intercultural factors on all phases of the product life cycle and, correspondingly, develop project and risk management strategies.

L. Critical Factors in Establishing and Maintaining Trust in Software Outsourcing Relationships

In software engineering we see that software outsourcing relationship is a comparatively innovative area of research. There is a growing awareness that understanding the dynamics of building and observance of expectation's between clients and sellers, who frequently need preceding affiliations and usually from different social backgrounds [21].

M. Bridging Gaps between Developers and Testers in Globally distributed Software Development

Authors believe that advances in addressing these issues can result in more efficient and actual methodologies for distributed software development and testing [22]. Results from the entire research/study expose more about retailer's needs that should be observed by full responsibility as well as experienced by clients in order to protect long term associations. Moreover, flexible behavior in terms of changing needs of client definitely comfort the advantage and preserve expectation time to time. Authors plan to conduct further empirical research by interviewing representatives of some clients of the companies' participating in our training. In software outsourcing relationships the conclusions will permit us to increase an understanding of client's expectations [21].

N. Global Software Development: Where are the Benefits?

This study tells us the benefits of GSD that are most important for an organization. GSD play an imperative role for the progress of any organization. But there are major valuable aspects of GSD. But our study is clearly defined that these are not clear. There may be the awareness of the risk that is related to GSD. But do not assume that the overall expenses will reduced as the wages are comparing with the higher management. Pure follow-the-sun software development the progress seems very unusual. Other companies like to make models instead of taking advantage of developers placed in various times. Rapid growth for progress there is seeking of employees. Share of information may be risky so do not share with their colleagues or do not trust on them. Taking advantage of closeness to foreign markets leads to a number of cultural problems which have to be addressed [23].

O. Improving Distributed Software Development in Small and Medium Enterprises

This paper is related to challenges that are related to DSD and how to overcome these challenges. And also define the strategies and methods that are used to overcome the challenges. In these methods and strategies which one is the best form all of them. Every industry has its own rules and regulation and it depends upon them how they distribute the work. Every industry has its own needs. These are the key factor to success. But the application of maturity models

(CMMI) which provide a good source through which to carry out variation near DSD [24].

The process should be automated through a tool which provides a proficient communication between members an organization. The use of a right PML and the use of environments such as Spreamint, Rational Method Composer or Eclipse Process Framework Composer for the model definition are essential to the generation of structured process guidelines which will facilitate training of human resources [24].

IV. PROPOSED METHODS AND DEFINING STRATEGIES TO MINIMIZE GSD PROBLEMS

A. From the work of Gabriela N. Aranda1, Aurora Vizcaíno and Mario Piattini

Discussed strategies minimize the problems about time zones in different countries, language understanding problem, types of team and culture difference by training of cultural difference in high and intermediate degree, to minimize the language problem in high and intermediate degree by acquaintance of communication initiator. By knowing the nature of people and culture regarding to their environment can minimize the communication problem. GSDs projects should deal with language difference as people have different mother language, so English language should be used for communication by stakeholders for better understanding the concept of their domain during the requirement gathering and all other phases. Ontological play a vital role of understanding for sharing vocabulary that is common to everyone because some words may have different meanings [25].

B. From the work of Qingfei Min, Zhenhua Liu and Shaobo Ji

Different countries have their own time zones which are different from others, so the time for teams which are at different places all over the world may overlap the time hours. So Verticalness of Global Team effects the management in selection of tool and teams and also effect the communication. Culture of a nation affects the Global Virtual Teams because the people sitting in other countries have their own culture ethics and in Global Software Development they should contact with other people who have different culture and communicate for sharing data and information about the projects which is going to be build. In GVT people interact with those who are very different from other and there will be gap in understanding the domain of the software [26].

1) Task Characteristics

Task characteristics effect the GVTs communication as shown in Fig. 5. Team can easily finish the simple task within short time but if the tasks are complex or a new project totally then it is very important to communicate the members of GVT frequently and it is only done by video conference in which all members can share their ideas to the whole team. If the tasks have further subtasks then it is essential to communicate each member with others members to share the subtasks because every member may have different subtasks of same task, and their communication effects their subtasks allocation. By poor communication subtask may be re-executed.

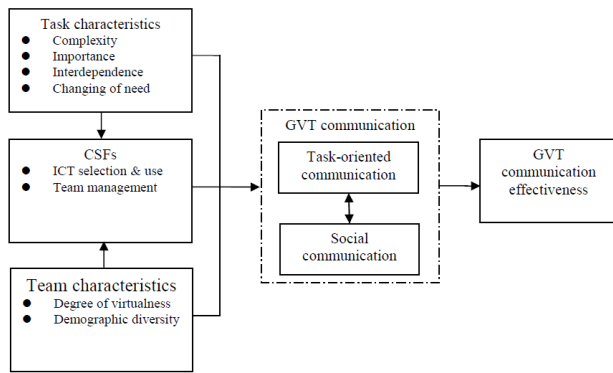


Fig. 5. Revised GVT communication effectiveness model. From the work of Martin Nordio et al.

During development some modules or requirements may be changed from the customer's side. In this situation the customer communicate with developer for specific changing, so the communication channel must be frequent to avoid the communication difficulties and to manage the cost and time of newly requirement. Most important and risky tasks may need more attention from the developer to develop because there is no chance of mistake because mistake may harm time, money or lives etc. In these types of tasks GVTs asked to pay their full attention on that task and for these the communication channel must be strong [27], [28].

Social communication in GVT members have great impact but at the starting of the work members only talk about the working for specific task for which they are connected but after spending some time with each other they become in relationship in social media and may got more chances of work from outside the organization and can help each other in some extend. The people who have relations with other can share work and ideas that become beneficial for both of them in on their initiative. But in the GVTs they have some private data and they can't share their data on social media, social communication is allowed but limited because then spend their working time in social communication. But social communication in spare time may encourage them. Following figure shows that how the task characteristics affect the factors that may influence GVTs communication and member's relationships and how they become beneficial for them [29].

C. From the work of Martin Nordio, H.-Christian Estler, Bertrand Meyer, Julian Tschannen

1) Analysis

As the time zones of all team members' locations are different to each other where teams are working on the specific phase of the development, Authors analyze the total time that teams expend on the projects and the time of their communication that due to geographical distance. Authors estimate the time of all phases and all members averagely regarding to team members size [30].

2) Communication in two-location and three-location Projects

We analyze the time of all members that they have spent on project adding the time that they have spent on their communication due to geographical distance by comparing

and finding their time ratio. Comparing by ratio can decrease the difference in their results points [31].

3) Reply Time of Projects in Different Time Zones

By comparing the time of email reply in Large, Medium and Small time zone ranges. Find that the time reply in large time zone range is maximum as compared to the Medium and Small time zone ranges [32].

D. From the work of "Mikko Korkala, Minna Pikkarainen and Kieran Conboy"

Utilizing the cooperation information from project A, Authors developed month to month communication and coordination systems [35]. Such an example proposes, to the point that a specific gathering of designers are at the focal point of the coordination activities and the trading of data among engineers. The rest of the engineers appear to depend exclusively on associations with the midway situated designers for organizing their assignments. A similar pattern was covered over each of the 39 months secured by the information. The solid center outskirts designs were logically affirmed utilizing Borgatti and Everett's [35] strategies for fitting system examples to a center fringe structure. The normal fit, in view of the consistent model, over each of the "39 months was 0.721 with a minimum fit of 0.568 and a greatest one of 0.858". Another vital finding delineated in Fig. 1 is that the center gathering, made out of specialists from every one of the three areas, appears to go about as doors or guards to other topographical areas for the developers in the periphery. Fig. 2 demonstrates the coordination organize from project B comparing to the primary organization of the overview. The general example of coordination conduct varies essentially from project A. There are couples of people that go about as "bridges" between geological areas. Indeed, those examples stem, Authors contend, from the meaning of formal parts to deal with cross site correspondence that were built up in project B. These two differentiating designs bring up fascinating issues, would one say one is example of coordination superior to the next? Assuming this is the case, which one and under which criteria? Past research has featured the basic part "contacts" people play in the execution of groups and improvement projects [33], [34]. The utilization of "contact" or "Gatekeepers" to deal with the conditions between groups has additionally been proposed as an instrument for encouraging coordination in geographically distributed software development [36]. In any case, a key issue in software development is the recognizable pieces of proof of the important specialized and assignment conditions. On the off chance that gatekeepers are deliberately implanted in the coordination systems, they could conceivably obtain the vital learning to find the imperative conditions and, therefore, give a profitable contact part. In any case, the distinguishing proof of the significant arrangement of conditions may require broad comprehension of the executed software code, learning that is commonly gained by being personally associated with the improvement exertion. The following segment analyzes in the connection between organize position and commitments to the improvement exertion.

In both projects, modification requests (MRs) spoke to a noteworthy segment of the development effort. Thus, the

quantity of MRs settled speaks to a decent measure of an engineer's commitment to the task. The longitudinal idea of the datasets renders customary direct relapse models insufficient for measurable investigation. Consequently, a multi-level model [37], additionally referred to in the writing as mixed regression models, was utilized to look at the impact of correspondence and coordination designs on singular level execution and its development after some time. The detail of a multi-level model incorporates settled and arbitrary impacts that might be connected to different factors for a given stream of longitudinal information. Along these lines, Authors represent the impacts of individual-level elements, qualities of the improvement work that are particular to an advancement aggregate and additionally occasional and other time-related fluctuation in our populaces.

E. From the work of Marcelo Cataldo and James D. Herbsleb

Face to face communication happened just between the two arranged Irish groups [38]. The sub improvement group had an on location client who additionally was bringing organized list of the requirements to the sprint arranging meetings. We have one business support that essentially finances all the work that Authors do. He is giving facilities to various different business and they may have different prerequisites and different needs, so we'd work in a group accord with reference to what Authors would do next. After each sprint, the group additionally displayed the results to the on location client to get criticism about the work that they have done "once They have something accessible, Authors likewise complete a demo for the clients [on-site customer] with the goal that They know that Authors are in good shape" (Developer). This finding shows that the client was engaged with the advancement process and was giving important input to the groups. It was moderately simple for the designers of the Irish sub improvement group to get to important data at whatever point required: "What was simple about the correspondence is I can just stroll down the hall and address some individual" (Developer). Obviously, there was no data covering up depicted in [38] display. In addition face to face communication, also wiki and email were effectively utilized.

Our perceptions support the contention that the agile practices are the best in the conditions where quick communication is empowered. The thought is that quick communication is probably going to cut down the measure of time spent on significant decisions the U.S. client group was associated with the basic leadership just in the start of the project when the objective was to characterize fixed up-front requirement for the overall product. After the first round of requirements definition, the client gather did not take an interest to general cycle arranging exercises or every day gatherings [20]. This led into the circumstance in which the requirements must be refined by the Irish primary group in isolation in light of what they trusted that the client required.

In any case, the regular agile meetings (Sprint arranging, discharge, and every day gatherings) were held inside "So it's a telephone call meeting, one individual in India bringing in also". In those continuous meetings, the reason for existing was to choose the objective and substance of the iteration. Although, the client groups from Units 1 and 3 were associated with the greater part of the coordinated gatherings

helping the advancement groups to settle on quick choices about the objectives and client stories that would be developed during the next iteration. Every one of these gatherings among Ireland and India were held through phone. In addition, wiki and email were effectively utilized for trading data inside the Unit 3. Normally, the act of having an on location client can be viewed as a key component in dynamic client cooperation. In any case, the Irish planner filled in as a client likewise for the Indian QA. Authors didn't discovered significant difficulties, for example, data hiding and absence of customer involvement, in the correspondence between the Indian QA association and the Irish fundamental group. For this situation agile meetings expanded the straightforwardness of the work and empowered data sharing amongst India and Ireland. Videoconferencing was only occasionally used [38]. Videoconferencing was just once in a while utilized. The dotted line between U.S. furthermore, IRE shows that this correspondence relationship has just been dissected in [38], [39].

F. From the work of Julia Kotlarsky and Ilan Oshri

We will present SAP and LeCroy case study results in this section. A study which is based on empirical evidence which shows that social ties and knowledge sharing contributed to successful collaboration in the companies [16]. On the base of the data which Authors analyzed, Authors claim that in globally distributed software development, teams, knowledge sharing and social ties improved collaboration [16]. To prove this argument three level of evidence will be discussed in this section. On first level all statements which made by the interviewees associated with the concepts which Authors are investigated. On second level, frequency of these statements are checked and on third level all number of instances present in which all social ties, collaborative tools and knowledge sharing were linked to successful collaboration. LeCroy and SAP evidence analysis suggested that there were two phase who supported the build-up of social ties: 1) before face to face; 2) after face to face. Empirical evidence analysis suggested that there were some tools which applied on the projects [16]. SAP interviews consider prior to face to face meeting for building social ties. LeCroy managers also consider initial activity before the face to face meeting for built the social relationship. Non-hierarchical communication is also important for social relationship. So far all of the evidence which is important for social aspects in globally distributed teams has been presented.

G. From the work of Fatma Cemile Serce, Ferda-Nur Alpaslan, Kathleen Swigger and Robert Brazile

1) Overview of communication behaviors in groups

Over all the twenty worldwide programming development learning ventures groups, an aggregate of 1985 correspondence episodes were investigated. In the event that the conduct was definitely not display in a correspondence episode, it was doled out a score of 0; on the other hand, if a correspondence behavior(s) was available in a posting, at that point it was appointed the code or then again codes for that conduct. As an unwavering quality check, a second coder examined similar talks. Between rater unwavering quality between coders for the association's practices was adequate [38], [39].

H. From the work of Mark Grechanik, James A. Jones, Alessandro Orso and Andr'e van der Hoek

Conventional programming cost models depend on the supposition that everybody engaged with a product venture is headed to make it fruitful also concurs on the objectives and strategies to make progress. Nonetheless, distinctive group members see a definitive achievement of the venture diversely in light of their own objectives. This is particularly valid in settings that include performing artists from various associations, as it is regularly the case in conveyed improvement. Authors trust that new refined financial models are required to examine programming ventures as no cooperative amusements to reveal concealed reasons for disappointments of programming undertakings and propose approaches to settle them [40]. A watchful examination of these financial variables of new programming advancement models will be basic for the achievement of profoundly disseminated advancement rehearses.

I. From the work of Miguel Jiménez, Aurora Vizcaíno and Mario Piattini

We recommend an approach to DSD in SME environments, by taking the limited complexity and budget of these organizations which usually results to applying basic methodologies, giving precise responsiveness to their organizational configuration. Not all of the activities proposed by the common standards "(ISO/IEC 12207 [41])" are always suitable for these environments, which also apply lower levels of maturity in association to larger companies.

1) Communication

This theory is established on the idea of taking out communication through structured models that will display the candidates in the organization of information to increase communication by decreasing the number of essential communications. This method should be used in all formal communication between concentrated members, improving the overall knowledge of the status of the project and keeping the information produced in a mutual source, thus helping avoid identical discussions. Developers may also need to communicate to other remote developers who are working on different parts of the software. It is not always possible to know which person is to contact so it is beneficial to take out communications through the local sub-director who need to accomplish the overall communications for that site and for that project. For locating members the distribution of organizational charts [42] which identify the location of members must also be taken into considerations and the use of ideas [43] is also recommended. Moreover, it is also essential to temporary informal communication, which will concluded the use of direct messaging e-mail's and programs. The tools used are Asynchronous communication tools based on recommendations and traditional E-mails Synchronous traditional tools (video-conferences and chats).

V. MOTIVATION

The major concern of this research is to do a review of different existing Literature to identify the main factors that introduces the problem in GSD and then focus on the specifically the factors effecting the geographically distance

and communication challenges in GSD and we will compare the guidelines and solutions to solve the issues causing problems.

VI. CONCLUSION

Global software development (GSD) is a phenomenon that is receiving significant interest from all over the companies in the world. In GSD, stakeholders from different national and organizational cultures are involved in developing software. However, GSD is technically and organizationally complex and presents a variety of challenges to be managed by the software development team. The number of organizations distributing their software development processes globally keeps increasing and this change is having a deep impact on the way products are considered, designed, constructed, tested and supplied to customers GSD takes several forms. Geographical Distance creates many challenges in communication, coordination, organization, project planning and follow up, and work allocation. Communication technology and tools have carried GSD in focus. I section three and four all systems are discussed in detail with respect to the geographical distance and communication challenges. In this paper we will do the detailed study on geographical distances and communication challenges in GSD their inter dependencies and also the proposed solutions and guidelines to address these challenges that are very critical in success of GSD projects.

REFERENCES

- [1] Timothy Haig-Smith and Maureen Tanner, "Cloud Computing as an Enabler of Agile Global Software Development," in Issues in Informing Science and Information Technology, vol. 13, pp. 121-144,2016.
- [2] Adrián Hernández-López, Ricardo Colomo-Palacios, Ángel García-Crespo, Pedro Soto-Acosta," Trust Building Process for Global Software Development Teams. A review from the Literature", 66 International Journal of Knowledge Society Research, 1(1), 66-83, January-March 2010.
- [3] Sami ul Haq, Mushtaq Raza, Asraf Zia, M. Naeem Ahmed Khan, " Issues in Global Software Development: A Critical Review", in J. Software Engineering & Applications, 4, pp. 590-595, 2011.
- [4] Mansooreh Zahedi a , Mojtaba Shahin b , Muhammad Ali Babar ,"A Systematic Review of Knowledge Sharing Challenges and Practices in Global Software Development".
- [5] Calefato, F., Damian, D., Lanubile, "An Empirical Investigation on Text-Based Communication in Distributed Requirements Workshops" In: Proc. of the Int. Conf. on Global Software Engineering, pp. 3–11,2007.
- [6] Cataldo, "Dependencies in Geographically Distributed Software Development: Overcoming the Limitations of Modularity. in PhD Dissertation, School of Computer Science, Carnegie Mellon University, 2007.
- [7] Julia Kotlarsky and Ilan Oshri,"Social ties, knowledge sharing and successful collaboration in globally distributed system development projects",in European Journal of Information Systems,pp 37–48,Vol 14, 2005.
- [8] Layman, L., Williams, L., Damian, D. and Bures, H. "Essential communication practices for Extreme Programming in a global software development team." Information and Software Technology, vol. Volume 48, pp. 781-794 (2006).
- [9] A. Al-Rawas, and S. Easterbrook, "Communication problems in requirements engineering: a field study," In: First Westminster Conference on Professional Awareness in Software Engineering, London, pp. 47-60, 1996.
- [10] Gabriela N. Aranda1, Aurora Vizcaíno2 and Mario Piattini, "Analyzing and Evaluating the Main Factors that Challenge Global Software

- Development”, in The Open Software Engineering Journal, Vol 4, pp. 14-25, 2010.
- [11] Qingfei Min , Zhenhua Liu and Shaobo Ji “Communication Effectiveness in Global Virtual Teams: A Case Study of Software Outsourcing Industry in China”, in Proceedings of the 43rd Hawaii International Conference on System Sciences , 2010.
- [12] Martin Nordio, H.-Christian Estler, Bertrand Meyer, Julian Tschannen, Carlo Ghezzi, Elisabetta Di Nitto, “How do Distribution and Time Zones affect Software Development? A Case Study on Communication”, in Sixth IEEE International Conference on Global Software Engineering, 2011.
- [13] Jo Hanisch, Brian J. Corbitt, “Requirements Engineering During Global Software Development: Some Impediments to the Requirements Engineering Process - A Case Study”, in European Conference on Information Systems, 2004.
- [14] Mikko Korkala, Minna Pikkarainen, Kieran Conboy, “A case study of customer communication in globally distributed software product development”, ACM International Conference Proceeding Series, 2010.
- [15] Marcelo Cataldo, James D. Herbsleb, “Communication Patterns in Geographically Distributed Software Development and Engineers’ Contributions to the Development Effort”, in Proceedings of the 2008 international workshop on Cooperative and human aspects of software engineering, Pages 25-28, 2008.
- [16] Julia Kotlarski and Ilan Oshri, “Social ties, knowledge sharing and successful collaboration in globally distributed system development projects”, in European Journal of Information Systems Vol 14, PP. 37-48, 2005.
- [17] Fatma Cemile Serce, Ferda-Nur Alpaslan, Kathleen Swigger, Robert Brazile, George Dafoulas, Victor Lopez, Randy Schumacker, “Exploring Collaboration Patterns among Global Software Development Teams”, In Fourth IEEE International Conference on Global Software Engineering, 2009.
- [18] Felix Köbler, Marilyn Tremaine, Jan Marco Leimeister, Helmut Krcmar, “Non-Optimized Temporal Structures as a Failure in Virtual”, In wirtschafsinformatik proceedings, 2009.
- [19] Sadhana Deshpande , Ita Richardson , Valentine Casey , Sarah Beecham , “Culture in Global Software development - a Weakness or Strength?”, in Global Software Engineering (ICGSE), 2010 5th IEEE International Conference, 2010
- [20] Eve MacGregor, Yvonne Hsieh, Philippe Kruchten, “THE IMPACT OF INTERCULTURAL FACTORS ON GLOBAL SOFTWARE DEVELOPMENT”, in Electrical and Computer Engineering, 2005. Canadian Conference on, 2005.
- [21] Phong Thanh Nguyen, Muhammad Ali Babar, June M. Verner, “Critical factors in establishing and maintaining trust in software outsourcing relationships”, in ICSE '06 Proceedings of the 28th international conference on Software Engineering, PP 624-627, 2006.
- [22] Mark Grechanik, James A. Jones, Alessandro Orso, Andr e van der Hoek, “Bridging gaps between developers and testers in globally-distributed software development”, in Proceedings of the FSE/SDP workshop on future of software research, PP. 149-154, 2010.
- [23] Eoin   Conch ur, P r J.  gerfalk, Helena H. Olsson, and Brian Fitzgerald, “Global Software Development: Where are the Benefits?”, in Communications of the ACM - A Blind Person’s Interaction with Technology, Vol 52, issue 8, PP. 127-131, 2009.
- [24] Miguel Jim nez, Aurora Vizcaino and Mario Piattini, “Improving Distributed Software Development in Small and Medium Enterprises”, in The Open Software Engineering Journal, Vol 4, PP. 26-37, 2010.
- [25] Prikladnicki, R., J. L. N. Audy and R. Evaristo (2003). "Global Software Development in Practice Lessons Learned." Software Process Improvement and Practice 8(4): 267 - 279.
- [26] P. Banerjee, “Narration, Discourse and Dialogue: Issues in the Management of Intercultural Innovation,” AI & Society, Vol. 17, pp. 207-224, 2003.
- [27] G. Walsham, “Globalization and ICTs: Working across cultures”, University of Cambridge, Cambridge, UK, 2001.
- [28] S. Krishna., S. Sahay, and G. Walsham, “Managing cross-cultural issues in global software outsourcing”, Communications of the ACM, 47 (4). 62-66, 2004.
- [29] J. L. Gibbs, “Loose coupling in global teams: tracing the contours of cultural complexity,” Ph. D. dissertation, University of Southern California, Los Angeles, CA, USA, 2002.
- [30] P Anderson, J.C. and Narus, J.A. “A Model of Distributor Firm and Manufacturer Firm Working Partnerships” ,in Journal of Marketing, 54 (1). pp. 42-58.
- [31] Loh, L. and Venkatraman, N. , “Diffusion of Information Technology Outsourcing: Influence Sources and The Kodak Effect”, in Information Systems Research, 3 (4). pp. 334-358.
- [32] Miles, R.E. and Snow, C.C. , “Causes of Failure in Network Organizations”, in California Management Review, 34 (4). pp. 53-72.
- [33] W. Aspray, F. Mayades, and M. Vardi, “Globalization and Offshoring of Software”, in ACM, 2006.
- [34] “Test Driven Development: By Example. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2002.
- [35] X. Cai and M. R. Lyu, “The effect of code coverage on fault detection under different testing”, In Profiles, ICSE 2005 Workshop on Advances in Model-Based Software Testing (A-MOST, pages 1-7, 2005.
- [36] C. Kaner, J. Bach, and B. “Pettichord. Lessons Learned in Software Testing”, in John Wiley & Sons, Inc., New York, NY, USA, 2001.
- [37] Y. W. Kim, “Efficient use of code coverage in large-scale software development”, In CASCON '03: Proceedings of the 2003 conference of the Centre for Advanced Studies on Collaborative research, pages 145-155. IBM Press, 2003.
- [38] Herbsleb, J. D., Mockus, A., Finholt, T. A. and Grinter, R. E., “Distance, dependencies, and delay in a global collaboration”, in ACM 2000 Conference on Computer Supported Cooperative Work, ACM Press, NY, 319-328.
- [39] Molokken, K., and Jorgensen, M.” A review of surveys on software effort estimation.”, In Proceedings of the 2003 International Symposium on Empirical Software Engineering.
- [40] K. Cho and D. Jonassen, “The Effects of Argumentation Scaffolds on Argumentation and Problem Solving,” Educational Technology: Research & Development, 50, 3, 2002, pp. 5-22.
- [41] W. Lloyd, M. B. Rosson, and J. Arthur, "Effectiveness of elicitation techniques in distributed requirements engineering," In: 10th Anniversary IEEE Joint International Conference on Requirements Engineering, RE'02, Essen, Germany, pp. 311-318, 2002.
- [42] K. Narayanaswamy , and N. M. Goldman, "A flexible framework for cooperative distributed software development," J. Syst. Softw., vol. 16, no. 2, pp. 97-105, 1991.
- [43] M. Paasivaara, and C. Lassenius, "Collaboration practices in global inter-organizational software development projects," Softw. Process Improv. Pract., vol. 8, no. 4, pp. 183-199, 2003.

Gaze Direction based Mobile Application for Quadriplegia Wheelchair Control System

Dr. Muayad Sadik Croock¹, Dr. Salih Al-Qaraawi², Rawan Ali Taban³

Computer Engineering Department
University of Technology
Baghdad, Iraq

Abstract—People with quadriplegia recruit the interest of researchers in introducing automated movement systems for adopted special purpose wheelchairs. These systems were introduced for easing the movement of such type of disabled people independently. This paper proposed a comprehensive control system that can control the movement of Quadriplegia wheelchairs using gaze direction and blink detection. The presented system includes two main parts. The first part consists of a smartphone that applies the propose gaze direction detection based mobile application. It then sends the direction commands to the second part via Wi-Fi connection. The second part is a prototype representing the quadriplegia wheelchair that contains robotic car (two-wheel driving car), Raspberry Pi III and ultrasound sensors. The gaze direction commands, sent from the smartphone, are received by the Raspberry Pi for processing and producing the control signals. The ultrasound sensors are fixed at the front and back of the car for performing the emergency stop when obstacles are detected. The proposed system is based on gaze tracking and direction detection without the requirement of calibration with additional sensors and instruments. The obtained results show the superior performance of the proposed system that proves the claim of authors. The accuracy ratio is ranged between 66% and 82% depending on the environment (indoor and outdoor) and surrounding lighting as well as the smart phone type.

Keywords—Gaze direction detection; mobile application; obstacle detection; quadriplegia; Raspberry Pi microcomputer

I. INTRODUCTION

Paralysis is most often caused by damage in the nervous system, especially the spinal cord injuries as results of a traffic accident or from some diseases. These diseases can include cerebral palsy or multiple sclerosis (MS) which lead to quadriplegia problems. Quadriplegia typically occurs as a result of an injury at the thoracic spinal levels (T1 or above) causing a loss of sensation and movement in all four limbs [1]-[2].

Paralyzed people are in need of an assistant to move them in their wheelchair, therefore they often feel powerless and burden on others [3]. Consequently, it became important to designing and developing a semi-automatic wheelchair, wherein the movement of the wheelchair can be controlled by movements of the different organs of the human body [1].

Different interfaces have been developed for dealing with people suffering from different disabilities; such as joystick control, hand control, head control and voice control. Other

systems employ infrared head-pointing devices, mouth-actuated joysticks, tongue-movement analysis, or head-mounted cameras as interfaces [4].

Communication abilities for severe quadriplegic people are practically limited to "Yes" and "No" responses using eye movements or blinks because the only muscles around the eye can be fully and easily controlled by them [5]. Thus, not all people are able to move around using the above-mentioned interfaces.

Currently, eye tracking technology has been used to design wheelchair control systems and other useful communication systems to help disabled people. Different techniques for eye tracking are developed such as Electrooculography (EOG), search coil method or Infrared oculography using head-mounted cameras. Most of the methods are considered to be uncomfortable as they have a drawback of using attached devices to the user's face [5], [6].

The aim of this paper is to design and implement a system to help quadriplegic people to be self-reliance in driving their wheelchair without being a burden on others. Therefore, an affordable eye tracking based control system for driving a wheelchair prototype model is presented. The system adopts the front camera of the mobile phone as a sensor to capture a real-time video for further processing using a designed mobile application that detects and tracks the eye movements. The detected action is sent as a command to a robotic car with two motors as a prototype instead of the quadriplegic wheelchair due to lack of availability. This is to present the move maneuvers, such as forward, backward, right, and left or fully-stop it, as shown in Fig. 1.

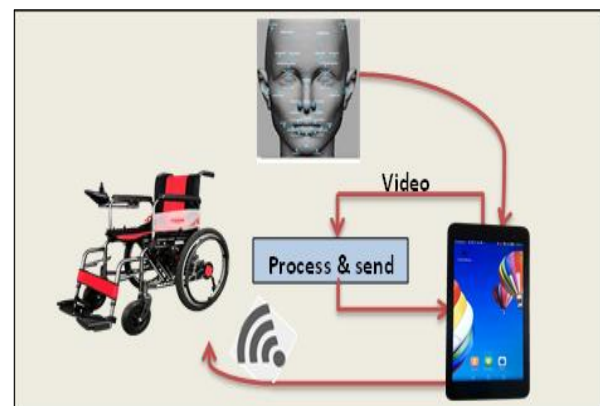


Fig. 1. Basic block diagram.

This paper is organized as follows. Section two presents the related works with wheelchair controlling systems. Different techniques for controlling the wheelchair are described in this section too. The proposed prototype control system based on eye-movement tracking and blink detection is shown in section three. Tests and experimental results are introduced in Section four. Section five and six give the limitation and a brief conclusion about the overall system respectively.

II. RELATED WORKS

There are many types of research concerned with designing control systems to help physically disabled people in driving their electric wheelchair.

In 2014, Biswajeet Champaty and et al. [7] proposed a prototype of an electric wheelchair control system based on Electrooculography (EOG) signals to move it in different directions. The system includes the blink detection in addition to gaze direction detection. The system was considered to be uncomfortable because of the face-attached EOG electrodes.

In 2017, Rana F. and Hani S. [8] presented a head direction based control system for disabled people to help them drive their wheelchair. The system used Viola Jones algorithm to detect four head directions. MATLAB is used for processing and sending the signal to the Arduino microprocessor for control the movement of a designed prototype. This system cannot be used by some people with severe paralysis.

In 2017, Xiang Gao and Lei Shi [9] designed a control system based on hand gesture detection for a robotic wheelchair to help the aged and the disabled people. The hand motion was detected using the Kinect which is a line of motion sensing input devices. The processing operations are performed on a PC placed in front of the user and then sending them to a DSP 28335 microcontroller to control the wheelchair. This system can be used only by Paraplegic and aged people because it needs a fully controlled hand.

In 2018, Raju Veerati, et al. [1] used a webcam with IR illuminators to track eye gaze direction with the help of MATLAB environment for processing. Detected gaze directions are used to control a wheelchair prototype which occupied by an ultrasonic sensor for obstacles avoidance to help people with disabilities. The results proved that the processing operations are quite slow; however the system runs with good accuracy. There are many types of research

concerned with designing control systems to help physically disabled people in driving their electric wheelchair.

III. PROPOSED SYSTEM

In general, the proposed system's structure consists of two main parts: the eye-sight placed mobile phone and the mini electric car which is used as a prototype of the quadriplegic wheelchair. The user's eye is tracked to estimate different actions and use it as input technique to the mobile. Hence, this application replaces the need of button pressing by finger touching. Based on the estimated eye actions, a related command is transferred via Wi-Fi connection to a microcomputer attached to the car for real-time motors controlling.

Different tools have been used to accomplish the work on the introduced system including:

The hardware tools:

- An Android smartphone.
- Two-wheel drive (2WD) electric car.
- Raspberry Pi microcomputer.
- Two ultrasonic sensors.
- Step down voltage regulator.
- Battery.
- Breadboard and connectors.

The software tools:

- Android Studio (2.3.3) IDE.
- OpenCV (2.4.9) library.

This section can be divided into three sub-sections for easing the reading flow of the presented work.

A. Proposed Algorithm

Most of the processing operations are performed at the software side which is the mobile application. In addition, the hardware side (the microcomputer) is responsible for generating the controlling signals used for car driving. Fig. 2 shows the proposed algorithm as a flowchart. All the detection operations are run in the background of the application which are the face, eye and pupil detection.

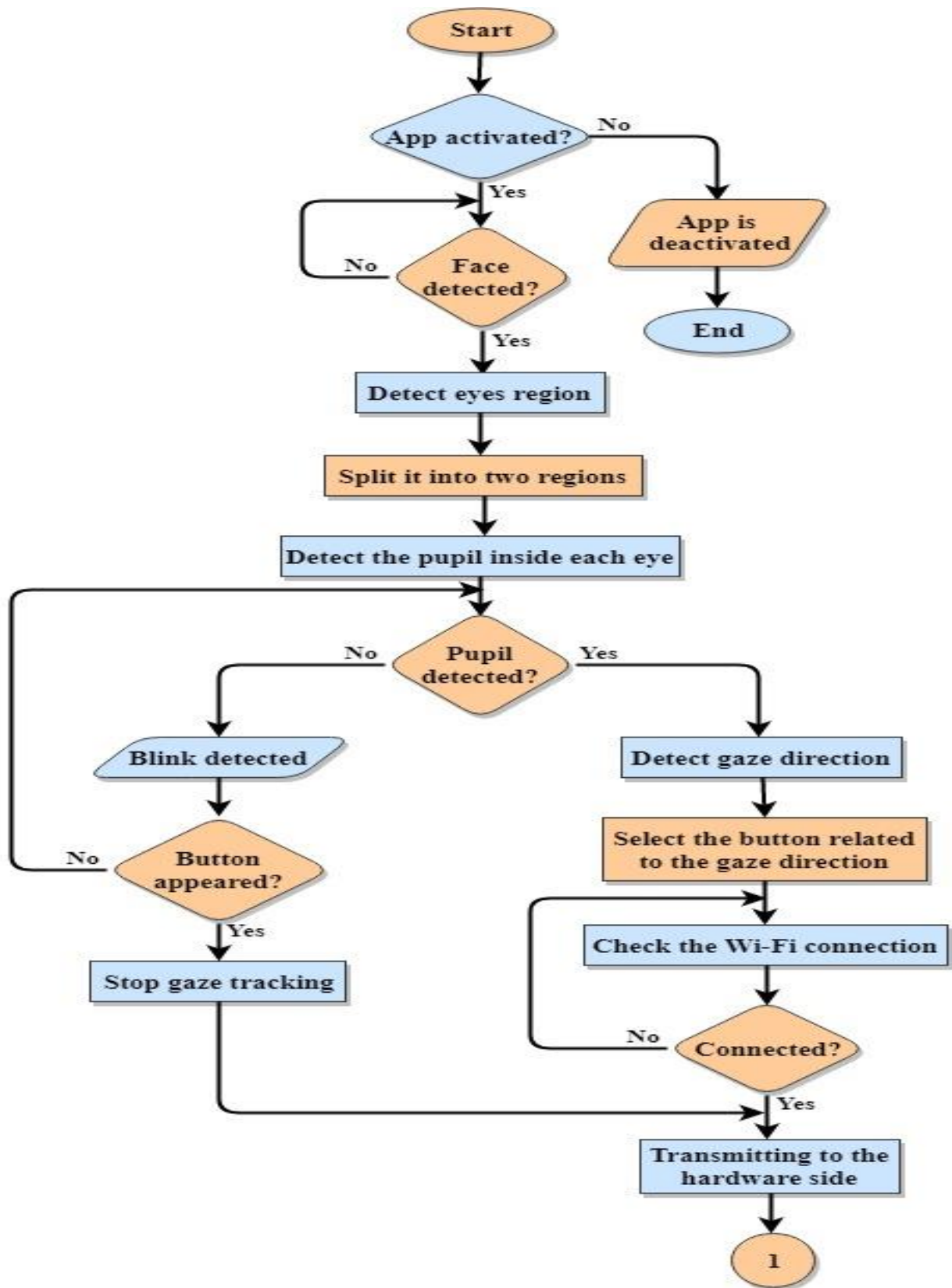


Fig. 2. (a) The software side.

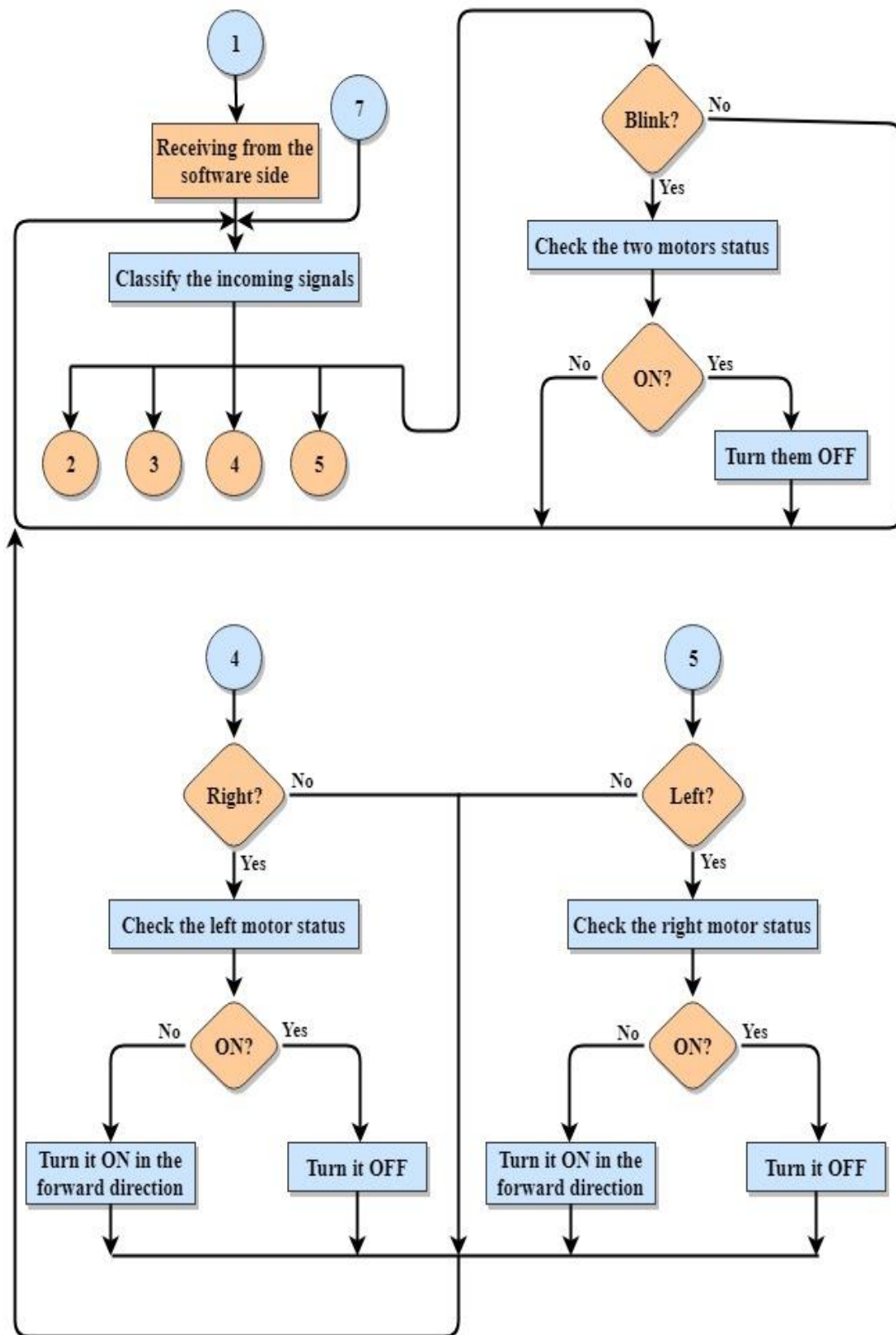


Fig. 2. (b) The hardware side: the blink, right and left gaze directions responses.

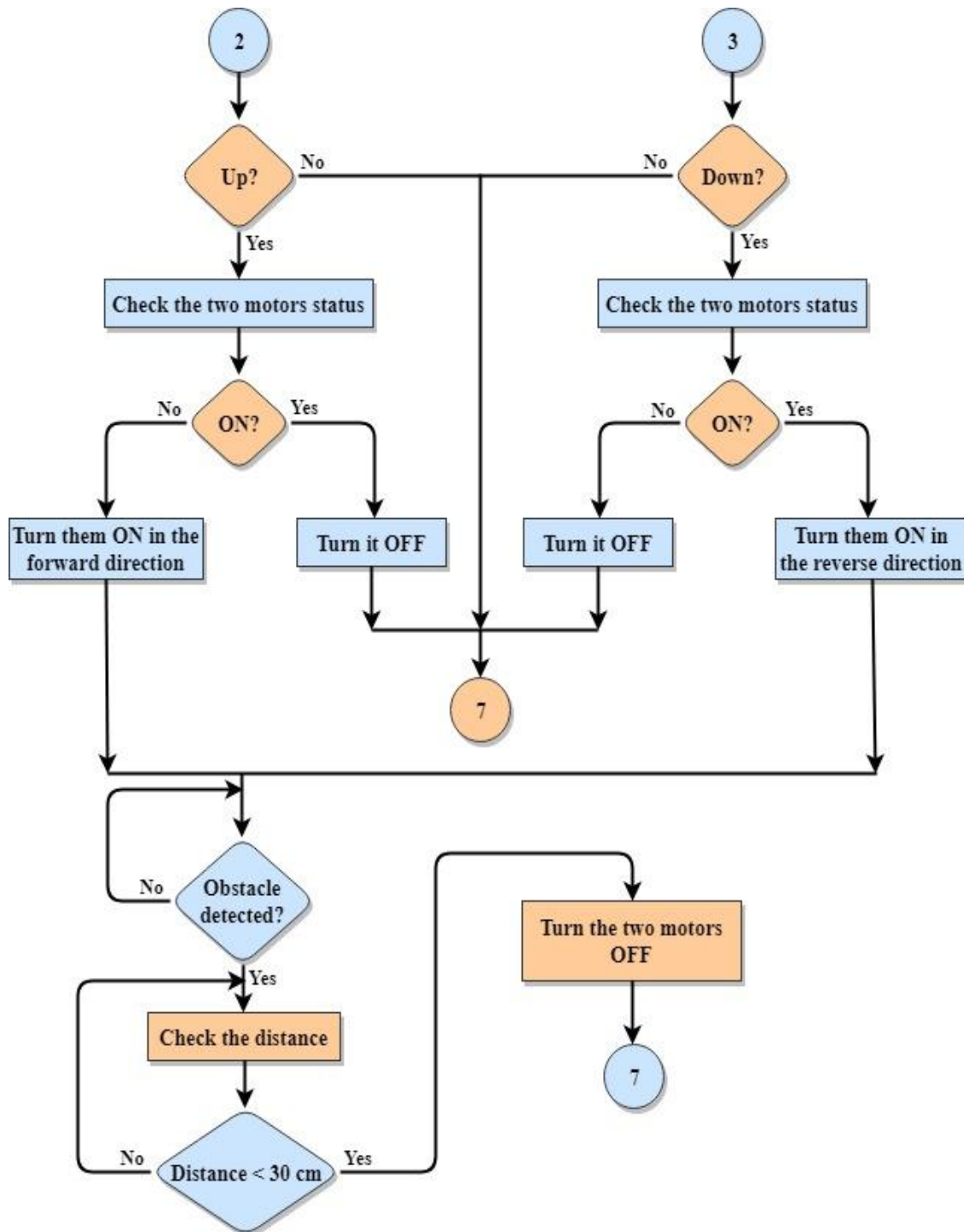


Fig. 2. (c) The hardware side: the up and down gaze directions responses.

The application can detect two types of eye actions including the gaze direction detection and eye blinking detection which then send using Wi-Fi connection to the hardware side for controlling purposes. The "right" and "left" gaze directions used for turning the car to right or left respectively Whereas the "up" and "down" gaze directions used for moving the car in forward or backward direction respectively. The blinking action is used for eye-tracking and car-controlling deactivating. Two ultrasonic sensors have been

used for obstacles detection during forward and backward car movements.

B. The software side

It is the core part of the proposed system which is a mobile application developed to track the human eye and detect different actions from it. The designed application interface is depicted in Fig. 3.



Fig. 3. The application user interface.

The application is designed under the android studio IDE to be applied on android based smartphones and tabs. A machine learning based algorithm which is viola-jones algorithm has been adopted in this application processing. It is an algorithm used to train a classifier on numerous positive and negative images. It is used to detect a specific object in these images by applying scanning process several times at different scales to detect different sizes of the object [10].

At the other hand, Haar-like features based cascade classifiers have been used because of the short time needed to compute an extensive set of features with efficient results of detection related to using this classifier [11]. As aforementioned, working on the application image processing is accomplished using the following stages respectively:

1) *Face detection stage*: A trained Haar cascade classifier for face detection has been used in this stage. Wherein, the human face is detected and surrounded by a rectangle in real-time video as shown in Fig. 4(a).

2) *ROI calculation stage*: In this stage, the Region of Interest (ROI), the eyes-region, is approximately determined by calculating the coordinates of the upper left corner of the eye region (x, y) , the width and the height depending on the dimensions of the detected face rectangle by using (1), (2) and (3), respectively, assuming that the region of eyes occupies $1/3$ of the face height and $3/4$ of the face width [12]:

$$EyeArea_{(x,y)} = (FR.x + \frac{FR.w}{8}, FR.y + \frac{FR.h}{4.5}) \quad (1)$$

$$EyeArea_{width} = \frac{3}{4} FR.w \quad (2)$$

$$EyeArea_{height} = \frac{FR.h}{3} \quad (3)$$

Thus, the ROI can be determined and located as shown in Fig. 4(b). The next step is to split the resulted eye region into two regions as shown in Fig. 4(c). This is for distinct the two eyes from each other and shrink the ROI as a consequence, easing the detection of both eyes separately in the next stage. This is accomplished by following (4), (5) and (6) [12]:

$$RightEye_{(x,y)} = (FR.x + \frac{FR.w}{2}, FR.y + \frac{FR.h}{4.5}) \quad (4)$$

$$RightEye_{width} = \frac{7}{16} FR.w \quad (5)$$

$$LeftEye_{(x,y)} = (FR.x + \frac{FR.w}{16}, FR.y + \frac{FR.h}{4.5}) \quad (6)$$

Where, FR is the face rectangle upper left corner point. Moreover, w and h are the width and the height of the rectangle respectively. The width of the left eye is computed using (5), whereas the heights of both eyes are separately computed using (3).

3) *Pupil detection stage*: The proposed algorithm applies another Haar cascade classifiers on the two regions, detected in Fig. 4(c), for separately eyes detection. The right and left eyes are detected without the eye-brow and the unwanted skin regions [12]. After that, the pupil's point is obtained which is the darkest point in that region. Then, it is surrounded by a circle to distinguish and track it as shown in Fig. 4(d).

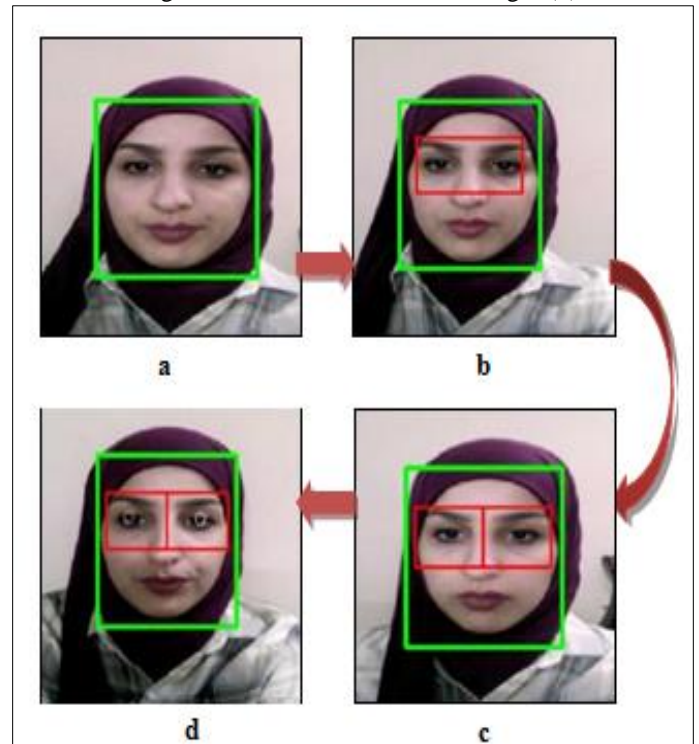


Fig. 4. The first three stages of processing: (a) is the first stage, (b-c) for the second stage and (d) is the third stage.

4) *Eye-actions detection*: It is the final stage in the application processing. The two eyes regions detected in Fig. 4(c) are separately divided into a 5x5 grid, as depicted in Fig. 5, to use them as thresholds of the pupil's movement directions range, used in detecting the following two actions:

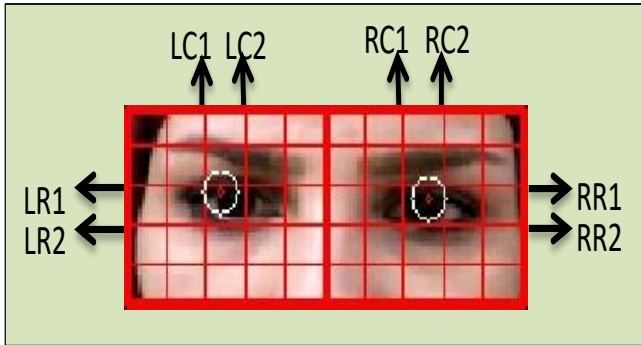


Fig. 5. The division of both eyes regions into a 5x5 grid.

a) *Gaze direction detection*: Gaze direction detection: Four eye gaze directions have been detected at this stage including up, down, right and left. This is accomplished by depending on the taken thresholds of the right eye solely, as expressed in the pseudo code of Fig. 6. Ultimately, a specific button is pressed after detect its related gaze direction.

b) *Blink detection*: The absence of the two eyes pupils means that the eye blinking action is performed. Nonetheless, this type of actions can be obtained involuntarily by the user. Therefore, this action is detected under the condition of the non-existence of pupil circles of both eyes with some conditions that can be summarized as:

- The user's gaze must be at the forward position first before blinking, to avoid the overlap with the gaze directions action.
- The user's eye must be continually closed for 2sec.
- The blink must be hard enough to guarantee the eyelashes disappearance. This is to avoid overlapping with the condition of detecting the pupil depending on

the darkest point which is the eyelashes in the case of the closed eye.

```

/*detect the eye gaze direction depending on the right eye*/
Let: RP = right eye pupil
    LP = left eye pupil
    RPC = right pupil circle's center
    GD = gaze direction
/*in addition to the rows and columns names mentioned in figure (5) */
If RP && LP are detected then
If RPC >= RC2 then
    GD = "right direction"
elseif RPC <= RC1 then
    GD = "left direction"
elseif RPC <= RR1 then
    GD = "up direction"
elseif RPC >= RR2 then
    GD = "down direction"
endif
endif
endif
endif
    
```

Fig. 6. Pseudo code of gaze direction detection method.

C. *The hardware side*:

This part mainly consists of a robot car that is controlled using the attached Raspberry Pi microcomputer. This section can be divided into three sub-sections as follows:

1) *Car aggregation and setup*: As illustrated in Fig. 7, a two-wheel motor driver (2WD) robotic car is adopted in this system with other hardware peripherals that can be listed as follows:

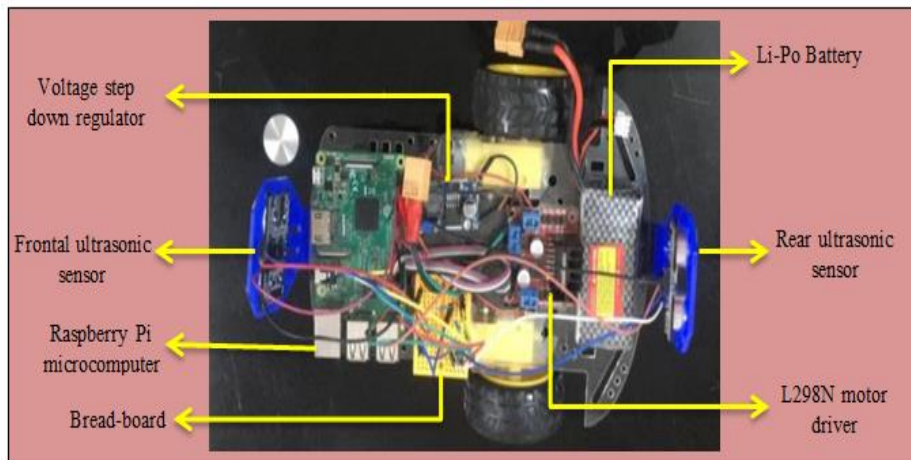


Fig. 7. Car module hardware requirements.

a) *The Raspberry Pi III model B microcomputer:* As shown in Fig. 8, it is a small, credit card-sized computer developed in the United Kingdom by the Raspberry Pi Foundation. The module is well equipped with different kinds of input/output ports and it performs quite well in the tasks that are related to the proposed system. It is used to control the movement of the car's motors [13].



Fig. 8. Raspberry Pi III model B.

b) *The L298N motor driver:* It is a high current, high voltage dual full-bridge driver designed to accept standard TTL logic levels and drive inductive loads such as relays, solenoids, DC and stepping motors [14]. Fig. 9 shows the adopted motor driver model.



Fig. 9. L298N motor driver.

It is used in this system to drive the two-car motors in different directions. Li-Po battery is used as a power supply for the motor driver.

c) *The ultrasonic sensors:* For obstacles detection, two ultrasonic sensors are used one in the front and back of the car.

d) *The voltage step down regulator:* It is used to supply the Raspberry Pi with its required voltage from the motor driver. The used model of this regulator is depicted in Fig. 10.



Fig. 10. Voltage step down regulator.

2) *Controlling process:* The central part of controlling in this system is the Raspberry Pi. Its built-in Wi-Fi communication feature is exploited, for establishing a Wi-Fi connection to be used as an access point, connected to the mobile phone. Raspberry Pi receives the command of the

detected eye action from the software side and then controls the movement of the car using the motor driver.

The car can move forward and backward as well as turning left and right depending on the incoming command of the detected gaze direction. Receiving the same command twice respectively, regardless of the time period between the two commands, means that the first one is for moving toward the intended direction and the second one is to stop this movement.

In addition, two different stopping commands can be applied to the car. The first one terminates the current movement direction of the car which accomplished by receiving the same command twice. The second one stops the car which accomplished in the case of blinking receiving.

3) *Obstacles avoidance:* As mentioned above, two ultrasonic sensors are used to detect the obstacles by measuring the distance between the car and the faced obstacles. The Raspberry Pi fully stops the car if the measured distance is below than (30 cm) to be waiting until the obstacles are removed or the user sends another command. The system can detect the obstacles during the forward and backward directions movements.

IV. EXPERIMENTAL RESULTS

Different experiments are conducted to the determination of the limits, rules and some guides for using the introduced system. Because the proposed system deals with a camera as a sensor, one of the important factors to be test is the intensity of light. In addition, the distance between the user's eye and the camera is another important factor that must be checked and determined. This is due to the architecture of the system that deals with the pixels positions, to determine the gaze direction, which affected by this distance.

Before illustrating the adopted experiments, a major point is necessary to be stated which is the patient's head position. It needs to be fixed straight forward in front of the device's camera to obtain perfect detection accuracy. In addition to that the device is also must be fixed by an adjustable holder.

The proposed system was tested and evaluated on Huawei T1-701u tablet with Quad core 1.2 GHz, 2MP camera and 600x1024 resolution which running Android 4.4.2 version.

Different cases has been studied and tested to check the effect of both the light and the distance factors which can be concluded as follows:

- The distance between the patient's eye and the device's camera: To determine the optimal distance between the eye and the camera, different values for this distance are supposed and applied on the system with no change on the ambient light intensity which is the room's florescence lamps. The florescence lamps are not directly faced the user's eye nevertheless it already placed near the room's roof.

Case 1: 45cm, florescence lamps

In this case, 45cm is the distance between the patient's eye and the device's front camera under only the florescence lamps.

The obtained overall accuracy is 52%. The blink detection was better than the gaze direction detection because the gaze direction detection depends on the difference between pixels' positions which vary depending on the distance between the face and the camera.

Case 2: 50cm, florescence lamps

In this case, the distance is increased by 5cm amongst the first case with maintaining the same light intensity of the first case. In this experiment, the obtained accuracy is 78%. Wherein, the gaze direction detection and the blinking detection dependently give acceptable accuracy.

Case 3: 55cm, florescence lamps

Due to the increased distance between the camera and the user's eye, and hence the difficulty of detecting the pupil point, the obtained accuracy of this experiment was 62% which is lower than the accuracy of the previous experiment.

- The intensity of the ambient light: The determining of the intensity of the ambient light during using the system is another challenge in the system test. To do that, a table lamp with 2 levels of intensities, adjusted manually, is used. From the previous experiments, because the optimal accuracy is obtained when the distance is 50cm, it was considered as the optimal distance and used as a constant factor during testing the

intensity of light effect. The lamp is placed in front of the user's face above the device to determine the effect of the high intensity of light when it directed to the eye.

Case 1: 50cm, (florescence lamps and table lamp's low intensity)

In addition to the ambient light of room's florescence lamps, the low intensity of the table lamp is used which faces the user's face. Under the conditions of this case, the overall system accuracy is 56%.

Case 2: 50cm, (florescence lamps and table lamp's high intensity)

This experiment proved that as much as the intensity of light increased, the accuracy of the system will decrease. In this experiment, the system is in worst case because of the increased brightness on the pupil and hence the difficulty of detecting the required pupil.

A. Car Locomotion Test

As mentioned previously, the car can move in four different directions (forward, backward, right turn and left turn) depending on the user gaze direction detected at the mobile application side (software side). Fig. 11-14 express the proposed system of the four detected gaze directions and the response of the application interface with the car movement indicated by the white LED brightness.

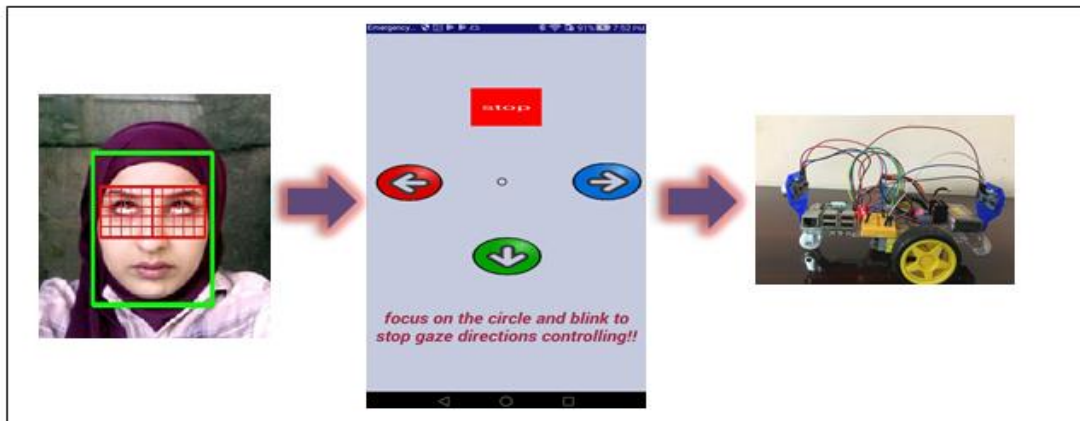


Fig. 11. Forward movement direction.

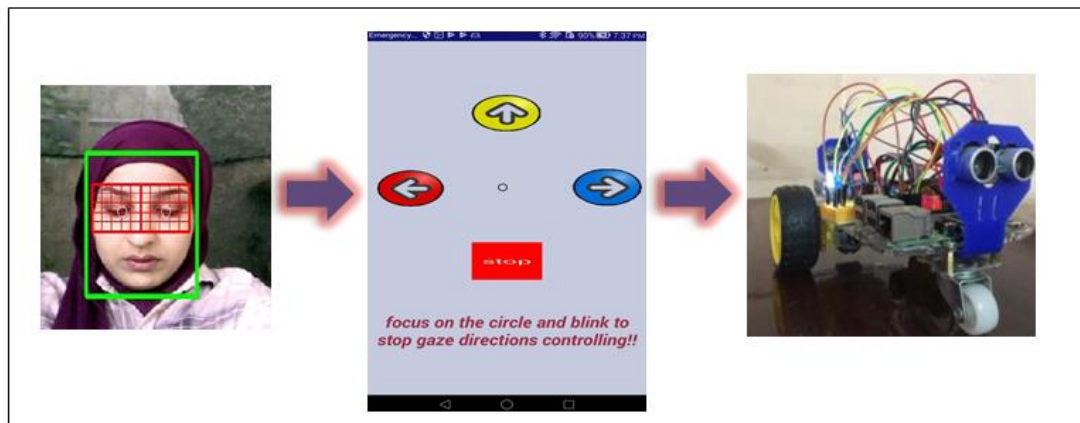


Fig. 12. Backward movement direction.

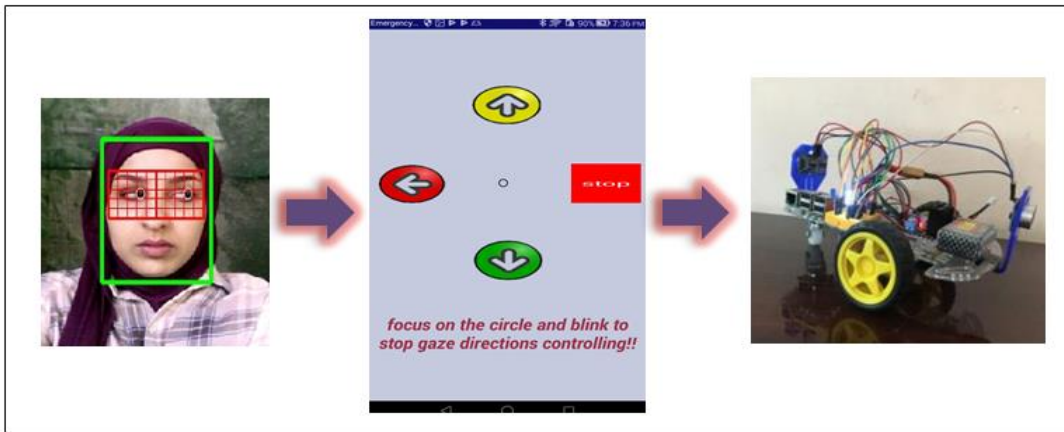


Fig. 13. Right-turn movement direction.

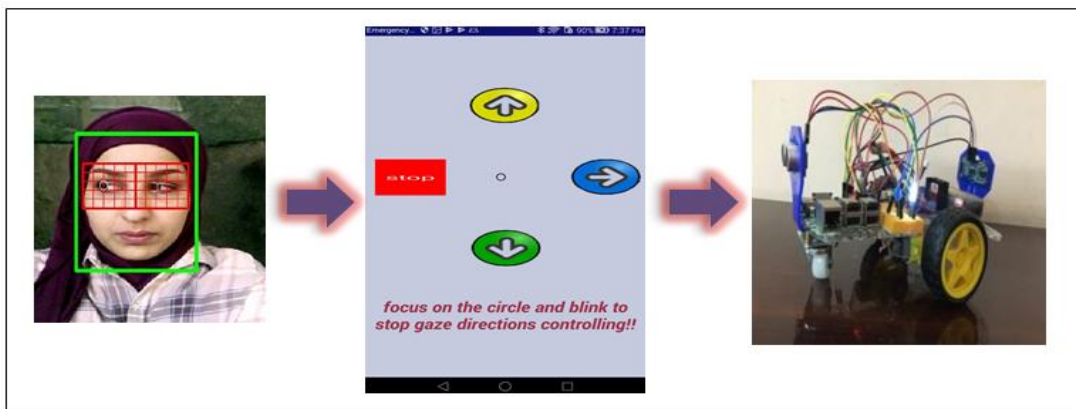


Fig. 14. Left-turn movement direction.

B. System Deactivation Test

If the user needs to deactivate the overall system and being free in his/her eyes motion, he/she needs to use the blinking action. In this case, the system is deactivated and the car has not received any moving action until the next blinking received to reactivate the overall system. The blinking action is depicted in Fig. 15.

C. The Obstacles Stop Test

From the experiments, the car has the ability to perfectly detect the obstacles during forward and backward movement directions using front and rear ultrasonic sensors. Fig. 16 shows the car is stopped after obstacle detection in two states forward and backward which indicated by the blue LED brightness (the white LED OFF state indicate the OFF state of directions movement).

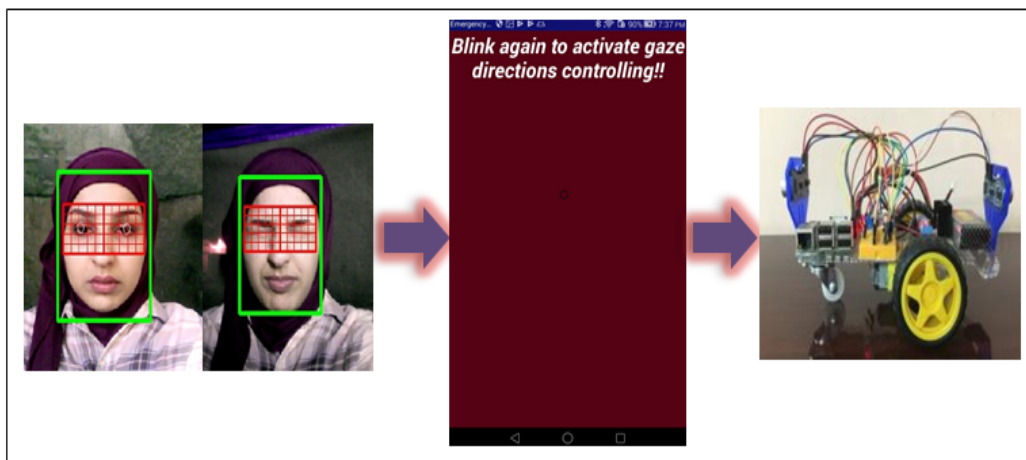


Fig. 15. Stop car movement.

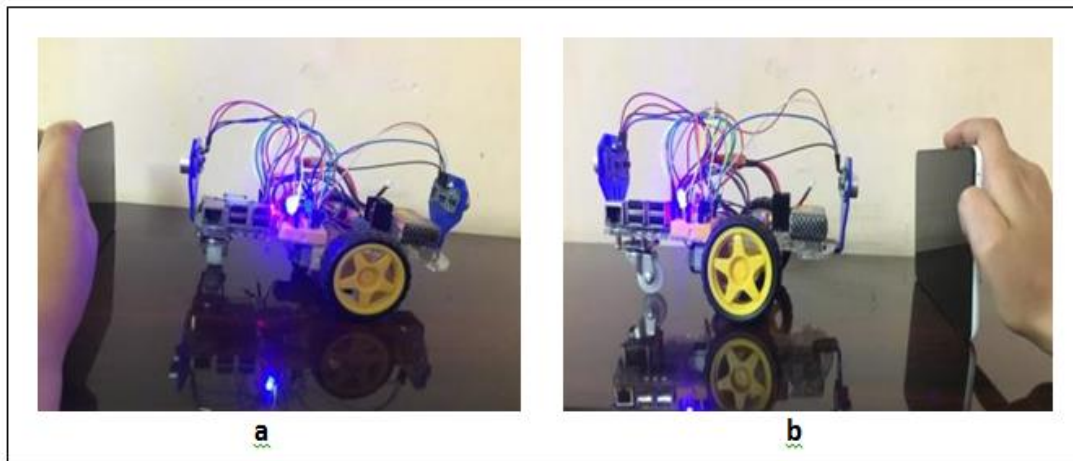


Fig. 16. Obstacle avoidance: (a) the frontal sensor and (b) the rear sensor.

D. Efficiency Evaluation of the Proposed System

There are different factors can impact on the system efficiency. Some of them have effects on the interface between the user's eye and the application, while others impact on the communication between the application and the car. Ultimately, all of them have effects on the system performance in direct or indirect impaction ways. As illustrated in the previously mentioned experiments, *the light intensity* is a very important factor that can affect the efficiency of the application and hence the overall system. Different light intensities have been tested during the experiments; ultimately, it is proved that the intensity of light must not be very high due to its effects on the size of pupil which causes difficult pupil detection [15]. Furthermore, the increased intensity of light can make the pupil brighter which causes incorrectly detection. This is because the pupil does not remain dark as needed. In addition, the low light intensity causes low detection accuracy because of the darkness of the shadow affected on the eye region and thus make the pupil difficult to detect. Therefore, the ambient light, during the system usage, must be neither directed to the eye nor a poor intensity light. Instead, it preferred to use a light source far from the eye and with acceptable intensity such as the florescence lamps.

Since, the efficiency of the overall system is affected by the light intensity; its accuracy can vary depending on the ambient environment (outdoor or indoor). Therefore, several experiments have been applied on the system in outdoor to test its accuracy. As illustrated in Table I, the confirmed results prove that the system accuracy is decreased in the outdoor environments due to the increased light reflected by the eye especially under the sun light. The solution for this issue can be performed by installing an umbrella above the wheelchair.

Different experiments reveal *the optimum distance* required between the user's eye and the mobile's camera which must be almost 50cm indoor and 54cm outdoor to obtain optimal eye-actions detection.

Although the system runs with accepted accuracy results under low device's camera resolution and processing speed like the Huawei phone used in the experiments of section 4.2, the high camera resolution and fast processing speed are preferred

to obtain better and perfect results. Therefore, in addition to the Huawei tab, a second device has been used through system testing which is a GALAXY tab with Quad core 1.6 GHz, 2MP camera and 800 x1280 resolution of the camera running Android 4.2 version. As shown in Table I, the obtained accuracy, using this device, outperforms the accuracy using the Huawei tab because of the powerful mobile specifications (especially its camera efficiency). Thus, *the device's processing speed* and its *camera resolution* have another impact on the system efficiency.

There are different factors can impact on the system efficiency. Some of them have effect on the interface between the user's eye and the application. Others impact on the communication between the application and the car. Ultimately, all of them have an impact on the system performance regardless of whether it was direct or indirect impact.

TABLE I. THE OBTAINED ACCURACY RESULTS.

Action	On Galaxy phone		On Huawei phone	
	Indoor accuracy	Outdoor accuracy	Indoor accuracy	Outdoor accuracy
Right gaze direction	90%	70%	90%	70%
Left gaze direction	100%	80%	90%	60%
Up gaze direction	80%	70%	70%	50%
Down gaze direction	70%	70%	70%	60%
Blinking	70%	90%	70%	90%
The whole system accuracy:	82%	76%	78%	66%

From this table, we can observe that the low detection accuracy ratio of the up and down directions as compared to the right and left directions is appeared because the horizontal line of the eye is wider than the vertical line. As a consequence, the horizontal movement range of the eye from left to right and vice versa is faster to detect than the vertical movement range between up and down directions. Moreover, this table also

shows a difference in overall accuracy when the device's specifications are varied. In addition, a noticeable difference is seen between the usage in outdoor and indoor environment accuracy. Wherein, the system wins accuracy in the indoor environment of about (6%-12%) as compared to the gained accuracy in the outdoor environment. The high accuracy is obtained by using the system indoor with powerful device's specifications which was almost 82% that shows the superior performance of the proposed system considered being optimal for disabled people's usage.

At the other hand, another major factor that can limit the system accuracy is *the user's head position*. The proposed system does not employ a feature for supporting the tilted head in the tracking of the eye. Therefore, the user's head must be kept straight forward in a stable manner.

The speed of the Wi-Fi connection and the microcomputer processing speed are other critical points in the system efficiency testing. Because the system is mainly designed for controlling aspects for a very sensitive group of the community (disabled people), it must be fast in receiving the control command and directly process it with fast latency after signal transfer without any noticeable delay. In the proposed system, the used Wi-Fi connection strategy, accomplished using socket connection makes the system active in real-time, and responsive with fast latency time between sending and receiving the command. In addition to that, the Raspberry Pi III model B is run with 1.2 GHz processor speed which is fast as compared to the old versions of related embedded systems.

V. CONCLUSION AND SUGGESTIONS FOR FUTURE WORK

An affordable control system based on eye-tracking technology was designed and implemented to control the movement of a mini robotic car as a prototype of the wheelchair. The system consisted of two main parts: the eye tracking application (software side) and the robotic car (hardware side). The eye tracking application was designed under an Android operating system that adopted the device's front camera as a sensor that captures a real-time video. The captured videos were processed using Viola Jones algorithm with Haar cascade classifiers to detect four gaze directions in addition to the blinking action. The resulting signals were sent via Wi-Fi to the second part of the system which is Raspberry Pi microcomputer attached to the car. Depending on the incoming signals, the Raspberry Pi controlled the movement of the used car motors. Wherein, the right and left gaze directions turned the car to the right and the left whereas up and down gaze directions moved the car forward and backward. In addition, the blink action has been used to fully stop the car. Experiments on the system have shown that it is easier to control the car using only the eye gaze directions with high performance in term of accuracy and precision. Moreover, it is important to note that the implementation cost of the proposed system was really low in comparison with manufactured similar systems. This due to the use of low cost embedded system and the proposed mobile application can work with low specifications phones or tablets.

To enhance the system by adding additional features as a future work, different suggestions are listed below:

- Making the system responsive with the people who are in black skin and wearing eye-glasses on their eyes by developing the proposed algorithm.
- Providing a feature on the application that enable it to calculate the distance between the device's camera and the user's eye and prevent running the system unless this specified distance is regulated.
- Integrating the iris recognition technology alongside with the application to obtain a full system security.
- Designing the same application using swift programming language to be able to support MAC OS based devices such as iPhones and iPads.
- Integrating a system on the wheelchair side that responsible of auto-driving the wheelchair of the patient from place to place automatically and based on information received from the improved application.

REFERENCES

- [1] R. Veerati, E. Suresh, A. Chakilam and S. P. Ravula, "Eye monitoring based motion controlled wheelchair for quadriplegics", *Microelectronics, Electromagnetics and Telecommunications*. Springer, Singapore, pp. 41-49, 2018.
- [2] J. A. Athanasou, "Spinal Cord Injury", *Encountering Personal Injury. Studies in Inclusive Education*, pp. 157-164, 2016.
- [3] B. Buvanswari and T. Kalpalatha Reddy, "Eye scrutinized wheel chair for people affected with tetraplegia", *International Journal of Computer Science, Engineering and Information Technology*, vol. 5, no. 2, pp. 15-24, 2015.
- [4] R. Cerejo, V. Correia and N. Pereira, "Eye controlled wheelchair based on arduino circuit", *International Journal of Technical Research and Applications*, vol. 3, pp. 94-98, 2015.
- [5] J. J. Magee, M. Betke, J. Gips, Matthew R. Scott, and Benjamin N. Waber, "A human-computer interface using symmetry between eyes to detect gaze direction", *IEEE transactions on systems, man, and cybernetics—part a: systems and humans*, vol. 38, no. 6, pp. 1248-1261, 2008.
- [6] H. R. Chennamma and X. Yuan, "A survey on eye-gaze tracking techniques", *Indian Journal of Computer Science and Engineering*, vol. 4, no. 5, pp. 388-393, Nov. 2013.
- [7] B. Champaty, J. Jose, K. Pal and T. A., "Development of EOG base human machine interface control system for motorized wheelchair", *IEEE Annual International Conference in Emerging Research Areas: Magnetics, Machines and Drives (AICERA/iCMMD)*, pp. 1-7, 2014.
- [8] R. F. Ghani and H. S. Hassan, "Human computer interface for wheelchair movement", *Baghdad Science Journal*, vol. 14, pp. 437-447, 2017.
- [9] X. Gao and L. Shi, "The design of robotic wheelchair control system based on hand gesture control for the disabled", *IEEE International Conference on Robotics and Automation Sciences*, pp. 30-34, 2017.
- [10] V. Paul and J. Michael, "Rapid object detection using a boosted cascade of simple features", *Computer Vision and Pattern Recognition, IEEE Computer Society Conference*, vol. 1, pp. 511-518, 2001.
- [11] A. Kasinski and A. Schmidt, "The architecture and performance of the face and eyes detection system based on the Haar cascade classifiers", *Pattern Analysis and Applications*, vol. 13, no. 2, pp. 197-211, 2010.
- [12] B. Javidi, *Image recognition and classification: algorithms, systems, and applications*. CRC press, 2002.
- [13] RASPBERRY PI, 7 March 2018 at 05:00 AM, Available on: <https://www.raspberrypi.org/products/raspberry-pi-3-model-b/>
- [14] L298N Motor Driver Board, 7 March 2018 at 05:20 AM, Available on: https://www.geeetech.com/wiki/index.php/L298N_Motor_Driver_Board
- [15] Resolution of an eye, 8 Feb. 2018 at 05:30 PM, Available on: https://www.wikilectures.eu/w/Resolution_of_an_eye.

A Study on Usability Awareness in Local IT Industry

Mahmood Ashraf¹

Department of Computer Science
Federal Urdu University of Arts Science & Technology,
Islamabad, Pakistan

Lal Khan²

Department of computer science and IT
University of Lahore, PakPatan, Pakistan

Muhammad Tahir³, Ahmed Alghamdi⁴,
Mohammed Alqarni⁵

Faculty of Computing and Information Technology
University of Jeddah, Jeddah, Saudi Arabia

Thabit Sabbah⁶

Faculty of Technology and Applied Sciences
Al-Quds Open University, Ramallah, Palestine

Muzafar Khan⁷

Department of Computer Science
COMSATS University, Islamabad, Pakistan

Abstract—Usability awareness receives more consideration by industry professionals and researchers throughout the world, but it is limited in Pakistan. This study reports survey results of the current state of usability awareness in the local Information Technology (IT) industry. Forty participants – IT practitioners from IT industry – were involved in the study. We used Usability Maturity Model (UMM) and content analysis methodology to discover the current status of usability awareness. The results indicate that 1) almost half (18 out of 40) of the participants were unaware of the term usability and related concepts, 2) there is shortage of HCI/Usability professionals in organizations, 3) most of the software companies were at unrecognized level of UMM and 4) they were also not interested in usability because of limited or no budget for it. The study also reveals a gap between usability awareness and its perceived usefulness among IT professionals.

Keywords—Usability; usability awareness; human-computer interaction (HCI); HCI practitioners; Pakistan IT industry

I. INTRODUCTION

Human Computer Interaction (HCI) deals with the knowledge to assist human's physical and mental skills for the ever-developing technology [1]. In Europe and United States, HCI is playing a vital role in IT industry since 1980 [2]. One major topic of HCI is the study and practice of usability. Usability is a quality characteristic that includes many factors e.g. a product should be easy to learn and use [3]. The scope of usability is not limited to user interface only rather it deals with the entire system [4].

A lot of work related to usability has been done e.g. the study [5] provides an overview of different usability evaluation methods used for web applications. Another study [6] also presents the various user experience (which includes such as usability) evaluation technologies for software applications. Similarly, in a recent systematic literature review [7], reporting mechanisms for usability defects are summarized and discussed.

Despite of all the work done, awareness about the usability methods and practices is not good and the role of HCI practitioners is overlooked [8]. According to [9], usability awareness means “designing for a sustainable world”. In general, awareness guides towards the maturity, eight levels of usability maturity (like software process maturity) for an organization are proposed in [10]. If the system is developed by the organization having the high level of usability maturity, the end user will use the system without any training required for it [11].

There is continuous appreciation and acceptance for usability in organizations. Many studies e.g. [11]-[14] have been conducted to assess the usability maturity of industry in different countries. These studies are mostly conducted in developed countries like Germany, Japan, and Israel. There is a need to conduct similar studies in other countries particularly the developing ones to assess the usability awareness and maturity in general. This study aims to explore the current state of usability awareness in the local IT industry of Pakistan. The findings of this study may guide the academia and other concerned authorities of the developing countries (Pakistan in particular) for better planning to cater the usability needs of the industry.

II. RELATED WORK

Many studies have been conducted about the usability awareness and practices in different parts of the world. A survey was conducted in Malaysia about the usability awareness of IT and non-IT practitioners [8]. Out of total 72 participants, 23 IT practitioners, 27 IT scholars and 22 non-IT experts participated in this study. The results revealed no major differences of usability awareness among IT practitioners, IT Scholars and non-IT professionals. Study participants also considered usability as God-gifted skill and common-sense knowledge for both IT and Non-IT staff. Another survey was conducted in UAE where the participants were IT managers, marketing professionals and end users [11]. The results revealed that the participants had introductory

knowledge about usability and they did not study it as a significant role-player for software development. Furthermore, no user involvement in design phase and unavailability of usability staff were reported.

A total of 72 participants were involved in a survey conducted in Korea [15]. The participants were software developers, usability and user interface professionals. The results showed that usability had not been applied in projects, but the increased awareness of usability was stated. The lack of usability professionals, time and cost-effective usability methods were also reported as the main problems. A study was conducted in France to find the perspective of HCI professionals about usability methods [16]. The results revealed that professionals were aware of usability methods after many years of working in industry. Results also showed that young experts didn't use usability methods and approaches in developing user interfaces at their earlier stages of profession. The specialists from engineering schools were more aware of usability as compared to other graduates.

The primary aim of the study was to find the practices, awareness level and perceptions of User Experience Professionals (UXPs) about web accessibility in Turkey [17]. An online survey study was performed to meet the primary goal of the study. The finding indicates that UXPs have confidence in that they have enough education and training regarding web accessibility. But, they were not aware with web accessibility standards and were not considering to apply them in their projects. They think that considering the web accessibility is the responsibility of the project managers. The study conducted in German institutes highlighted the current practices, usability awareness, perceptions of usability and networking strategies in Germany [18]. The institutions showed the openness to accept recommendations and suggestions to optimize the interface designs/usability in general. The shortage of budget was considered the main problem to carry out the usability related tasks.

The aim of the study was to find the current practices in the Nigerian software industry [19]. A survey study and semi-structured interviews were performed in local software houses. The results show intermediate level of usability awareness and the limited knowledge of HCI practices. Another study was conducted to find the understanding and awareness level of HCI [20]. The results show that HCI practices and processes are at their beginning level in most of software houses. The results also indicate that, during software development, the end user involvement was also not considered. Furthermore, it was found that there was lack of HCI knowledge transfer to university students.

In a survey conducted in Brazil, the attempt was to find out the opportunities and challenges in HCI education [21]. 109 participants contributed in the survey. One of the biggest challenges was to get a good HCI position in industry after completing the degree. The findings also highlighted the importance of continuously updated knowledge about the latest technology development and active collaboration among the faculty members. Another study in New Zealand indicated the absence of proper schooling, information and skills about usability approaches, procedures and practices between

designers and developers [22]. According to a study, usability community in Russia had been facing many problems e.g. lack of professional training and standards and insufficient awareness among professionals [23].

III. METHODOLOGY

A survey was conducted to assess the level of usability awareness and maturity in the local IT industry. For this purpose, a questionnaire was designed based on the Usability Maturity Model (UMM) [24] which helped to assess usability as the ability of an organization. UMM defines six levels of usability maturity i.e. Level X (unrecognized), A (recognized), B (considered), C (implemented), D (integrated), E (institutionalized). Level X is the lowest level which indicates the organization's negligence towards usability whereas Level E is the highest level of usability maturity within the organization. The following questions were asked from participants to find out the current state of usability awareness and the results were mapped on UMM maturity levels.

- ~ Is there any HCI/Usability staff in your organization?
- ~ Users are involved in design phase of the system?
- ~ Is there any budget allocated for usability related activities in your organization?
- ~ Does your organization consistently produce usable products?
- ~ Does the top management of your organization focus on design for human use?

A. Participants

The questionnaire was sent through email to professionals of 65 organizations located in different cities of Pakistan. In total, 40 participants responded with their feedback. The detail of their professional roles is given in Table I. Most of the participants (about 82%) were less than 30 years of age. The remaining (about 18%) participants were between 31-40 years of age. In terms of working experience, 2 participants (5%) had 10 years of job experience; 13 participants (32.5%) with 5 years; 9 participants (22.5%) had an experience of 3 years; 7 participants (17.5%) with 2 years of experience; and 9 participants (22.5%) had 1 year of job experience.

TABLE I. PARTICIPANTS' PROFESSIONAL ROLES AND AVERAGE EXPERIENCE

Professional Role	Number of Participants	Average Experience (in years)
Software developers	11	4.90
Software engineers	9	3.33
Managers	6	6.66
Senior executives	4	3.25
Software testers	4	3.5
System analysts	3	4.33
Graphic designers	3	3
Total	40	

IV. RESULTS AND DISCUSSION

The data we received as a result of survey is mostly qualitative and thus sits valid for content analysis [25]. The specific context of this data describes the IT industry's overall perception and awareness regarding usability.

A. Usability Awareness in General

Based on the separate question; "Have you ever heard of usability?" 18 participants out of 40 (45%) were not very familiar about the usability. Interestingly, most of them were computer programmers. The possible reason is their inclination towards coding rather than designing user interfaces which are not given the due priority (by the programmers themselves and also by the software development organizations). The remaining 22 participants (55%) were well-aware about the usability for product design.

B. Organizations' Capability to Handle usability Issues

The results of this study show that many organizations have no staff for handling issues in usability or in user interface design. In such organizations, software developers perform all tasks. Usability issues are not properly handled due to lack of usability practitioners. Furthermore, there seems to be no formal usability training, for software developers, to handle usability related issues.

C. Development of Usable Products and Budget for Usability

The development of usable product is directly related with the budget allocation for this purpose. In absence of budget allocation, it is difficult to hire usability experts or performing other usability related activities. More than 50% (23 out of 40) participants indicated that their organizations had no budget for handling usability related activities.

D. Usability Experts Inside Organization

It is important to know the availability of HCI/usability experts in an organization which indicates the seriousness of that organization towards usability. 15 participants (37.5%) informed about the unavailability of usability professionals in their organizations. It reveals that usability is not considered an important aspect in those organizations. The remaining 25 participants informed that their organizations have HCI professionals. These HCI professionals were hired for different positions and levels. 9 participants (out of 40, 22.5%) informed about the availability of user experience professional in their organizations. Five participants (12.5%) shared that use the titles of interaction designer and usability expert in their organizations. Two participants informed that the titles of usability engineer and HCI expert were used in their setups.

E. End user Involvement in System Design and Top Management Commitment Towards Usability

User participation and approval is valuable for system success [26]. There is a low chance for a system to be user-friendly if users are not involved in the designing phase of the system. An active and frequent user participation throughout the system development is a basic principle for User-Centered System Design (UCSD) [29]. 13 participants (32.5%) mentioned that their organizations did not involve users in system design. The rest of participants (67.5%) claimed that

their organizations focused on users' participation during the design process.

V. UMM MAPPING

Each question, described above, carried a weight equal to 1. The participants answers (Yes or No; and/or the indication of HCI/Usability staff in either Yes or No) were assigned equal weights and then summed up to maximum value of 5. The companies which achieved the score of 5 have the highest level of usability maturity. Fig. 1 describes the percentage of IT companies and their achieved level of UMM maturity. It explains that 31% of the companies achieved the average level of UMM maturity while there were 6% of the companies who were at unrecognized level (X). If we consider level C as the threshold then 57% of the companies were below that threshold (by counting the results of level X, A, B and C respectively).

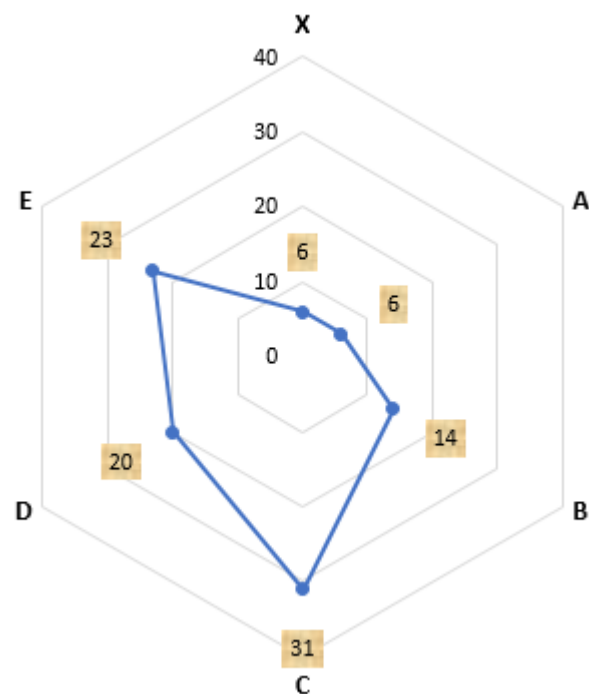


Fig. 1. The IT companies and their achieved level of UMM maturity.

VI. DATA ANALYSIS

This analysis provided main themes, categories and general insight on the usability awareness of the local IT industry. We performed the analysis on one part of survey data i.e. for the question; your perception of the concept of the usability. The scope of the question is limited to the data received from the professionals (population of this study) of local IT industry. We focused on these professionals because they work in industry where there is a need of applying usability and its principles. The other factors that describe the population include age, gender, profession, education and domain experience.

For the data (i.e. responses), received from participants, we used emergent and a priori coding techniques [27] to analyze it and found out interesting themes. In emergent coding we

performed the analysis of the responses and extracted the themes without considering existing HCI/Usability theories or models. It means, the resultant themes were actually the subjective list of participants-defined words. In a priori coding we performed the analysis of the responses and extracted the themes based on the usability and Human Computer Interaction concepts commonly found in the relevant literature. Two coders coded the responses independently in order to avoid the bias in the analysis. The result of this analysis is described in the Table II.

The data was then further cleaned and few entries were removed from responses. This was based on the participants' answers which were not very clear. For example, few responses were like this; "I don't have any idea", "Good", "don't know", etc. These responses were discarded because otherwise they would have created the anomalies in the final result. Thus, the total of 27 valid entries (Table II) were considered for further analysis. Both coders, independently, coded all 27 entries and found out similarity in 12 emergent themes and in 9 a priori themes (highlighted in Table II). While, only 7 entries were coded similarly in both types of coding.

TABLE II. EXTRACTION OF THEMES FROM PARTICIPANTS' RESPONSES

Parti ci.	Coder 1		Coder 2	
	<i>Emergent coding</i>	<i>A priori coding</i>	<i>Emergent coding</i>	<i>A priori coding</i>
3	importance	usability	usability is important for easy use	ease of use
4	user-friendly, ease of use	user-friendly	ease of use by focusing users' real tasks	ease of use
6	quality improvement	usability	improved quality product and user's comfort	user satisfaction
7	ease of use	ease of use, understandability	easy to understand and use	ease of use
8	ease of use	ease of use, understandability	easy to understand and use	ease of use
10	customized-definition of usability	usability	best way to perform the task	efficiency
11	customized-definition of usability	usability	effective and efficient way to achieve goals	efficiency
12	customized-definition of usability	usability	ease of use and learnability	ease of use
13	readability	readability	understandability of design	understandability
14	user understands without much help	Standards	product development according to standards; no user training	learnability

			for product use	
15	customized-definition of usability	Usability	ability to use	ease of use
17	the way to apply usability	Categorization	product dependent	quality in use
18	quality improvement, usability in design process	quality vs usability	quality product	quality in use
19	subjective opinion	software development	important for product development	quality in use
20	ease of use	user-friendly	easy to use	ease of use
23	customized-definition of usability	Usefulness	best utilization	efficiency
24	understandability	user interface design	easy to understand	understandability
26	importance	Usefulness	useful	useful
27	utility, usable product	efficiency, satisfaction	efficiency, effective, satisfaction	efficiency
29	satisfaction	Satisfaction	user satisfaction	satisfaction
30	importance	user-friendly	user friendly and easy to use	ease of use
32	user-focused	user-centered-design	user centered design	user centered design
33	quality improvement	quality vs usability	good quality	quality in use
35	user understands without much help	Interaction	easy to understand	understandability
38	subjective opinion	Usability	more awareness in people	awareness
39	system enhancement	quality vs usability	better functionality of system	quality in use
40	system development	Usability	important for product development	user centered design

As the results from this qualitative data analysis and its interpretation could involve researchers (coders) subjectivity, thus we need to be sure to describe the data reliability. This was achieved by performing inter-coders reliability [25] in terms of % of agreement and disagreement as:

% agreement= the number of themes coded the same way by different coders/the total number of themes

$$\% \text{ agreement (emergent)} = 12/27 = 44\%$$

$$\% \text{ agreement (a priori)} = 9/27 = 33\%$$

The results indicate that both coders agreed on less than half of the themes. This further strengthens our argument of performing the coding independently and thus obtaining the results without bias. Furthermore, it describes that the identification of themes is probably dependent on the number and richness of responses along with the coders and participants knowledge in the same domain.

In order to converge our results to few categories we then compared the themes and extracted the occurrences of similar themes. This also helped in identifying and defining the relationships between themes, and thus creating a set of categories which is called the code list (or nomenclature) [25]. This code list provides a hierarchical structure of the themes with multiple levels of details. In the word cloud (Fig. 2) and code list (Fig. 3, [28]), the unique occurrences i.e. the most prominent words in word cloud are; usability, ease of use, understandability, and quality improvement. Some themes are reported more than once and do not add up to the total number described in the parent node of the code list. This word cloud and the code list helped us in identifying the usability-relevant terms and describe the overall attitude, trend and inclination of IT industry towards usability awareness in-situ.



Fig. 2. Word cloud: the relevant themes popular within local IT industry.

In Fig. 3, code list created using emergent and a priori coding, describes the concepts currently exist in local IT industry. The numbers in parentheses represent the number of occurrences counted in participants responses’.

We also calculated top-2 box score [30] for the question; “How useful do you think is the usability for you (and your organization)?” This question was asked after the survey and was measured on 5-point rating scale value (5- very useful1- not at all useful). Top-2 box score describes the participants who strongly and somewhat strongly agree to the statement/question being asked. The results for this question indicates that almost 47% of the participants considered applying usability is a useful measure, though, 55% of the participants were well aware of usability. This result when combined with the result of the question; “Have you ever heard of usability?” helps in comparing awareness and usefulness. This returns an encouraging result i.e. there exist a smaller usability awareness vs usefulness gap (Fig. 4) which could be filled by following the suggestions/guidelines presented in conclusion.

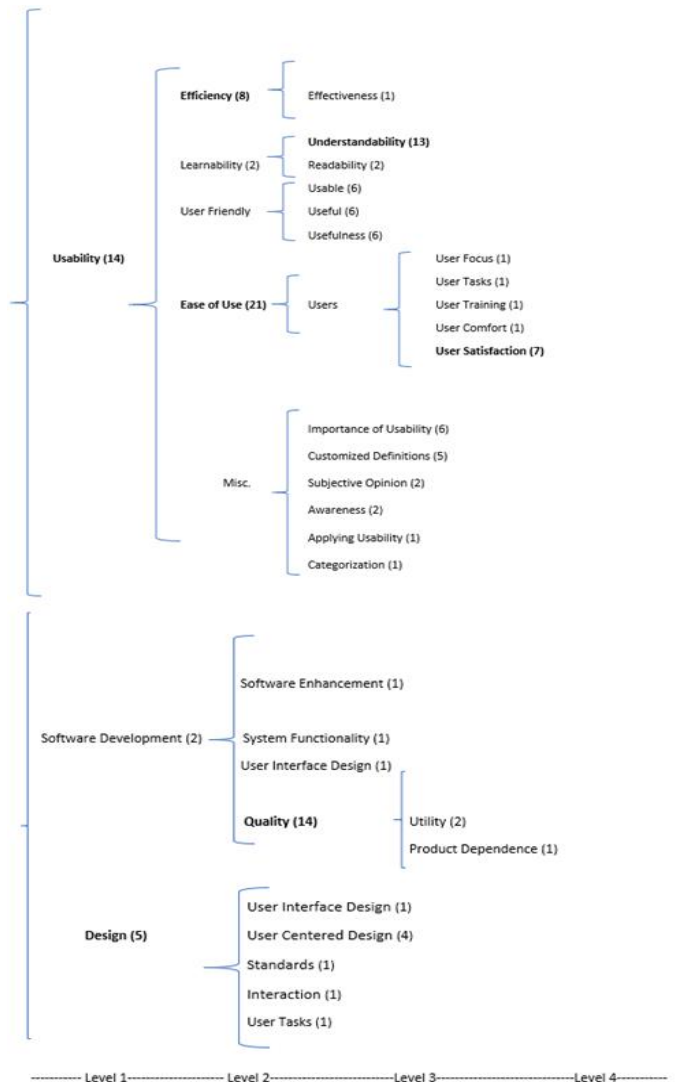


Fig. 3. The concepts (code list) currently exist in local IT industry.

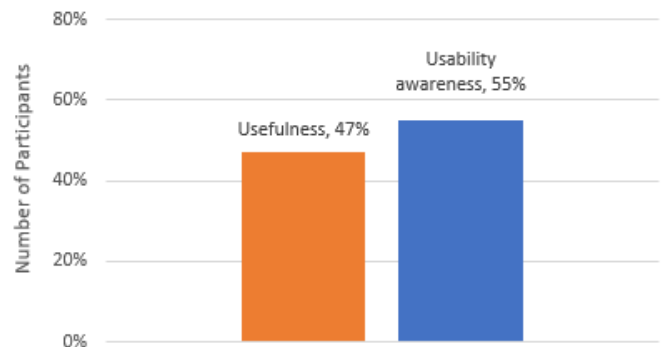


Fig. 4. Usability awareness vs usefulness gap.

VII. LIMITATIONS OF THE STUDY

Some themes in emergent and a priori coding columns (for both coders) look closely related but we didn’t count in similarity because both themes represent different levels of detail. For example, for participant 6, one theme is usability

(which is a general concept) while the other is user satisfaction (which is a specific concept within usability). The data analysis, performed by two coders, does not generate an acceptable level of agreement (this is <40% in average for both types of coding) which indicates the need of involving more coders in the analysis process and thus achieving a better reliability score. Furthermore, the word cloud and code list provide the relevant themes currently exist in IT industry, but they do not truly reflect an individual's or organization's intention or struggle towards achieving the goal (i.e. applying usability at their work place). Exploring this struggle and efforts - currently going on in the industry - is an important future research direction.

VIII. CONCLUSION

Usability awareness is currently limited in the local IT industry of Pakistan. Before performing this survey study, we assumed that IT professionals automatically gain usability perception and awareness during their work - because they have to deal with the evaluation of interactive systems - but the assumption was not true. It has been found that it depends upon their type of work and the organization's interest in applying usability techniques. Furthermore, in most of the organizations the end users were not involved in product user interface design phase. Lack of HCI/Usability professionals is also observed in most of the organizations. Therefore, based on the participants responses' few organizations could be ranked at recognized level of UMM. It was found that most organizations are not interested in usability because they have no budget for it.

We conclude that the usability perception, awareness and its importance can be achieved by 1) conducting training workshops describing the advantages of usability especially on the Return On Investment (ROI) of the company and 2) consulting higher education institutes and asking them to train students accordingly by focusing more on usability and 3) requesting top management to reserve budget for usability tasks. By following these suggestions, we are hopeful, that this will develop a rich usability culture in IT industry of Pakistan.

REFERENCES

[1] Shackel B (2009) Human-computer interaction-Whence and whither?. *Interacting with computers*, 21(5-6):353-366

[2] Organization for Economic Co-operation and Development (2004) ICT diffusion to business: National Peer review

[3] Bevan N (2009) International standards for usability should be more widely used. *Journal of Usability studies*, 4(3):106-113

[4] Seffah A, Metzker E (2009) On usability and usability engineering. In: *Adoption-centric Usability Engineering*. London (UK): Springer; pp.3-13.

[5] Fernandez A, Insfran E, Abrahao S (2011) Usability evaluation methods for the web: A systematic mapping study, *Information and Software Technology*, 53(8):789-817

[6] Rivero L, Conte, T (2017) A systematic mapping study on research contributions on UX evaluation technologies. In: *XVI Brazilian Symposium on Human Factors in Computing Systems*, Joinville, Brazil, pp. 1-10

[7] Yusop N, Grundy J, Vasa R (2017) Reporting usability defects: a systematic literature review, *IEEE Transactions on Software Engineering*, 43(9):848-867

[8] Hussein I, Mahmud M, Tap M (2011) A survey of usability awareness in Malaysia IT industry. In: *IEEE International Conference on User Science and Engineering*

[9] Rosenzweig E (2006) World usability day: a challenge for everyone. *Journal of Usability Studies*, 1(4):151-155

[10] Nielsen J (2018) Corporate UX Maturity: Stages 5-8. [accessed 2018 Feb 1]. <https://www.nngroup.com/articles/ux-maturity-stages-5-8/>

[11] Hindi M, Khalil A (2011) Usability practice and awareness in UAE. In: *IEEE International Conference and Workshop on Current Trends in Information Technology*

[12] Idyawati H, Seman A, Mahmud M (2009) Perceptions on interaction design in Malaysia. In: *Springer International Conference on Human-Computer Interaction*

[13] Eliav O, Sharon T (2011) Usability in Israel, In: *Global Usability Human-Computer Interaction Series Part II*, pp. 169-194

[14] Kurosu M (2011) Usability in Japan, In: *Global Usability Human-Computer Interaction Series Part II*, pp. 195-209

[15] Ji G, Yun H (2006) Enhancing the minority discipline in the IT industry: a survey of usability and user-centered design practice. *International Journal of Human-Computer Interaction*, 20(2):117-134.

[16] Rache A, Lespinet V, Andre M (2014) Use of usability evaluation methods in France: the reality in professional practices. In: *IEEE International Conference on User Science and Engineering*

[17] Inal Y, Rızvanoglu K, Yesilada Y (2017) Web accessibility in Turkey: awareness, understanding and practices of user experience professionals. *Universal Access in the Information Society*, pp.1-12.

[18] Peissner M, Rose K (2002) Usability engineering in Germany: situation, current practice and networking strategies. In: *European UPA conference*

[19] Ogunyemi A, Lamas D, Adagunodo R, Rosa B (2015) HCI practices in the Nigerian software industry. In: *Springer International Conference on Human-Computer Interaction*

[20] Ogunyemi A, Lamas D, Adagunodo R, Loizides F, Rosa B (2016) Theory, practice and policy: an inquiry into the uptake of HCI practices in the software industry of a developing country. *International Journal of Human-Computer Interaction*, 32(9):665-681

[21] Boscaroli C, Bim A, Silveira S, Prates O, Barbosa J (2013) HCI education in Brazil: challenges and opportunities. In: *Springer International Conference on Human-Computer Interaction*

[22] Apperley D, Nichols M (2011) Usability in Aotearoa/New Zealand. In: Douglas I, Liu Z, editors. *Global usability*. London (UK): Springer, pp. 237-245

[23] Burmistrov I, Kopylov A, Dneprovsky P, Perevalov Y (2004) HCI and usability in Russia. In: *Conference on Human Factors in Computing Systems*

[24] Earthy J (1998) Usability Maturity Model: human centeredness scale. *Information Engineering Usability Support Centres; Report No. INUSE D5.1.4s*

[25] Lazar J, Feng J, Hochheiser H (2017) *Research methods in human-computer interaction*, 2nd ed. Morgan Kaufmann Publishers Inc.

[26] Kujala S (2003) User involvement: a review of the benefits and challenges. *Behaviour and Information Technology*, 22(1):1-16.

[27] Corbin J, Strauss L (2014) *Basics of qualitative research: techniques and procedures for developing grounded theory*. 4th fourth ed. Sage Publications

[28] Feng J, Lazar J, Kumin L, Ozok A (2010) Computer usage by children with down syndrome: challenges and future research. *ACM Transactions on Accessible Computing*, 2(3):1-44.

[29] Gulliksen J, Goransson B, Boivie I, Blomkvist S, Persson J, Cajander A (2003) Key principles for user-centred systems design. *Behaviour and Information Technology*, 22(6):397-409.

[30] Tullis T, Albert W (2008) *Measuring the user experience: collecting, analyzing, and presenting usability metrics*. Morgan Kaufmann Publishers Inc.

Monitoring Vaccine Cold Chain Model with Coloured Petri Net

Fatima Ouzayd¹

Smart Systems Laboratory, Rabat IT Center,
ENSIAS, University Mohamed V,
Rabat, Morocco

Hajar Mansouri²

Team EAS, LISER,
ENSEM, University Hassan II
Casablanca, Morocco

Manal Tamir³, Raddouane Chiheb⁴

ADMIR, Rabat IT Center,
ENSIAS, University Mohamed V,
Rabat, Morocco

Zied Benhouma⁵

CEREP Research Unit,
Higher School of Commerce
Sfax-Tunisia

Abstract—To protect and prevent vaccines from excessively high or low temperatures throughout the supply chain, from manufacturing to administration, it is necessary to monitor and evaluate vaccine cold chain performance in real time. Therefore, today, the need for smart tracking is a requirement that is accentuated with critical systems, such as the vaccine supply chain. In this article, we propose a model for instant cold chain monitoring using a colored Petri net (CPN). This model focuses on the central storage of vaccines and takes into account certain WHO (World Health Organization) recommendations. The simulation and the key performance indicators obtained can be useful for decision-makers in order to measure the effectiveness and efficiency of vaccine storage.

Keywords—Vaccine cold chain; monitoring; World Health Organization (WHO); colored Petri net (CPN); performance

I. INTRODUCTION

Immunization programs have a major impact on the world's population by preventing many infectious diseases through vaccination [1]. The vaccine chain is made up of the activities and providers supplying, manufacturing, transport and administration of vaccines at the right time, to the right patient, of good quality [2]. Vaccines are biologically responsive products with their environment. In addition, any change in storage or transport conditions can have a negative impact on their usefulness and even cause undesirable effects [3]. The first parameter is the temperature at which the vaccines are exposed. This temperature should be monitored, recorded and reported throughout the supply chain, from its point of origin, the manufacturer, to its point of arrival, the place of vaccination. In order to prevent vaccines from their lifespan, they must be protected against exposure to too high or low temperatures throughout the supply chain, from manufacturing to administration. It is a shared responsibility between all actors in the chain. In our work, we will focus on the cold chain in the central store.

The aim is to provide a monitoring system to monitor the temperature of the hospital. This system detects and corrects technical issues in the cold chain components or other operational issues to maintain vaccine quality throughout the

supply chain. The final goal is to make sure the cold chain is working properly according to the recommended standards and to avoid vaccination.

In this article, we propose a colored petri nets for the description of the storage process in a central warehouse. The model incorporates the verification of the transported vaccines, the dispatching of received vaccines according to their category as well as their storage in the refrigerator taking into consideration the maintenance function and the recommended storage time in the coolant-pack. We also proposed performance indicators based on a data set.

The remainder of this paper is organized as follows: in Section 2 we present the importance of vaccine supply chain in the immunization system. In Section 3, we identify the principles processes of vaccine supply chain based on the reference model SCOR. In Section 4, we propose our Coloured Petri Net model for storage and monitoring vaccine cold chain in central warehouse, Section 5 discusses our findings and, finally, Section 6 concludes the paper.

II. IMMUNIZATION SYSTEM

The vaccination or immunization plan is one of the best ways to save the lives of millions of children worldwide in both developed and developing countries. For [4], the Immunization is unquestionably one of the most cost-effective public health interventions available because that allows long-term decreasing in illnesses and disabilities, as well as reducing health spending, etc.

In order to improve its accessibility to children worldwide, the World Health Organization (WHO) launched the Expanded Program on Immunization (EPI) in 1974 with the objective to prevent seven of the most serious diseases [5]. Through, this objective, every year, GAVI buys vaccines for more than US\$ 1 billion. In 2015, several developing countries paid the additional co-financed vaccine (more than US\$ 130 million) [6].

For this purpose, all stakeholders, WHO, UNICEF, Bill & Melinda Gates Foundation and GAVI, has given more

attention to the different constraints that condition the flow of vaccine throughout the immunization system and particularly the supply chain. Today, the improvement of the supply chain is one way to ensure that all vaccines stay safe and effective, and reach the children who need them. This requires a system to achieve the six rights of supply-chain management [2]: 6 Rights (product, quantity, condition, place, cost and temperature).

For 2020 horizon, GATES Foundations with other stakeholders set as objective to prevent more than 11 million deaths, 3.9 million disabilities, and 264 million illnesses [7]. However, all vaccines strategies continue to face delivery challenges in terms of supply chain. Unfortunately, gaps in vaccine cold chain and logistics (CCL) systems are one of the common factors limiting full and equitable access to the benefits of immunization. This is because such gaps undermine the availability and potency of vaccines at the point of administration, prevents the introduction of new life-saving vaccines, and waste precious human and financial resources [8]. The vaccines cold chain can't be efficient without four elements (a) pertinence estimation of needs system, (b) secure delivery vaccines system, (c) optimization cost and delay and (b) performance cold chain network (sustainably closing the immunization coverage gap, introducing new vaccines and securing sustainable funding) [9]. In the following, we will focus on the details of the vaccine supply chain.

III. VACCINE SUPPLY CHAIN

In this section, we first extend the processes of vaccine supply chain, and then we propose the level 2 of SCOR model that describes the processes and sub-processes of vaccine supply chain.

According to WHO [10], the vaccine supply chain contains the following processes: estimating of needs, storage, distribution, monitoring and supervision.

A. Processes of Vaccine Supply Chain

1) Estimating of Needs

The aim of this process is to ensure an adequate supply of vaccines, diluents and safe-injection equipment by assuring quality to every immunization service. Indeed, the effective management and storage of supplies must optimize costs, prevent high wastage rates and stock-outs, and improve the safety of immunizations. There are two methods that are commonly used to estimate vaccine and safe-injection equipment needs at the provincial level:

- Estimating vaccine and injection equipment needs based on the target population.
- Estimating vaccine and injection equipment needs based on previous consumption.

2) Storage

This process is about how to select and maintain cold-chain equipment, how to estimate the total volume of vaccines and safe- injection equipment to be stored and how to manage the storage of these items. About vaccines storage condition, WHO proposes that each vaccine has its own specific storage requirements, so it is extremely important to know how long

and at what temperature each vaccine can be stored. Also, for selecting appropriate cold chain equipment, WHO recommends to have information about reliable electricity and local situation. Therefore, there is another activity in this process. It's about storage capacity. That concerns fixing the total volume available to store vaccines.

3) Distribution

The aim of distribution systems for vaccines and safe-injection equipment is "to ensure continuous availability of adequate quantities of potent vaccine and safe-injection equipment" [10]. Otherwise, the distribution systems must allow having a well-functioning distribution system and clearly establishing: the supply period for each level and the corresponding quantities of vaccines and safe-injection equipment to be supplied, and the suitable route and transport needed to distribute the vaccines and safe- injection equipment.

4) Monitoring and Supervision

According to [8], in order to conserve its potency and safety, each vaccine must be strictly maintained within a specific temperature range from the manufacturer to the recipient.

In Table I, we present an example of temperatures of vaccines in central Warehouse or regional Warehouse according to [11].

TABLE I. STORAGE TEMPERATURES OF VACCINE IN CENTRAL WAREHOUSE AND REGIONAL WAREHOUSE

Vaccine	Central Warehouse	Regional Warehouse
DTC	+2°C à +8°C	+2°C à +8°C
HB	+2°C à +8°C	+2°C à +8°C
BCG	+2°C à +8°C	+2°C à +8°C
VAT	+2°C à +8°C	+2°C à +8°C
Hib	+2°C à +8°C	+2°C à +8°C
VPO	-15°C à -25°C	-15°C à -25°C
VAR	-15°C à -25°C	-15°C à -25°C
RR	-15°C à -25°C	-15°C à -25°C

Indeed, the activity of monitoring and supervision takes place monthly. The aims of monitoring vaccines and safe-injection equipment are:

- Ensuring the availability of adequate quantities and the required quality of each item;
- Ensuring appropriate use in service delivery;
- Enabling the timely detection of management problems in the implementation of immunization activities so that corrective action can be taken;
- Guiding the planning process.

IV. SCOR MODEL OF VACCINE SUPPLY CHAIN

A. Modelling Techniques of Vaccine Supply Chain

The literature represents various modelling techniques among which we quote: UML, Petri Networks, BPMN, SCOR, ABM, ARIS, mathematical programming ...

For the vaccine supply chain, many authors have used modelling methods to measure the performance of the different links of this chain. For example, [12] reviews the literature on model-based supply chain network design to identify the applicability of these models for the design of a vaccine supply chain. Author in [13] has opted for the mathematical programming to model the vaccine distribution network.

B. SCOR Model of Vaccine Supply Chain

The Supply Chain Operations Reference (SCOR) model was proposed by the Supply Chain Council in 1996. The SCOR model present an approach, processes, indicators and the best practices for evaluate and diagnose the Supply Chain. This methodology based on the client is generic, rigorous, complete and structuring [14]. The SCOR model focuses on the supply chain management function from an operational process perspective and includes customer interactions, physical transactions, and market interactions [15].

The main global processes of SCOR are: Plan, Source, Make, Deliver, and Return. For Level 2, SCOR recommend to describe core processes according the production strategy. There is the "Make-to-stock" category, the "Make-to-order" category and the "Engineer-to-order" category. Level 3 of the SCOR model specifies the best practices of each process.

In Fig. 1, we define the level 2 of SCOR model of Moroccan vaccine supply chain [11]. The upstream chain concerns the producers who supply the needs of the Moroccan partners. The central or regional warehouse replenishes the hospitals according to the national immunization plan (M1 bloc: manufacturing for storage, S1 bloc: Supply for storage, D1: Deliver for storage and DR1: Return for).

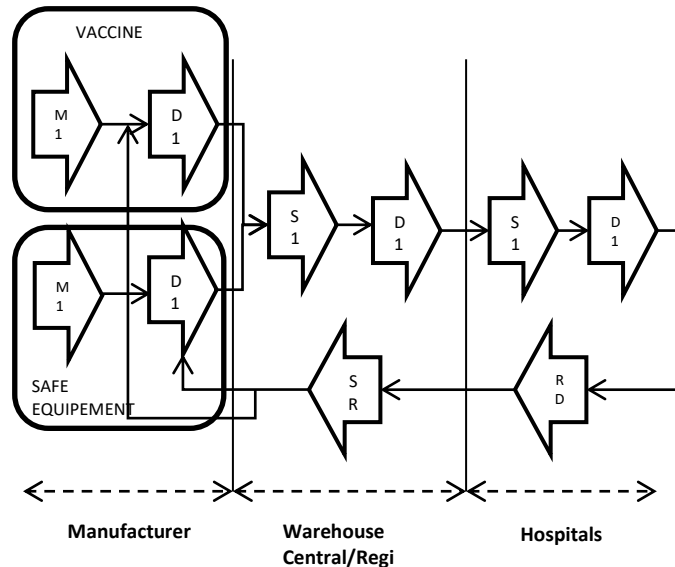


Fig. 1. SCOR Model of vaccine supply chain.

The objective of the SCOR model is to present an overview which describes the global operation of the chain in terms of: partners, global processes, critical processes which include the cold chain process. In fact, this modelling will facilitate the study of the performance of the vaccine chain through that of the cold chain.

V. PERFORMANCE OF VACCINE SUPPLY CHAIN BY CONTROLLING THE COLD CHAIN

The achievement of the goals set out in the vaccination requires a supply chain and logistics from end to end functional. Indeed, the supply chain is constrained effectively manage the increase in the introduction of new vaccines, from adapt to the needs of new delivery or delivery strategies take advantage of new technological advances in the field of cold chain equipment to improve their efficiency. This requires the establishment of a system to obtain the six supply chain management criteria: the right product, right amount, right condition, right place, right time and the right cost. The evidence indicates problems in seven areas: limitations in supply chain system design, insufficient human resources, Inefficient use of data for management, weak distribution systems, inadequate budgeting and distribution systems and deficient cold chain equipment [2].

In this sense, GAVI [16] has proposed a supply chain strategy based on 5 axes: (1) Supply chain leadership, (2) Continuous improvement & planning, (3) supply chain data for management, (4) supply chain system design and cold chain equipment. For that, we propose monitoring model that allow evaluating complaint vaccine with Coloured Petri Net. The objective of the model is to automate tracking of storage and conservation vaccine in global warehouse.

A. Tracking Model of Storage Vaccine Cold Chain

1) Coloured Petri Net and CPN Tools

Coloured Petri Net offers a modelling framework that is perfectly for distributed and concurrent processes with both synchronous and asynchronous communication [17].

In addition, they are useful in modelling both non-deterministic and stochastic processes as well. We introduce for brief presentation of Coloured Petri Net theory, developed by [18]. CPN is a tuple $(\Sigma, P, T, A, N, C, G, E, I)$ satisfying the following requirements:

- a) Σ is a finite set of non-empty types, also called colour sets.
- b) P is a finite set of places.
- c) T is a finite set of transitions.
- d) A is a finite set of arcs such that: $P \cap T = P \cap A = T \cap A = \emptyset$.
- e) N is a node function. It is defined from A into $P \times T \cup T \times P$.
- f) C is a color function. It is defined from P into Σ .
- g) G is a guard function. It is defined from T into expressions such that: $\forall t \in T: [Type(G(t)) = B \wedge Type(Var(G(t))) \subseteq S]$.

h) E is an arc expression function. It is defined from A into expressions such that: $\forall a \in A: [\text{Type}(E(a)) = C(p)\text{MS} \wedge \text{Type}(\text{Var}(E(a))) \subseteq S]$ where p is the place of N(a).

i) I is an initialization function. It is defined from P into closed expressions such that: $\forall p \in P: [\text{Type}(I(p)) = C(p)\text{MS}]$.

CPN is tool and framework that allow design, specification, validation, and verification of systems [19].

B. Storage and Monitoring Coloured Petri Net of Vaccine Cold Chain

1) Description of the system

In this work, we focus on the storage process in a central warehouse. We present an example of this system with the following elements: Two operators, two refrigerators, thermometers, etc.

The extract of declaration system with CPN Tools is shown below in Fig. 2:

```

▼ Declarations
  ▶ Standard priorities
  ▶ Standard declarations
  ▼ colset Temp= int;
  ▼ colset CLot = with DTC | HB | BCG | VAT | Hib | VPO | VAR | RR ;
  ▼ colset Stime= real;
  ▼ colset Tunit=UNIT timed;
  ▼ var tu:Tunit;
  ▼ colset LotId= string;
  ▼ colset OPE=UNIT;
  ▼ var ope:OPE;
  ▼ var t: Stime;
  ▼ var tS:Stime;
  ▼ var mu_T:STRING;
  ▼ var t_faultHT:INT;
  ▼ var t_faultT:INT;
  ▼ var t_faultHT1:INT;
  ▼ var t_faultHT2:INT;
  
```

Fig. 2. Extracts of declaration system with CPN tools.

The used colors are illustrated in Table II; different variables were defined for the system. Some of these variables are used for the monitor and performance analysis and some are used for the model behavior verification. The process of vaccine storage is illustrated in Fig. 3, 4 and 5.

The token of BatchContent place represents ID of the batch, arrival time of the batch which is taken randomly from the exponential distribution and finally the content of the batch. The variable a is involved to take on the random token, from the place Temp to the transition Reception. The output of the transition is a timed token, with four different values. The reception control transition is responsible for the separation of the batches in two groups. Place Set 1 contains only batches containing type 1 vaccines whereas place set 2 contains batches containing type 2 vaccines. The control of the

temperature of the vaccines contained in each lot, allow the reject of any vaccine whose temperature does not correspond to the norm. The compliant vaccines will be stored in the refrigerator. Due to different temperature standards for vaccines 1 and 2 we distinguish two types of refrigerator. In case the refrigerator thermometer is out of order the vaccines will be transferred to the coolant-pack. Therefore, we identify an example of place with closet and type.

TABLE II. NAME AND TYPE OF COLSET

Place Name	Colset	Type
Vaccine Name	Colset	Enumerated
BatchName	Colset	Enumerated
BatchContent	ClotVaccin	Product timed
Operator	OPE	UNIT
Set 1	ColsetVaccin2	Product timed
Set 2	ColsetVaccin2	Product timed
Compliant Set1	ColsetVaccin2	Product timed
Compliant Set2	ColsetVaccin2	Product timed
NoComplaint2	ColsetVaccin2	Product timed
Refrigerator 1	ColsetVaccin2	Product timed
Refrigerator 2	ColsetVaccin2	Product timed
Reorder	ColsetVaccin2	Product timed
Thermometer good	Unit	
Thermometer broken	Unit	
Cool-Packs1	ColsetVaccin2	Product timed
Cool-Packs2	ColsetVaccin2	Product timed

The output of the transition is a timed token, with four different values. The receive control transition is responsible for separating batches into two groups. Set 1 contains only batches containing type 1 vaccines (temperature +2°C à +8°C), while set 2 contains batches containing type 2 vaccines (temperature -15°C à -25°C) (see Table III). The temperature control of the vaccines contained in each batch allows the rejection of any vaccine.

TABLE III. STANDARD TEMPERATURES OF 2 SETS VACCINES

Vaccine	Standard Temperature	Set
DTC-HB-BCG-VAT-HiB	+2°C à +8°C	SET1
VPO-VAR-RR	-15°C à -25°C	SET2

In this model, we have 2 levels of monitoring. The first is controlling the batches after reception. This control, we allow identifying the number of the non-complaint vaccine through transportation (PLACE NO-COMPLAINT2). The second is monitoring the respect of temperature under warehouse from storage (validation, verifications and conservation) to replenishment of the regional warehouses. If the products correspond to the recommended temperatures then they are stored in refrigerators or they are rejected and classified as non-compliant vaccines (PLACE NOCOMPLAINT). We add also, in this place the case of the refrigerator (with thermometer) does not work then the new vaccine that arrives must be kept in the coolant-pack while waiting to correct the failure.

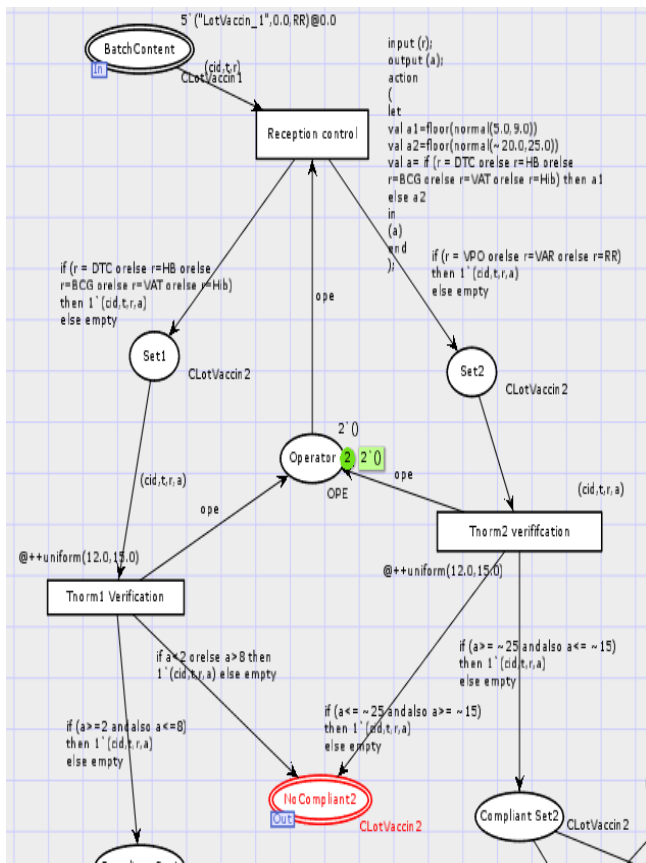


Fig. 3. Reception and control of temperature vaccine.

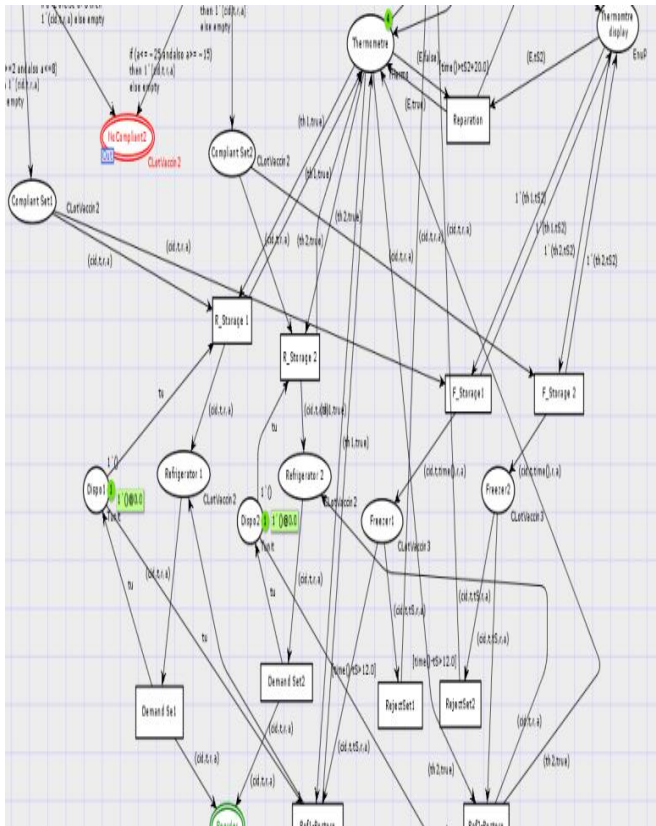


Fig. 4. Dispatching of vaccine in the refrigerators.

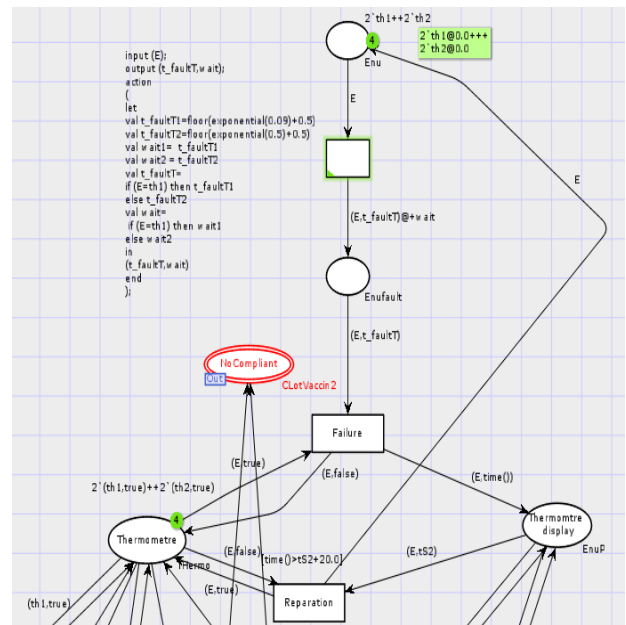


Fig. 5. Maintenance system of the refrigerators.

2) Simulation and Validation the Model

According to [20] there is a critical Cold Chain Logistics (CCL)-related issues such as delayed new vaccine introductions due to insufficient cold chain capacity, vaccine potency compromised by exposure to dangerous temperatures, and missed immunization opportunities from stock-outs. To better target such gaps and carry out corrective actions, National Immunization Programs (NIPs) need to assess CCL systems on their performance, adopting key performance indicators (KPIs) that evaluate each core CCL system task against NIP needs. Our work focuses on the performance evaluation of vaccine storage in warehouse. The aim objective of central warehouse is to:

- Reduce the risks of the break in the cold chain
- Monitor temperatures in real time
- Manage changes in product availability, inventory movements and delivery (number of insulated bags, ...)

The model that we propose allows evaluating several scenarios. We focused at first time to measure five indicators.

The first KPI represent the non-compliant vaccines due to the non-respect of the cold chain during transport. Indeed, as mentioned above, before storing vaccines in the warehouse, a temperature check is necessary. Non-conforming vaccines are rejected by stopping the process.

The “flowtime_Vaccins_Non_Compliant2” indicator represent the residence time of this type of non-compliant vaccines in the process.

The third KPI is the rate of the non-compliance vaccine due to refrigerator failures. As shown in Fig. 6, the vaccine service rate varies according to the failure rate of the refrigerator. For the failure rate between 0.1/hour for th1 and 0.9/h we get a service rate of 88% while for the failure rate between 0.001/hour for th1 and 0.009/h for th2 we get a service rate of

99%. This also has an impact on the residence time of this type of non-compliant vaccines in the process. Indeed, when the refrigerator's thermometer does not work, any new vaccine introduced into the warehouse will be stored in the coolant-pack until the thermometer is repaired. The last indicator represents the flow time of the compliant vaccines in the process.

For simplicity, we considered that the replenishment management rules are automatic.

Note that these statistics have been calculated for data that is not necessarily independent or identically distributed.

Untimed statistics					
Name	Count	Sum	Avrg	Min	Max
Percentage_transport_Non_Compliant	7210	1234	0.171151	0	1
flowtime_Vaccins_Non_compliant	2	54.415273	27.207637	25.577353	28.837920
flowtime_Vaccins_Non_compliant_2	1234	16663.172985	13.503382	12.000829	14.997730
flowtime_Vaccins_compliant	5976	81103.204572	13.571487	12.000006	36.290227
percentage_serviced	5978	5976	0.999665	0	1

Simulation steps executed: 76201
Model time: 843284.68658

Note that these statistics have been calculated for data that is not necessarily independent or identically distributed.

Untimed statistics					
Name	Count	Sum	Avrg	Min	Max
Percentage_transport_Non_Compliant	6452	1300	0.201488	0	1
flowtime_Vaccins_Non_compliant	693	19806.173226	28.580337	24.243625	45.372550
flowtime_Vaccins_Non_compliant_2	1300	17594.364584	13.534127	12.004182	14.998509
flowtime_Vaccins_compliant	5152	90664.141949	17.597854	12.000236	39.844453
percentage_serviced	5845	5152	0.881437	0	1

Simulation steps executed: 66410
Model time: 14508.3540798

Generated: Mon Mar 26 13:48:52 2018

Fig. 6. Results of rate service with 2 scenarios of failure.

VI. CONCLUSION

In this paper, we presented monitoring temperature model in vaccine cold chain with Coloured Petri Net. This automation model allows to supervising this physical parameter and to evaluate the performance cold chain in real time.

The decision maker can simulate the several scenarios to measure efficiency of system and propose some solutions. These models allow evaluating a performance of system, by fixed a means indicators: ComplaintVaccine, NoComplaintVaccine, utilization resource, Rate serveries satisfaction, etc.

We have integrated the maintenance function for evaluate the impact of reactive reparation on complaint vaccine. Indeed, we can adapt this model to any storage vaccine system in central warehouse because this model is generic monitoring temperature system with WHO recommendations.

In the next steps, we will propose a global architecture based on multi-agent systems for the implementation of the monitoring temperature model in vaccine cold chain with Coloured Petri Net.

We will also extend this model to include other global processes: supply, production, cold chain and distribution to ensure analysis and control of global vaccine chain performance.

REFERENCES

- [1] Lloyd.JS, "Improving the cold chain for vaccines," WHO Chronicle 1977 Jan;31 (1):13-8.
- [2] L. Evertje Duijzer, W.Van Jaarsveld,R.Dekker, "Literature review: The vaccine supply chain" European Journal of Operational Research 268 (2018) 174-192
- [3] J.Lloyd, J.Cheyne "The origins of the vaccine cold chain and a glimpse of the future" Vaccine 35 (2017) 2115-2120
- [4] Houari, M., (2017) Vaccins : Actualités et exigences en Contrôle qualité, Medecine Thesis, Mohammed V University, Morocco.
- [5] Ehreth J (2003) 'The global value of vaccination', Vaccine, Vol. 21, pp.596-600.
- [6] Ateudjieu, J., Kenfack, B., Nkontchou, B. W., & Demanou, M. (2013), 'Program on immunization and cold chain monitoring: The status in eight health districts in Cameroon', BMC Research Notes, Vol.6, N°1, pp.1,
- [7] Noni E. MacDonald and all, "Moving forward on strengthening and sustaining National Immunization Technical Advisory Groups (NITAGs) globally: Recommendations from the 2nd global NITAG network meeting," Vaccine 35 (2017) 6925-6930
- [8] Ashok, A., Brison, M., & LeTallec, Y. (2017), 'Improving cold chain systems: Challenges and solutions', Vaccine, Vol. 35, N°17, pp.2217-2223.
- [9] Humphreys G. (2011), Vaccination: rattling the supply chain. Bull World Health Organ, Vol. 89, N°5, pp.324-325.
- [10] Judith R.Kaufmann and all, "Vaccine supply chains need to be better funded and strengthened, or lives will be at risk" health affairs 30, no. 6 (2011): 1113-1121
- [11] Cohen, R. (2018) Guide Marocain Vaccinologie, [online], Université Cadi Ayyad, 2ème Edition, https://pharmacie.ma/uploads/pdfs/guide_marocain_de_vaccinologie.pdf (Accessed 21 March 2018).
- [12] Lemmens S., Decouttere C., Vandaele N., Bernuzzi M. (2016), "A review of integrated supply chain network design models: Key issues for vaccine supply chains", Chemical Engineering Research and Design, 109 (2016) 366-384
- [13] Sheng-I Chen, (2012), "Modeling the WHO-EPI vaccine supply chain in low and middle income countries", thesis, University of Pittsburgh, 2012.
- [14] John Paul, Jean-Jacques Laville, SUPPLY CHAIN MAGAZINE, VOL 13, pp.96 MARS 2007
- [15] Christiansen, B., (2015) 'Handbook of Research on Global Supply Chain Management', PryMarke, LLC, USA, pp.49
- [16] S.Ozawa and all, "Funding gap for immunization across 94 low and middle-income countires,"Vaccine 34 (2016) 6408-6416
- [17] Jensen, K. (1994), 'An Introduction to the Theoretical Aspects of Coloured Petri Nets', Computer Science, Vol. 803, Springer-Verlag, pp. 230-272.
- [18] Kristensen, L. M., Christensen S., Jensen. K. (1998), 'The practitioner's guide to coloured Petri nets', International Journal on Software Tools for Technology Transfer, Vol.2, pp.98-132.
- [19] Kurt J., Kristensen L.M., Wells L. (2007), "Coloured Petri Nets and CPN Tools for modelling and validation of concurrent systems", Int J Softw Tools Technol Transfer, Springer-Verlag, Vol. 9, Issue 3-4, pp 214.
- [20] Brison, M. and LeTallec Y., (2017), 'Transforming cold chain performance and management in lower-income countries', Vaccine, Vol.35, pp.2107-2109,

Framework for Rumors Detection in Social Media

Rehana Moin, *Zahoor-ur-Rehman, Khalid Mahmood
Department of Computer Science COMSATS University
Islamabad, Attock Campus, Pakistan

Mohammad Eid Alzahrani, Muhammad Qaiser Saleem
College of Computer Science and Information Technology,
Al Baha University, Al Baha, Saudi Arabia

Abstract—The development of social networks has led the public in general to find easy accessibility for communication with respect to rapid communication to each other at any time. Such services provide the quick transmission of information which is its positive side but its negative side needs to be kept in mind thereby misinformation can spread. Nowadays, in this era of digitalization, the validation of such information has become a real challenge, due to lack of information authentication method. In this paper, we design a framework for the rumors detection from the Facebook events data, which is based on inquiry comments. The proposed Inquiry Comments Detection Model (ICDM) identifies inquiry comments utilizing a rule-based approach which entails regular expressions to categorize the sentences as an inquiry into those starting with an intransitive verb (like is, am, was, will, would and so on) and also those sentences ending with a question mark. We set the threshold value to compare with the ratio of Inquiry to English comments and identify the rumors. We verified the proposed ICDM on labeled data, collected from snopes.com. Our experiments revealed that the proposed method achieved considerably well in comparison to the existing machine learning techniques. The proposed ICDM approach attained better results of 89% precision, 77% recall, and 82% *F*-measure. We are of the opinion that our experimental findings of this study will be useful for the worldwide adoption.

Keywords—Social networks; rumors; inquiry comments; question identification

I. INTRODUCTION

A rumor is an unverified claim about any event, transmitting from person to person. It may refer to an incident, object or problem of public concern. It may prove to be a social destructive phenomenon in any human culture. Usually, the social media rapidly transmits the unverified statements that may be harmful for anybody. Nowadays, social networks like Twitter and Facebook are more popular with regards to acquiring and propagating information. On social networks everybody is free to obtain and share information, anywhere at any time [1].

Besides, it has been reported that these social sites are capable to spread rumors [2]. In general, a rumor refers to the information that lacks source and its truthfulness. Ordinarily, it is generated in an emergency situation, leading to anxiety, disruption of social activities; thus, reducing the government credibility, even endangering the national security, for instance, on March 2011, after Japan Earthquake followed by tsunami and nuclear disaster. A rumor was propagated by microblog platforms, advising use of iodized salt for protection of people by nuclear radiation. Consequently, the public in general rushed to markets to buy salt, which was

totally untrue and unnecessary practice. In the future, to avoid such unfruitful happenings, at the earliest, rumor detection is essential.

Earlier, much work has been done on rumor detection using the Twitter. We did work on Facebook to address the problem of rumor detection. We selected Facebook reason being the most popular social network. In Oct 2012, Facebook was having one billion users per month. Cameron Marlow, one of the research scientists, considered Facebook as world's most powerful instrument for studying human society [3].

A framework diagram is developed for rumors detection, starting from Facebook data collection, preprocessing of data, extraction of English text, apply TopicRank to obtain keyphrases and based on those keyphrases (topics) extract the event data and detect assertion to filter assertive event posts and finally detecting the inquiry comments on assertive posts using our proposed ICDM approach. We used labeled data from snopes.com to check the validity of our proposed ICDM approach and to make comparison with machine learning techniques.

We aim to tackle the rumors detection problem using inquiry comments identification through ICDM approach. This comprises two steps. In the first phase, we identify questionable statements named as “inquiry comments”. We adopt both machine learning supervised approach like classifiers to detect questions and rule-based method to detect question marks, 5W1H words and regular expressions [4] which utilizes patterns to filter inquiries. In the second phase, we extract inquiry comments asking question about the event. We define the threshold to identify the rumors and test our ICDM model using labeled data from snopes.com.

Consequently, following research questions are formulated:

- How English text is separated from different languages?
- How to develop rumors detection framework that can correctly identify the rumors?
- How can we verify our proposed ICDM (Inquiry Comments Detection Model)?

Remaining part of this paper is organized as follows: Section II provides the related work; in Section III, methodology is presented; in Section IV, the results are presented; Section V concludes the whole work and addresses the future research directions of this study.

II. RELATED WORK

This research related work aims to explore the role of Social Media in real-world emergencies, news diffusion, and rumor detection approaches.

A. Role of Social Media in Real World Emergencies

The use of social media in emergencies and crisis has gone up many folds in recent years [5], by involving reports from the eyewitnesses. Furthermore, the use of social media in actual emergencies has also been studied. These studies have shown the importance of Social media for breaking news, information gathering and coordinating in different situations, including emergencies, protests [6], natural disasters like earthquake, floods, hurricanes and forest fires. During natural and un-natural emergencies, the social media has a significant role in both transmitting information to the situation affected people and getting live reports from eyewitnesses. The local infrastructure if intact can provide the situation information to the public in much faster rate on their mobile services, rather than using the traditional news media.

B. News Diffusion through Social Media

The social media is a medium for conversation besides a source of news for the public. Mostly the current topics on the social media are news related therefore, they can be used to detect breaking news. Breaking news reporting enables the public getting to know the current information through eyewitnesses [7]. Researchers aim to use this feature of social networks to develop the tools for latest news-gathering [8] and to report the current situation, analyzing the user-generated content (UGC) [9] and discovered the potential of social networks to give rise to citizen journalism as well as verification of reports posted in the social media platforms [10].

C. Rumor Detection Approaches

Nowadays, the discovery of social network services has led to the public to spread rumors at fastest rate. Castillo et al. [11] studied the Twitter reports of the public during 2010 Chile earthquake thus, analyzing user's behavior on microblogging platforms to reach credibility of such information and examined the retweets system to analyze the rumors propagation pattern on the Twitter. Qazvinian et al. [12] classified the rumors related tweets using the matching regular expression with the keyword query.

Zhao et al. [4] approach based on the assumption that rumors will provoke tweets from the user's inquiring about their reliability; it implies that such tweets are possibly rumor having a number of enquiring tweets. The author prepared a list of five regular expressions (such as "is (that | this | it) true") that are useful to identify the inquiry tweets.

There has been found little work on automatic rumor detection regardless of the extensive study to examine the rumors in social media and developing techniques to tackle this problem [13], [14]. A set of predefined rumors (e.g. Obama is Muslim) is fed to a classifier, which classifies new tweets as being linked to the predefined rumors or not (e.g. I think Obama is not Muslim would be about the rumor, while Obama was talking to a group of Muslims wouldn't).

Tolosi et al. [15] use feature analysis on different events datasets and found it difficult to distinguish between rumors and non-rumors as features change dramatically across events. Zubiaga et al. [16] resolved these discoveries at the tweet platform, by showing that generalizability can be attained by leveraging context of the events. In [17] studied the use of a crowdsourcing platform to detect rumors and non-rumors in social network.

III. METHODOLOGY

The framework diagram for rumors detection is presented in Fig. 1.

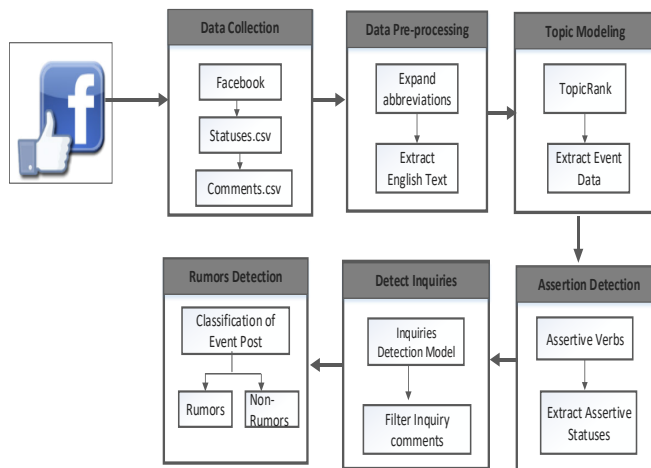


Fig. 1. Framework diagram for rumors detection.

A. Data Extraction

The Facebook data can be easily accessed and publicly available. Fig. 2 presents the flow diagram of data collection.

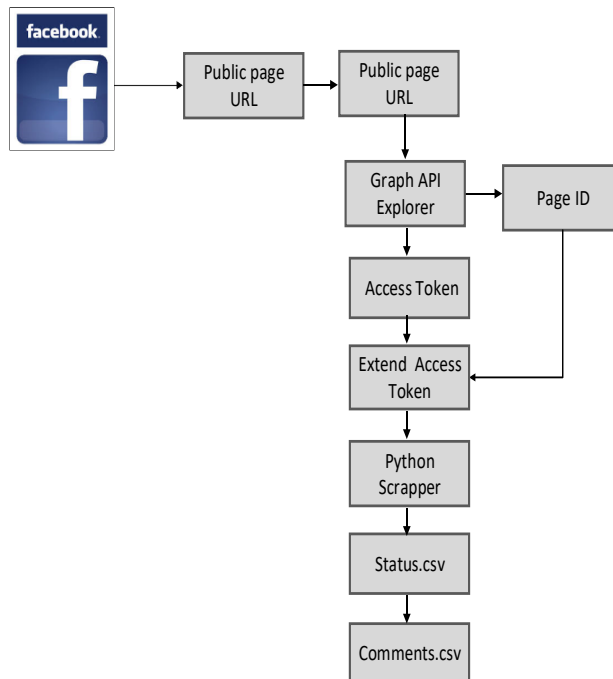


Fig. 2. Flow diagram of data collection.

Facebook data is collected using Python scraper and Facebook Graph API Explorer is used to get access token and page id of Facebook public pages. Dataset of the particular events (such as PIA flight crash PK-661 in Havelian and Pakistan Elections 2013) has been collected for the interval of one month. We have scraped all the posts along with the comments its related metadata from the Facebook news pages. The dataset is noisy and un-structured which needs pre-processing.

B. Data pre-processing

1) Data Transformation

Data pre-processing remove the un-necessary data because it degrades the system performance by making it difficult to classify the raw data. The Facebook data is pre-processed as follows:

- Expand the abbreviations.
- Extract the English text.
- Facebook content containing headlines with URLs and video links. All these links are removed as our concern is only to process the plain English text.
- Removal of such status messages having no comments

2) Expanding Abbreviations and Extraction of English Text

Public comments on the Facebook contain typos, misspelled words, unstructured and informal text. We designed an approach to tackle the typos by expanding the abbreviations to correct the short form of words and recognize the English text using dictionary-based approach. We created an abbreviation list of most commonly used abbreviations. We compared each word in a comment with the abbreviation list to extract full form of word and replaced it in the actual dataset. To separate English text, a dictionary-based approach is applied in which each word in the comment is checked in the English dictionary and extracted the English text in a separate file. For each comment, number of English words are counted and divided by the total number of words present in the comment to get its weight as presented in (1).

$$S_w = \frac{Eng_w}{total_w} * 100 \tag{1}$$

Where, S_w indicates sentence weight, Eng_w represents number of English words in a comment and $total_w$ represents total number of words in a comment. To detect English, threshold value is set as 80%. We check text if it contains 80% of English words then retain the row and if, no, delete the row and do not process it. Therefore, the comments with weight equal or above 80% are classified as English whereas those comments having weight below 80% are classified as belonging to non-English text. Repeat till file ends. After the completion of the whole process it generates file that contain only English Text. The flow of replacement of abbreviation with actual word is shown in Fig. 3.

Facebook data contains many abbreviations and short form of words having varying writing styles. We cannot apply text mining techniques directly to get better and acceptable results. The data need to be transformed into standard format before

applying any text mining techniques. The abbreviation replacement can improved the performance of the system significantly.

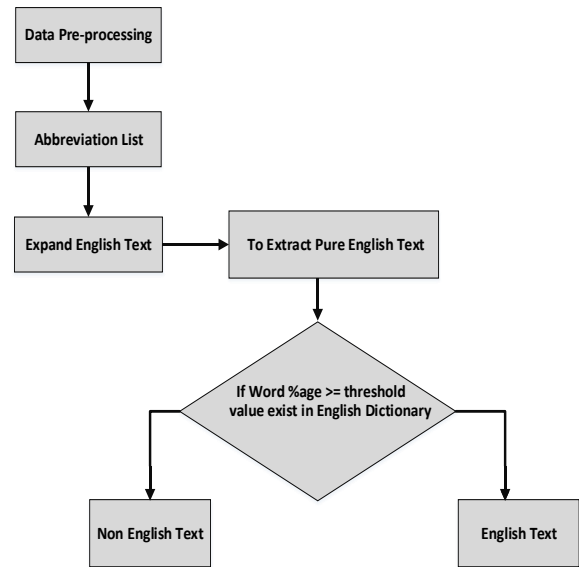


Fig. 3. Expanding abbreviations and extracting English text.

A list of most commonly used abbreviations is presented in Table I.

TABLE I. LIST OF ABBREVIATIONS WITH EXPANSION

Abbreviations	Expansion	Abbreviations	Expansion
Y	Why	Abt	About
V	Very	Plz	Please
K	Ok	thnx	Thanks
R	Are	Sry	Sorry
&	And	Sty	Stay
Ur	Your	Dept	department
w8	Wait	Ths	This
Fi9	Fine	B/W	between

C. Topic Modelling for Keyphrase Extraction

Topic Rank is an unsupervised, graph-based key phrase extraction method. In Social media, event detection is a prominent research topic. TopicRank is used to discover the topics that need to be manually analyzed during post-processing to select the best topic, describing about the events. It has become a challenge to obtain the event relevant posts, since posts may have event relevant terms but describing something other than the event such as a post containing the term “earthquake” could refer to an actual earthquake or to a conference on earthquakes. Fig. 4 presents the steps involved in keyphrases extraction from a document though TopicRank. Generally in a document, one noun phrase is sufficient to convey the topic. Therefore, some candidate keyphrases are redundant to represent the topic. Existing graph-based methods (Text Rank, Single Rank, etc.) do not take that fact into account. Candidate weighting, is assigned using a random walk algorithm. N-best selection, keyphrases contains the 10 highest scored candidates as (keyphrase, score) are extracted.

D. Assertive Detection

An assertion is a forceful statement of fact or belief. The assertion detection is used to get better understanding of intention and state of mind of the users behind event posts on Facebook. For example, a post “making a statement” can help us to track the assertions being made about events and can reveal a lot about the general attitude of users about that topic.

Automatic classification of the dialogue acts is a challenging task, traditionally a number of supervised methods are used to address these challenges; however, they require substantial manual time, effort and these depend on the availability of abundant training data.

We used an unsupervised approach, for classifying assertive dialogue acts of an event. We collected the analytical verbs from online¹ resource, a total of 76 words and we used 106 assertive verbs from the published² (Soroush Vosoughi, 2016). By using these verbs, assertive posts are extracted and afterward we scraped comments of assertive posts and detect inquiries from them.

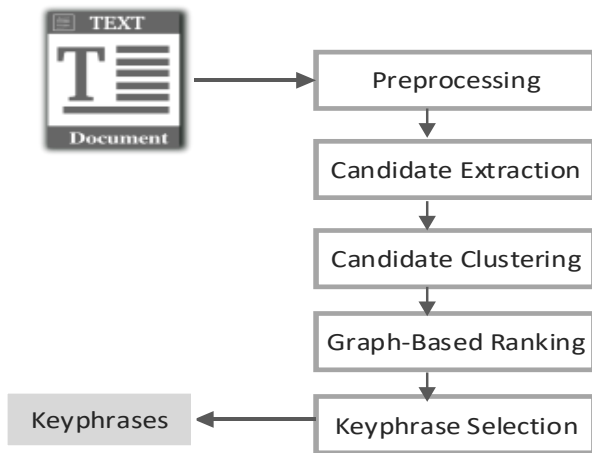


Fig. 4. Steps involved in Keyphrase Extraction through TopicRank.

E. Inquiry Comments Detection Model (ICDM)

We investigate the problem of inquiry comments detection in the textual environment; it involves detection of comments of assertive posts that contain questions (we call them “inquiry comments”). To detect inquiries from comments, a novel method is employed.

Inquiry Comments Detection Model has been implemented using Python. Algorithm employs two modules to classify sentences as inquiry comments which are useful in finding fact about some certain event propagating on the social media.

First module of our algorithm identifies a sentence as a question that satisfies following conditions:

- It ends with a question mark “?” and particularly such sentence does not contain any URL.
- It starts with an intransitive verb (like is, am, was, were, will, would, could, has, have, etc.)

- It starts with 1H-5W question words (How, What, Who, When, Why or Where).

Its second word is also checked, if it is not an intransitive verb, example: "where there is a will, there is a way." is not a question. But if it is “where are you going”, the word “where” is followed by an intransitive verb, so, it is a question.

Second module of our algorithm is based on regular expression given in Table II.

TABLE II. PATTERNS USED TO FILTER INQUIRIES

Regular Expression	Type
wh[a]*t[?!][?1]*	Question
is (that this it) true	Question
Real ?? really?? unconfirmed	Question

F. Rumors Detection and Verification System

After calculating the ratio between English and Inquiry comments of a post, the ratio is compared with threshold to determine whether it is a rumor or not.

The threshold of a post (T_p) is:

$$\text{Threshold } (T_p) = \begin{cases} \text{rumor} & \text{if } T_p \geq 0.1 \\ \text{non - rumor} & \text{if } T_p < 0.1 \end{cases}$$

If $T_p \geq 0.1$ post is considered as a rumor.

If $T_p < 0.1$ post is non-rumor.

For reliable rumors detection and testing our proposed model, we collect a set of labeled and verified rumors/non-rumors dataset from Facebook snopes.com page. Snopes.com is a popular resource for debunking and validating rumors. It cites and aggregate trustworthy external sources (news or governmental organization) for verification of rumors, urban legends, documents Internet rumors and other questionable statements.

IV. RESULTS AND EVALUATION

A. Experimental Setup

The rule-based model has been designed for inquiry comments detection. We performed various experiments to measure the effectiveness of ICDM with traditionally designed machine learning models and to test it using labeled data.

In our experiments:

- **TP (True Positive)** are inquiry samples that are accurately classified inquiry samples.
- **TN (True Negative)** are non-inquiry samples that are accurately classified non-inquiry samples.
- **FP (False Positive)** samples are inquiry samples that are misclassified non-inquiry samples.
- **FN (False Negative)** samples are as non-inquiry

¹ http://msweinfurter.weebly.com/uploads/5/4/3/7/5437316/analytical_verbs.pdf
² http://soroush.mit.edu/publications/vosoughi_roy_speechact_icwsm2016.pdf

samples that are misclassified samples inquiry.

Aiming to measure the efficiency of our proposed approach, we have used well-known metrics: precision, recall, and *F*-measure. Precision represents the ratio of predicted positive samples that are real positives and is calculated by (2). On the other hand, recall is the ratio of true positive samples that were correctly predicted as such. Recall is calculated by equation 3. *F*-measure is the harmonic mean of precision and recall. *F*-measure is calculated by (4). Accuracy is the number of true positive and true negative samples out of total number of samples. We have calculated accuracy by using (5).

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

$$F - \text{measure} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

$$\text{Accuracy} = \frac{(TP + TN)}{TP + TN + FP + FN} \quad (5)$$

B. Results

In result section, we have described the accuracy of our proposed framework of inquiry comments detection model (ICDM) for rumors detection and verification using labeled data from snopes.com.

1) Rumor Verification System

We retrieved from the Snopes.com two classifications of the rumors they have analyzed, the first is the reliability, which includes “true” and “false”, but also a range of intermediate instances i.e. partly true, mixture, unproven, miscaptioned and legend. A rumor was annotated as true or false if trustworthy source snopes.com confirmed it as such. The rumor and non-rumor dataset contains 500 posts having more than 5,000 comments. The target variable labels each post to be true or false is represented by rumor or non-rumors.

TABLE III. RESULTS OF VERIFICATION SYSTEM USING ICDM AND MACHINE LEARNING METHODS

Techniques	Accuracy	Precision	Recall	F-measure
ICDM	0.70	0.89	0.77	0.82
k-NN	0.26	0.21	0.88	0.33
Naive Bayes	0.21	0.21	1.00	0.35

Table III shows that using labeled data, our approach obtained encouraging results with a precision of 0.89%, a recall of 0.77%, and *F*-measure of 0.82%.

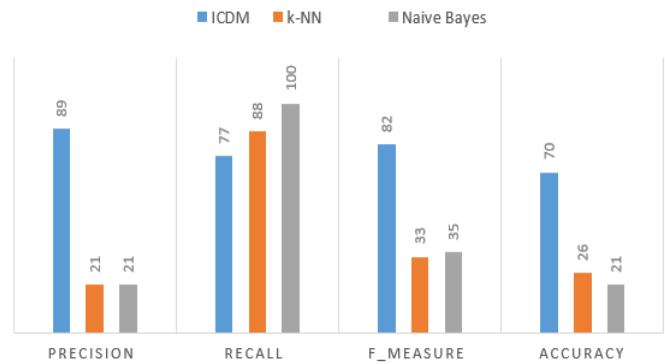


Fig. 5. Comparison between ICDM and machine learning methods.

Fig. 5 shows that on labeled data of 500 posts, our proposed ICDM approach obtained better results with a precision of 89%, a recall of 77%, and *F*-measure of 82% as compared to Machine learning approaches.

V. CONCLUSION

The current work explored the problem of rumors detection based on inquiry comments identification using textual content of social media especially Facebook. Data of Facebook is secure and inaccessible to access except the public pages. Therefore, we have searched for Facebook news pages and collected some event’s data. The scraped data is unlabeled, huge in volume and mix of multiple languages. An expanded list of abbreviation is prepared to remove inconsistencies during the pre-processing phase. For detection of event relevant post, we used topic modeling technique such as TopicRank to discover and select best topic describing the event. Once the system is able to mark topics, assertive posts about that event are extracted and relevant comments were scraped. A rule-based approach is developed to extract inquiry comments from the assertive post comments. We verified our ICDM approach using labeled data from snopes.com and achieved better results as compared to machine learning-based approaches, with 89% precision, 77% recall, and 82% *F*-measure. We believed that the experimental findings from this study will be useful in real-world inquiry classification problems.

In the future, our focus will be to address multilingual content to avoid removal of other language content that reduces the data resulting in the loss of vital information to get better understanding of people judgment about the event on the social network and containing short and informal questions where 5W1H words or question mark are likely to be absent. Our focus will be to handle informal online languages to explore question characteristics and devise an automated method to detect interrogative sentences based on syntactic and lexical features.

REFERENCES

- [1] A. Friggeri, L. Adamic, D. Eckles, and J. Cheng, “Rumor Cascades,” *Icwsn*, pp. 101–110, 2014.
- [2] J. Kostka and R. Wattenhofer, “Word of Mouth : Rumor Dissemination in Social Networks,” pp. 1–14.
- [3] G. Goggles, “WHAT emtech MIT,” no. August, 2012.
- [4] Z. Zhao, P. Resnick, and Q. Mei, “Enquiring Minds: Early Detection of Rumors in Social Media from Enquiry Posts,” *WWW ’15 Proc. 24th Int. Conf. World Wide Web*, pp. 1395–1405, 2015.

- [5] S. E. Middleton, L. Middleton, and S. Modafferi, "Real-time crisis mapping of natural disasters using social media," *IEEE Intell. Syst.*, vol. 29, no. 2, pp. 9–17, 2014.
- [6] S. D. Agarwal, W. L. Bennett, C. N. Johnson, and S. Walker, "A model of crowd-enabled organization: Theory and methods for understanding the role of twitter in the occupy protests," *Int. J. Commun.*, vol. 8, no. 1, pp. 646–672, 2014.
- [7] N. Diakopoulos, M. De Choudhury, and M. Naaman, "Finding and assessing social media information sources in the context of journalism," *Proc. 2012 ACM Annu. Conf. Hum. Factors Comput. Syst. - CHI '12*, p. 2451, 2012.
- [8] A. Zubiaga, H. Ji, and K. Knight, "Curating and contextualizing Twitter stories to assist with social newsgathering," *Proc. 2013 Int. Conf. Intell. user interfaces - IUI '13*, p. 213, 2013.
- [9] P. Tolmie et al., "Supporting the Use of User Generated Content in Journalistic Practice," *Proc. 2017 CHI Conf. Hum. Factors Comput. Syst. - CHI '17*, pp. 3632–3644, 2017.
- [10] J. Spangenberg, "News from the Crowd : Grassroots and Collaborative 3 . GRASSROOTS JOURNALISM - A NEW," *Www*, pp. 765–768, 2014.
- [11] M. Mendoza, B. Poblete, and C. Castillo, "Twitter under crisis," *Proc. First Work. Soc. Media Anal. - SOMA '10*, pp. 71–79, 2010.
- [12] E. Rosengren, D. R. Radev, Q. Mei, and A. Arbor, "Rumor has it : Identifying Misinformation in Microblogs," pp. 1589–1599, 2011.
- [13] S. Hamidian and M. T. Diab, "Rumor Identification and Belief Investigation on Twitter," *Acl*, pp. 3–8, 2016.
- [14] S. Hamidian and M. Diab, "Rumor Detection and Classification for Twitter Data," no. c, pp. 71–77, 2015.
- [15] L. Tolos, A. Tagarev, and G. Georgiev, "An Analysis of Event-Agnostic Features for Rumour Classification in Twitter," pp. 151–158, 2012.
- [16] A. Zubiaga, M. Liakata, and R. Procter, "Learning Reporting Dynamics during Breaking News for Rumour Detection in Social Media," 2016.
- [17] R. McCreddie, C. Macdonald, and I. Ounis, "Crowdsourced Rumour Identification During Emergencies," *Proc. 24th Int. Conf. World Wide Web - WWW '15 Companion*, pp. 965–970, 2015.

Modeling of Arduino-based Prepaid Energy Meter using GSM Technology

Uzair Ahmed Rajput¹
Department of Electrical Engineering
QUEST Nawabshah Sindh Pakistan

Khalid Rafique²
Director General Information Technology Board
Azad Jammu & Kashmir

Abdul Sattar Saand³, Mujtaba Shaikh⁴, Muhammad Tarique⁵
Department of Electrical Engineering
QUEST Nawabshah, Sindh, Pakistan

Abstract—It is realized that one of the defective subsystems adding to the tremendous budgetary loss in Power Supply Company is the conventional metering and charging framework. Mistakes get presented at each phase of charging the energy rates, similar to blunders with conventional meters, reading errors by human while noticing the consumed energy; and blunder during the preparation of paid and the due bills. The solution for this downside is a prepaid charging or billing framework of consumed energy. Most of the developing countries are shifting their conventional energy management practices to the modern one by replacing the old and conventional energy meters with the smart meters outfitted with the prepaid facility to quantify the power consumption so as to decrease the income deficits looked by utilities because of customer unwillingness to make consumed energy payments on time. Our proposed design embedded with Arduino and GSM technology is advancement over conventional energy meter, which enables consumer to effectively manage their electricity usage. The system performance is good with the acquired results. An earlier charging will undoubtedly get rid of the issues of unpaid bills and human mistakes in meter readings, along these lines guaranteeing justified income for the utility.

Keywords—Arduino; energy meter; smart meters; RFID; GSM

I. INTRODUCTION

Electric energy consumed by the power devices is measured by a gadget known as an energy meter. Since 1980s, the energy meter's journey started. When there were large energy meters which have been made with lots of innovations in energy meters to reduce its size and also the weight. Particularly, enhancement in accuracy, specifications and features of energy meters has been a topic of discussion.

In Pakistan, electromechanical energy meters were used for a long time. These meters work by checking and figuring the quantity of turns of an electrically directing metal plate which is made to rotate at a speed in respect to the power experiencing the meter. Those electromechanical energy meters are being supplanted by the newly digitized meters due to different problems like there is no way to upgrade those energy meters, its accuracy was limited and those meters were easy to manipulate because direction of revolving disc can be easily reversed. Nowadays, digital energy meters can measure

voltage, current and power also but electromechanical energy meters can only measure active power. Digital meters measure energy usage by highly integrated circuits, by capitalizing the voltage and current that gives the instantaneous power in watts. Digital meters show usage of electricity in digits on a liquid crystal display and those meters are highly accurate, inexpensive, theft reluctant, etc.

This work is intended to gather the information about the data which is consumed energy of a specific user or consumer through a wireless communication system (not required to visit consumer premises), and the system is called as AMR (Automatic Meter Reading). The AMR system is proposed to remotely accumulate the meter readings of a locale using a relating remote wireless system without individuals physically going to and taking note of the readings of the meters [1].

II. PROBLEM STATEMENT

As we mentioned above several advantages of digital energy meters, but always there are chances of innovation or modification in different instruments for ease of consumer and supplier. Following are some problems observed in those energy meters which should be rectified:

- Meter reading and other related tasks like bill payment are performed by a large number of staff i.e., large number of employees are required.
- An expansive number of staff is utilized for meter reading and other related assignments like bill payment.
- Billing errors due to carelessness of meter readers during meter reading and sometime billing estimation.
- Consumer has to stand in queue for hours for bill payment.
- Careless usage of electricity by consumer who is unaware of its cost.
- Consumers are not bound to pay bill on time.

III. HYPOTHESIS

As a solution of above mentioned problems, "prepaid energy meters" are being introduced. As per a current report

from Navigant Research, the overall introduced base of prepaid meters is required to add up to more than 85 million from 2014 to 2024 [2].

IV. EXISTING PREPAID SYSTEM

- Smart Card Based Prepaid Energy Meters

In this type of prepaid energy meter, there are two main components; one is smart card and another one is smart card reader. Smart card is like credit cards made of plastic and it consists of different components like CPU, ROM, EEPROM, etc. so basically integrated circuit is embedded on a smart card [3]. There is a whole smart card operating system through which data is controlled of a smart card. In this kind of scheme, the consumer must recharge his card as much number of units he wants. Later, that card is inserted into card reader which is embedded with energy meter like a whole package. Afterwards, card reader does its work and stores the units which are available in smart card, energy meter reduces the units as much electricity is being consumed. When unit reaches to zero, it disconnects the electricity until recharge.

- Smart meters (prepaid meters) with GSM technology.
- In this sort of scheme, a message is sent to smart meters via GSM network by the consumers after recharging their cell phone account. As much amount (rupees) is sent to energy meter, it purchases number of units and those units are stored in energy meter. As consumer use electricity units are reduced by energy meter and when the purchased units are exhausted electricity is cut-off. When next recharge is sent over to energy meter, it recognizes the mobile number and decodes message, add number of units in its storage so as electricity is restored [4].

V. RFID BASED SMART METER

Radio-frequency identification (RFID) is a programmed recognizable proof strategy, depending on putting away and remotely recovering information utilizing gadgets called RFID tags or also known as transponders. The development requires some level of the coordinated effort of a RFID reader and the tag. A RFID tag is associated with or melded into a thing, animal, or a person with the true objective of recognizable proof and following using radio waves. A few tags can be read from a few meters away and past the line of sight (LOS) of the reader [5].

Following this strategy, the individuals utilize the RFID cards issued by the power providers. The energy could be purchased by recharging the novel RFID cards while utilizing the code in the card. At the point, when the buyer needs to utilize the electricity, he needs to demonstrate the card to the reader, at that point the one of a kind code inside the card is perceived by the reader, and begins deducting the RFID card amount according to the quantized unit charge. After utilization of whole amount, the consumer needs to recharge the RFID card again [5].

VI. PROPOSED SYSTEM

A. Benefits of the Proposed System

There are many benefits for customer as well as for supplier of this project:

- Pay according to your current income situation.
- Reduce electricity consuming when income is tight.
- No billing errors.
- No debt money.
- No need of extra staff for meter reading.
- Customer will be responsible for disconnections.

B. System architecture (B.D & B.D Description)

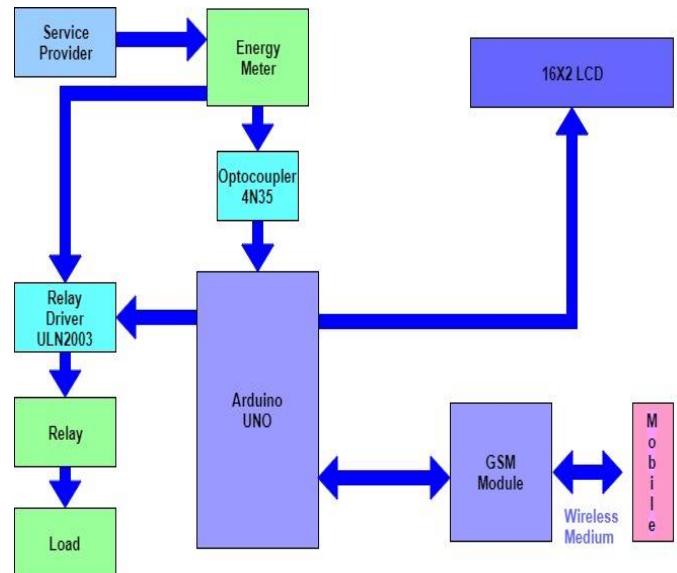


Fig. 1. Layout of prepaid energy meter.

This layout of Fig. 1 yields the idea behind prepaid energy meter scheme. It shows different hardware components incorporated, their connections and the data flowing in this scheme. In this layout, the energy meter is supplied by the service provider (WAPDA) and supplying the energy to the load through an electromechanical relay, which operates under certain conditions. Energy meter is also connected to the Arduino UNO in order to obtain the objective of automatic measuring the consumed energy and to perform the suitable action depends upon the current units available. Arduino UNO is additionally associated with the GSM Module for sending or receiving the SMS to the client for the status of their connection and to recharge the energy meter if it is required to do so. LCD is connected to the Arduino UNO board to display the current status of the connection so the customer remains in touch with the information about the load and their purchased units [6].

D. Methodology

The entire Circuit diagram has been shown below in Fig. 2. The power is measured by the energy meter with respect to time and is calculated by multiplication of voltage and current signals. The IC of energy meter generates pulses according to real power utilization. This energy meter calculates 1KWh for 3200 impulses, so rated as 3200imp/KWh, and there will be blinking of an LED for its every pulse. An Optocoupler has been connected to this LED so Optocoupler will be switched whenever LED blinks. We cannot directly connect energy meter's LED with Arduino because LED possesses analogue signals while we are feeding Arduino on the digital side. The pin number(D8) of Arduino is attached to the switching side of an Optocoupler for detecting pulses coming from energy meter. When a pulse occurs from energy meter, optocoupler is switched, pin D8 of Arduino detects a digital 0, otherwise it is not active and is in undefined state. There will be a count 1 to a data when there will be change on the state of the pin from digital 1 to 0. We have interfaced GSM module with Arduino UNO. The data communication pins are RX and TX, Arduino's RX pin is connected with GSM module's TX pin and vice-versa. Before connecting GSM module with Arduino, a valid SIM card must be installed in SIM card port of GSM

module. All ground pins GND are connected together. For switching purpose (ON/OFF) to supply a relay is being used. We cannot connect Arduino directly with relay because as Arduino has ATMEGA328P processor and its pins can supply roughly 25mA, Processor pins have large effective resistance and a high voltage will "drop" as increasing current is drawn and a low voltage will rise as load increases. Pins may be specific with a maximum short circuit current but at that point a high pin will be pulled low and a low pin will be pulled high so short circuit current has limited applicability. So, relay is connected with Arduino through ULN2003 IC or relay driver, ON/OFF instructions are sent over to relay driver by Arduino and it can turn ON/OFF relay. LCD is also interfaced with Arduino digital pins (7, 6, 5, 4, 3, 2) on which we can see how much units are purchased, remaining units and balance, etc. [7]. Fig. 3 shows the flow diagram of processes involved in prepaid energy meter scheme.

E. Circuit Diagram

See Fig. 2 below.

F. Work flow

See Fig. 3 below.

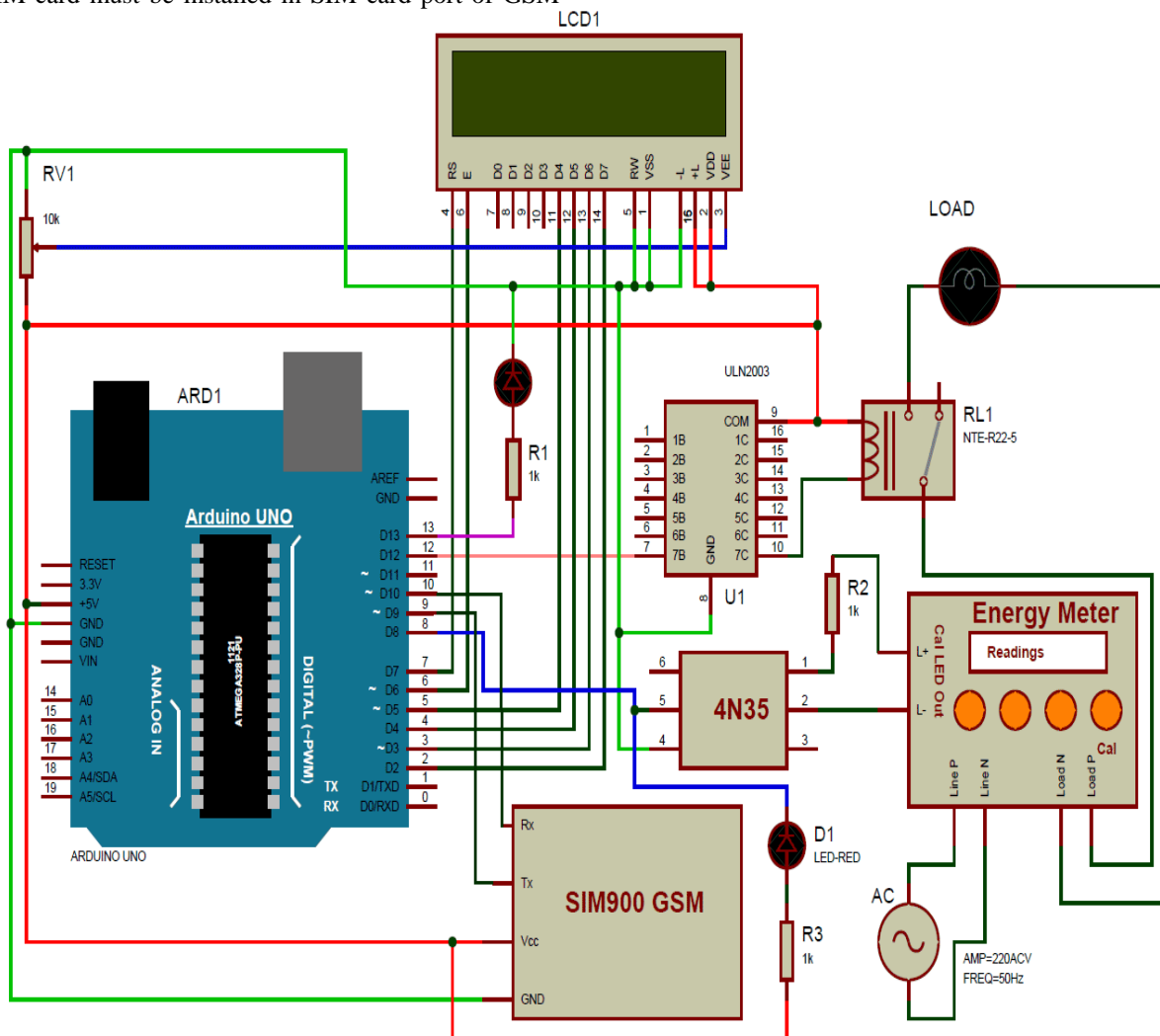


Fig. 2. Circuit diagram of prepaid energy meter.

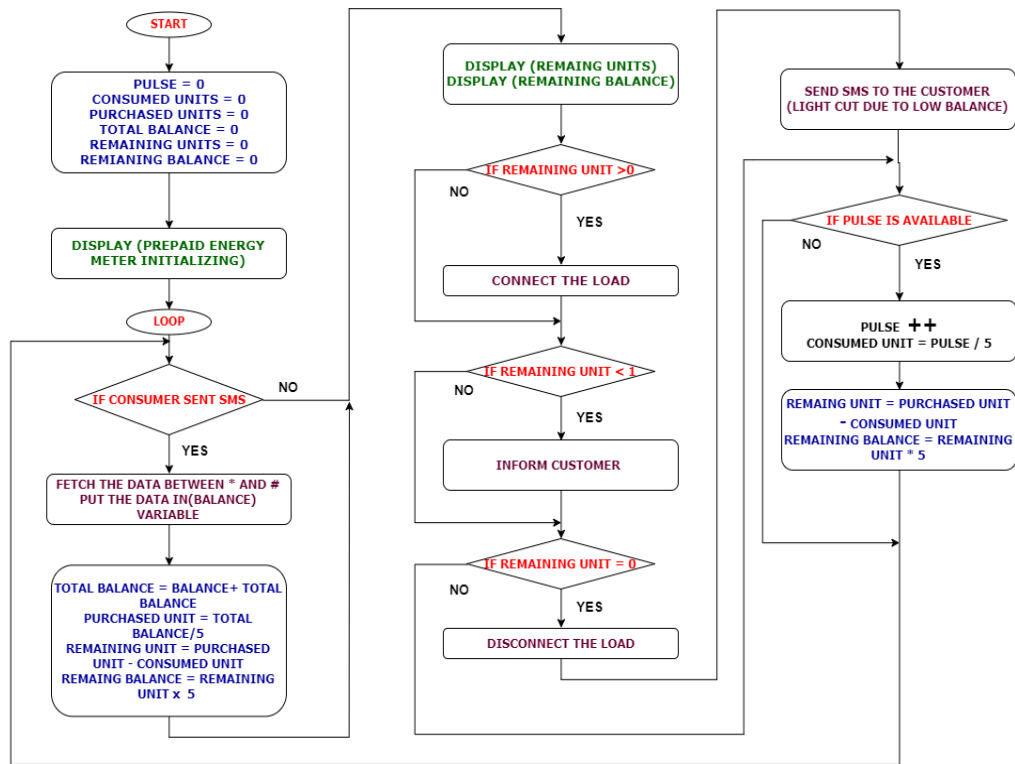


Fig. 3. Flow diagram of processes involved in prepaid energy meter scheme



Fig. 4. Display zero units and zero balance.

When there is no balance in energy meter as it would be at initial state or it can be when all the purchased units are consumed, the microcontroller will display zero units and zero balance as shown in Fig. 4. At the same time, it will inform the customer regarding no balance by sending an SMS through GSM module as shown in Fig. 5.

- When customer recharge some amount to the energy meter (Fig. 6)

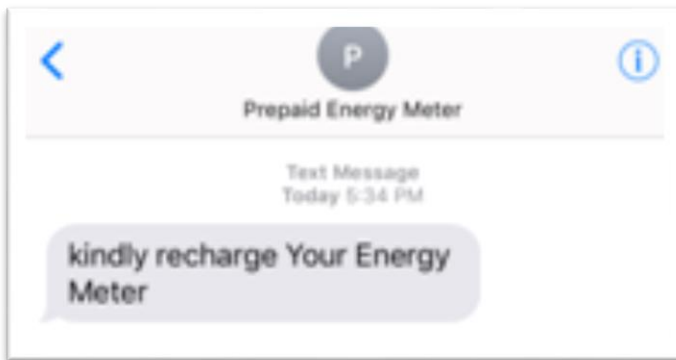


Fig. 5. No balance by sending an SMS through GSM module.



Fig. 6. When customer recharge some amount to the Energy Meter.



Fig. 7. Show credit.

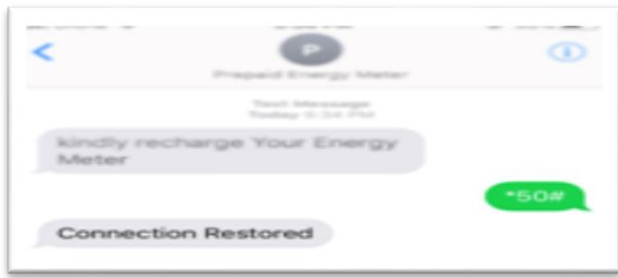


Fig. 8. An SMS is sent back to the customer.

When the customer recharges a certain amount by sending a message to the GSM module, it delivers that specific data to the microcontroller (i.e. Arduino UNO) so it can decode it and fetch the amount that customer wants to recharge his account. The recharged amount is displayed on LCD which is shown in Fig. 7.

After a certain process the microcontroller commands to connect the load to the supply as the balance in customer's account is sufficient to get the connection back, for this regard an SMS is sent back to the customer to inform him that their connection is restored as shown in Fig. 8.

- When energy meter cut off the load due to insufficient balance (Fig. 9)



Fig. 9. When energy meter cut off the load due to insufficient balance.

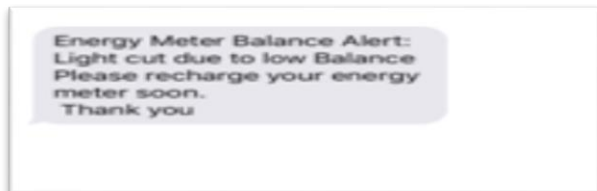


Fig. 10. SMS to inform the customer regarding the status of their connection.

The customer is bound to use as much amount of electrical energy as the balance he has in his account. As the balance in his account reaches zero value then it becomes insufficient to continue the electric supply to the load; that is why the microcontroller commands to disconnect the load from supply. A message about this disconnection of load is displayed on LCD and also sent by SMS to inform the customer regarding

the status of their connection as shown in Fig. 9 and 10, respectively.

VII. CONCLUSION

The advancement in power distribution system is non-stop process and new technology is always in progress. In this paper, an Arduino and a GSM based smart prepaid energy meter has been proposed. Units are purchased by using GSM technology and those units are deducted according to electricity usage. This project presents a single-phase energy meter for domestic consumers with prepayment billing method. The significant preferred standpoint is the capacity of this system to update the current conventional meters into smart prepaid meters with a connection of Arduino and GSM (Prepaid Module). This kills the need of totally supplant the energy meters. Cost is the main important factor of this work which is quite high but will reduce from 3 to 4 times after implementation of this project. Nowadays as power supply companies need labour for meter reading after implementing this, there will be no need of so many meter readers and lots of money will be saved.

The idea of prepayment electricity bill prior its usage is being gradually accepted around the world, and that's why the market for prepaid energy metering is growing. After having many advantages, this project still needs more safety check and modification especially the GSM module for the network coverage of SIM which is being used, should be strong so that the GSM can work properly.

REFERENCES

- [1] Dike, Damian O., et al. "Minimizing household electricity theft in Nigeria using GSM based prepaid meter." *American Journal of Engineering Research (AJER)* 1 (2015): 2320-0936.
- [2] <https://www.navigantresearch.com/newsroom/prepaid-electric-meters-are-expected-to-have-an-installed-base-of-more-than-85-million-from-2014-to-2024>
- [3] M. W. Raad, T. Sheltami2 and M. Sallout, "A SMART CARD BASED PREPAID ELECTRICITY SYSTEM," in *Pervasive Computing and Applications, 2007. ICPCA 2007. 2nd International Conference, Birmingham, 2007*
- [4] Omijeh, B. O., and G. I. Ighalo. "Modeling of gsm-based energy recharge scheme for prepaid meter." *IOSR Journal of Electrical and Electronics Engineering* 4.1 (2013): 46-53.
- [5] R. Teymourzadeh, M. I. S and A. J. A. Abueida, "RFID-BASED Prepaid Power Meter," in *2013 IEEE Student Conference on Research and Development (SCoReD, Putrajaya, 2013.*
- [6] Sheelasobanarani, Dr K., et al. "A Prepaid Energy Meter for Efficient Power Management." *International Journal of Emerging Technology and Advanced Engineering* 4.3 (2014): 593-595.
- [7] Rodrigo, W. D. A. S., et al. "A prepaid energy meter using GPRS/GSM Technology for improved metering and billing." *International Journal of Computer Science and Electronics Engineering (IJCSSEE) Volume 4* (2016)

Koch Island Fractal Patch Antenna (KIFPA) for Wideband Applications

Meryem HADJI, Sidi Mohammed MERIAH, Djamila ZIANI

Laboratory of Telecom of Tlemcen LTT,
Department of Telecom,
University of Tlemcen, ALGERIA

Abstract—In this paper, a new modified printed Koch Island Fractal Patch Antenna (KIFPA) is studied. The conception of such antenna is based on the combination of different techniques. The first, concerns the fractal geometry of the patch, while the second comprises modified ground-plane. The patch is etched according to Koch Island geometry with different iteration number ($n = 1, 2$ and 3) as inductive loading. It is proximity fed by a 50Ω micro strip line. The proposed antenna operates in the frequency band [6.03–12.62 GHz] with 70.7% for $S_{11} \leq -10$ dB. The antenna gain and radiation patterns within the operating band are simulated. The design was performed using the CST Microwave Studio Software and the results are presented, compared and discussed. Finally, the proposed antenna is fabricated and the reflection coefficient parameter is measured to validate simulation results.

Keywords—Fractal antenna; Koch Island fractal-shape; microstrip patch antenna; wideband antenna

I. INTRODUCTION

Nowadays, with the emergence of wireless communication technology and growing interest on their applications, several fractals have been widely deployed in antenna designs, due to the significant improvements added to their characteristics performance.

The theory of “Fractal” was first defined by Benoit Mandelbrot in 1975 [1], which was derived from the latin word “fractus”, signifying “broken” or “fractured”. The idea behind it was to describe nature’s geometry and classify complex geometries that were generated with an iterative procedure [2], whose dimensions were not whole numbers.

The concept of fractal geometry is very requested in the field of antenna design. For instance, the space filling properties of Giuseppe Peano [3] and Minkowski [4] fractals have been exploited for the miniaturization. While, the self-similarity property of Koch [5], Sierpinski [6], Minkowski [7] and circular fractal shape [8] can be used to resonate the antenna at a number of frequency bands to obtain multi-band behavior. Fractals proprieties are also investigated in [9] and [10 -11] to obtain various kinds of wide band and ultra-wide band antennas respectively.

In this paper, we propose a new wideband fractal patch antenna for X band applications.

A. Related Work

The obvious development of wireless communication applications requires a growing need for wide band and low profile antennas. In this context many examples clearly illustrate the importance of this kind of antenna.

In the literature, several methods are suggested to improve the antenna bandwidth such as, a printed Γ -shape Fractal Antenna [12]; A Spidron and Giuseppe Peano fractal slot antenna studied in [13] and [14], respectively; a combination of two fractal geometries (Koch-Minkowski and Koch-Koch) along with the slot of a rectangular printed patch antenna with partial ground plane [15]; a coaxial feeding technique of modified printed square antenna [16]; a printed U-shape antenna on a circular ground plane with an inverted U-shape slot [17]; a defected ground plane using CPW feeding technique with modified and octagonal fractal patch proposed in [18] and [19], respectively.

Consequently, during our research, we are going to focus, explore, and investigate the wide band fractal patch antenna with modified ground plane.

B. Contributions

In this study, a new type of printed fractal modify Koch Island patch antenna is proposed. The printed patch of the final designed structure is obtained from the third iteration of the half Koch fractal. Also, the ground plane is modified by inserting small steps. So the final proposed antenna is characterized by:

- Wide impedance bandwidth with the same direction of radiation which is suitable for wideband applications.
- Small size, lower profile and easy to fabricate.

A comprehensive parametric study has been carried out to understand the effects of various Fractal Iteration Number parameters in order to achieve the best performance possible of the final antenna with modified ground-plane. Good impedance bandwidth covering 6.03–12.62 GHz frequency band (determined from -10 dB return loss) used for X-band applications is achieved. The simulations results are performed using CST Microwave Studio software [20]. Thereafter experimental results are obtained for the third iteration order antenna for comparison.

II. DESCRIPTION OF THE PROPOSED ANTENNA

Fractal shaped antennas present very important characteristics which are related to the geometrical properties. The Koch fractal appeared in 1904 by the Swedish mathematician Helge von Koch [21]. It is a simple example of a fractal structures used to reduce the antenna dimensions and enhance bandwidth performance, due to the increase of the effective electronic length.

It is characterized by two important parameters that are the Iteration Factor and the Iteration number.

The Iteration Factor specifies the process rule of fractal structure generation, in this case $IF=1/3$ whereas the Iteration number describes how iterative processes are performed. When the number of iteration is zero, we have an equilateral triangle called the fractal generator as shown in Fig. 1. But if the iteration number converges to infinity meaning the process is executed an infinity of times, the Koch snowflake fractal is achieved [22].

The Koch Island is used in this antenna owing to the self-similarity property, the structure conserves the same form for any step as shown in Fig. 2(a). R is the radius of a circle which covers the entire fractal shape. So a large surface area is bounded by a circle of fixed radius R . Fig. 2(a) present the Koch Island fractal generated by equilateral triangle at the third iteration. More, the circular patch presents a very interesting characteristic for wide band applications compared to other structures.

Whereas, in our work, the antenna structure consists only of the half third iterative Koch Island patch generated from an equilateral triangle of side $a = 11.5\text{mm}$ which is very similar to that of a semi-circle patch studied in [23] and modified partial ground plane.

The construction of the designed final antenna structure is illustrated in Fig. 2(b). The antenna is printed on a 1.6 mm-thickness FR4 dielectric substrate of dimensions $30 \times 30 \text{mm}^2$ and relative permittivity 4.4. Where, a patch is placed co-planar with finite modified rectangular ground plane with size of $GL=10\text{mm}$ and $GW=30\text{mm}$. The separation distance between the patch and the ground is $d = 2 \text{mm}$. A 50Ω SMA connector is used for feeding the Koch patch antenna through the microstrip line section which is printed on the opposite side.

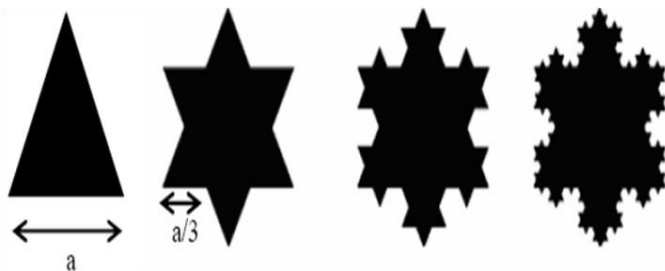


Fig. 1. Koch Island generation processing.

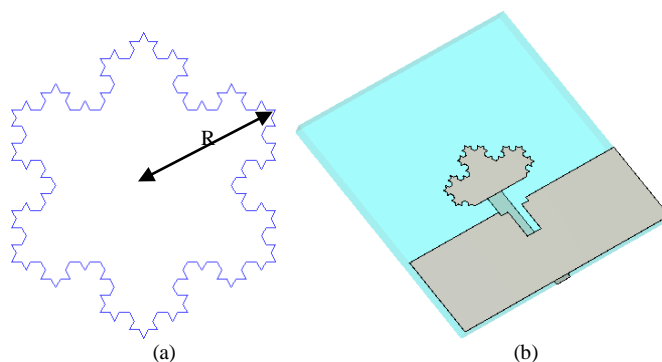


Fig. 2. (a) Koch fractal geometries; (b) Design of the proposed final antenna.

III. SIMULATION AND MEASUREMENT RESULTS

Initially, several different shapes for the patch antenna with simple ground plane were used. But in order to improve bandwidth characteristics; it was found that a patch and an optimized geometry of the whole structure gives the best possible results.

A. Effects of Iteration Number (n)

In order to see the effect of the variation of the iteration number n on the band width performance, the structures with Iteration Factor $IF=1/3$ and different $n=1, 2$ and 3 for antenna 1, antenna 2 and antenna 3 respectively as shown in Fig. 3 were simulated by CST Microwave Studio electromagnetic simulator.

Fig. 4 illustrates the reflection coefficient S_{11} comparison against the frequency of the designed antenna structures for the three iterations numbers presented in Fig. 3. It can be shown that for antenna 1, dual operating bands $6.08 - 7.57\text{GHz}$ and $8.74 - 12.23\text{GHz}$ with respect to -10dB are obtained. For the antenna 2 and antenna 3 with iteration number 2 and 3 respectively, the central frequency $fr1$ and $fr2$ of the lower and higher frequency bands got shifted toward lower frequency. We can also observe that for the structure 3, the two operating bands obtained $6.00 - 7.71\text{GHz}$ and $8.17 - 12.55\text{GHz}$ are more important than that of the antenna 1 and antenna 2. However, it entails degradation for level bandwidth adaptation for higher frequencies. This signifies that the impedance bandwidth is improved with increasing the number of iterations. But, owing to practical limitations and simulation results, the fourth iterative structure is not designed and simulated. Then, the antenna 3 described above was constructed for measurement use.

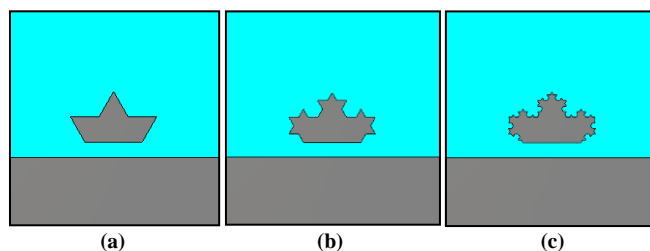


Fig. 3. Iterative Structures of the simulated antennas (a) Antenna 1; (b) Antenna 2; (c) Antenna 3.

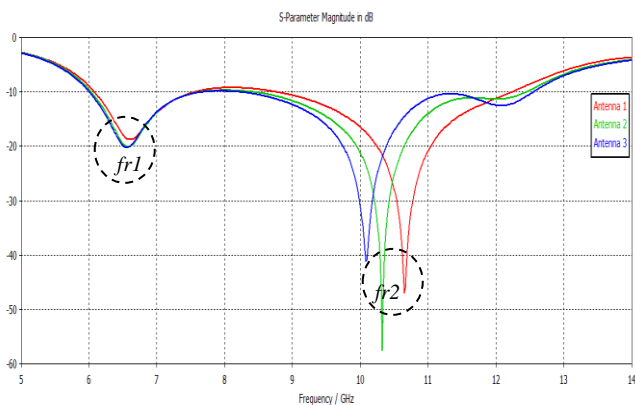


Fig. 4. Simulated Reflection coefficient versus frequency characteristics for iterative structures of designed antenna.

The fabricated prototype of the designed antenna structure is demonstrated in Fig. 5.

Fig. 6 shows the comparison between simulated and measured return loss for the proposed antenna 3. Good agreement can be observed. Exceptionally, for the frequency band [5.24GHz - 6.00GHz]. The frequency drift is produced by the error in the manufacture and measurement. Measured impedance bandwidths for -10 dB return loss of the two operating bands are 890MHz (5.94 – 6.83GHz) and 4.8GHz (9.06 – 13.86GHz).

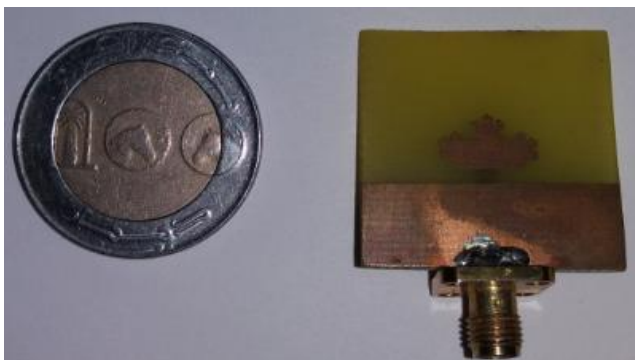


Fig. 5. Photograph of the antenna 3 fabricated prototype.

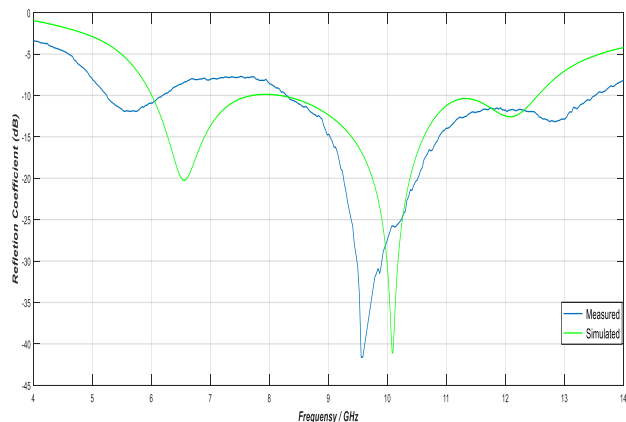


Fig. 6. Simulated and measured reflection coefficients of the proposed antenna 3.

The simulated 3D radiation patterns of the designed antenna 3 structure at some frequency points are illustrated in Fig. 7. It shows acceptable stability of the radiation selected frequency antenna. While the simulated gain (Fig. 8) of the proposed antenna 3 against frequency band [4GHz - 14 GHz] shows that the gain increases for higher frequencies. It is around 2.51dBi, 5.37dBi and 3.58 dBi at 6.54GHz, 10.00GHz and 12.18GHz, respectively.

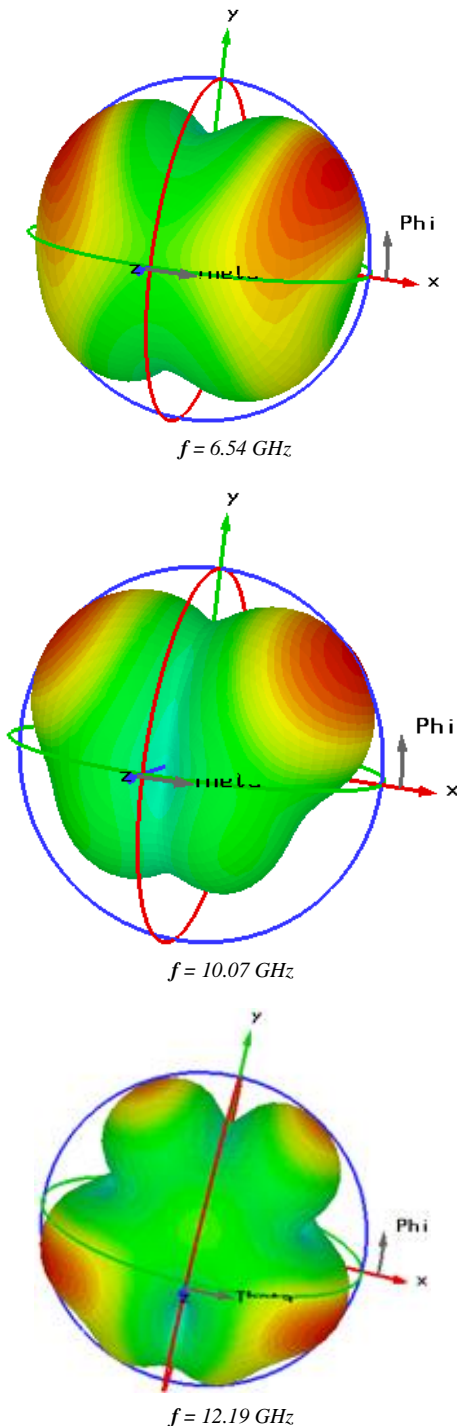


Fig. 7. Simulated 3D radiation patterns of the antenna 3 at some frequencies.

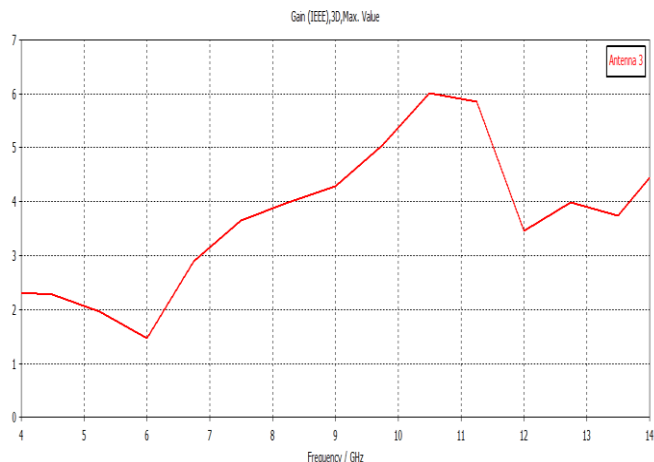


Fig. 8. Variation of simulated Gain in dBi against frequency for designed antenna 3.

B. Effects of Ground Plane (GND)

The rectangular partial ground plane used in the microstrip patch antenna design is one of the most important criteria. So, the modified ground plane has received much attention to ameliorate the antenna characteristics. This kind of study is presented in [24], [25].

The potential benefits of the proposed antenna 3 with the simple ground have been discussed previously. In this section the effects of inserting slots on the antenna 3 ground plane (GND) as shown in Fig. 9 are examined.

The comparison among the reflection coefficient characteristics of the antenna 3 with simple ground plane and final proposed antenna 4 is presented in Fig. 10. It is observed that for antenna 4 which the ground plane is modified by inserting slots, the two operating antenna 3 bands got merged resulting into a wide operating band of 6.59 GHz under the condition of $S_{11} < -10$ dB, from 6.03 GHz to 12.62 GHz with 70.7%. This signifies that the impedance matching improved over the entire band, enhancing the bandwidth of operation.

The final proposed antenna 4 with modified ground plane was constructed for measurement use. The photograph of prototype antenna is shown in Fig. 11.

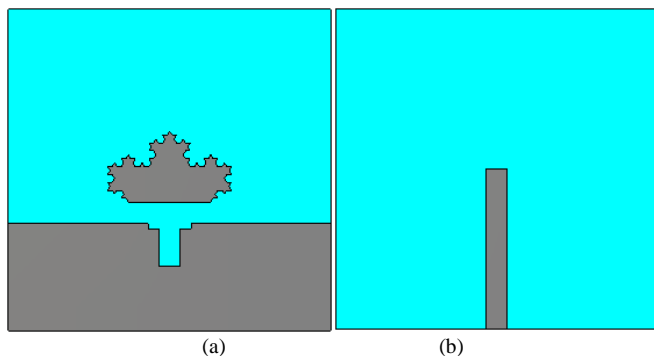


Fig. 9. Proposed final Koch Island antenna (Antenna 4) (a) Back view; (b) Front view.

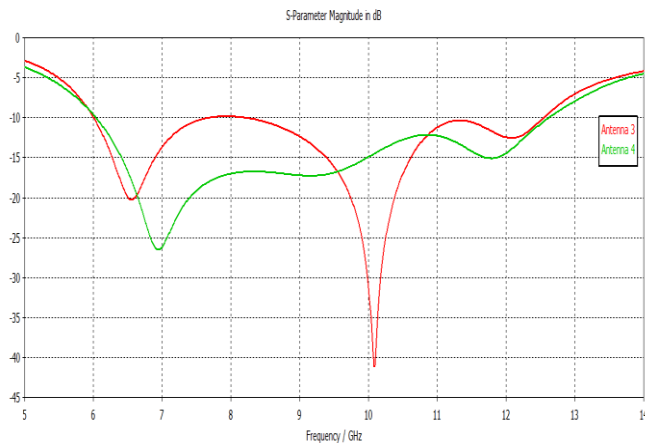


Fig. 10. Simulated reflection coefficient of antenna 3 with simple GND and Koch Island antenna 4 with modified GND.

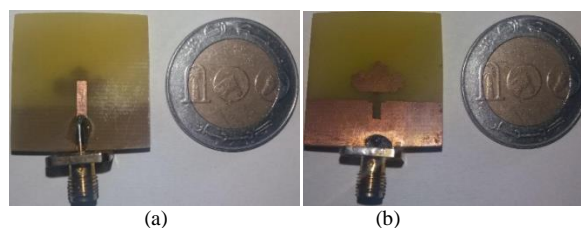


Fig. 11. Photograph of the antenna 4 fabricated prototype. (a) Front view; (b) Back view.

After the realization of the prototype of the proposed antenna 4, the input reflection coefficient is tested. The simulated and measured return loss against frequency of the proposed antenna 4 design is plotted in Fig. 12, show good agreement. Measured reflection coefficient of this wide band antenna indicates that the -10 dB operating bandwidth is 5.94 GHz extends on the interval [6.03 - 11.97GHz]. It is wide sufficient to cover the required bandwidths for x band [8 - 12 GHz] applications. A difference in the adaptation level was observed between simulated and measured results. This is due to the fabrication inaccuracy and measurement conditions.

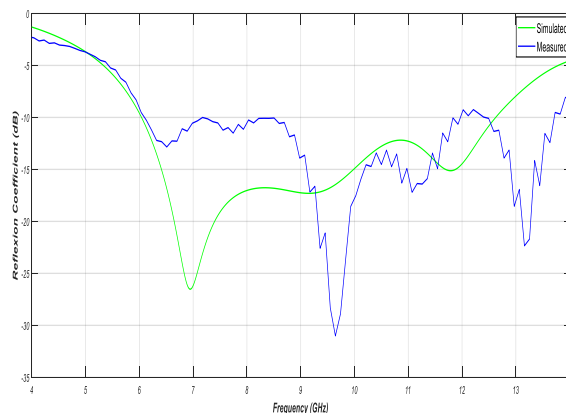


Fig. 12. Simulated and measured reflection coefficients of the final proposed antenna 4.

The simulated radiation patterns at different frequencies of the operating bandwidth shows that the antenna radiates in the

same direction for all selected frequencies as presented in Fig. 13.

Fig. 14 illustrates the variation of the simulated gain versus the frequency interval [5 - 14GHz]. It is observed that the antenna gain is varying from a minimum of 1.5dBi for 6.03GHz to a maximum of 6.28dBi for 10.5 GHz across the desired bandwidth.

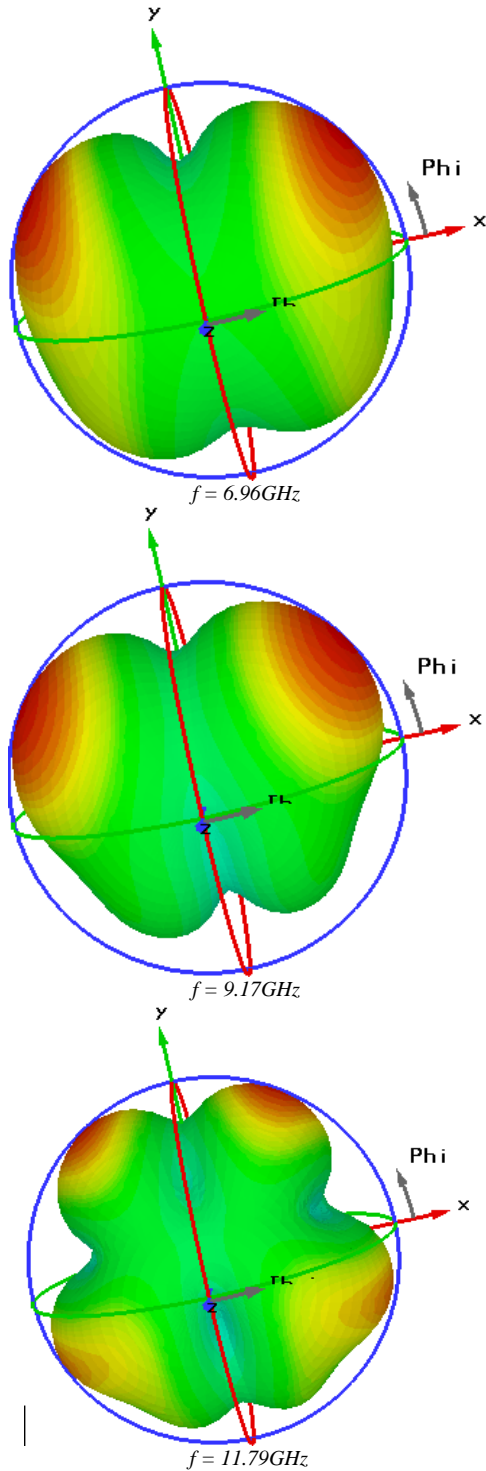


Fig. 13. Simulated 3D Radiation patterns of the proposed antenna 4 at some frequencies.

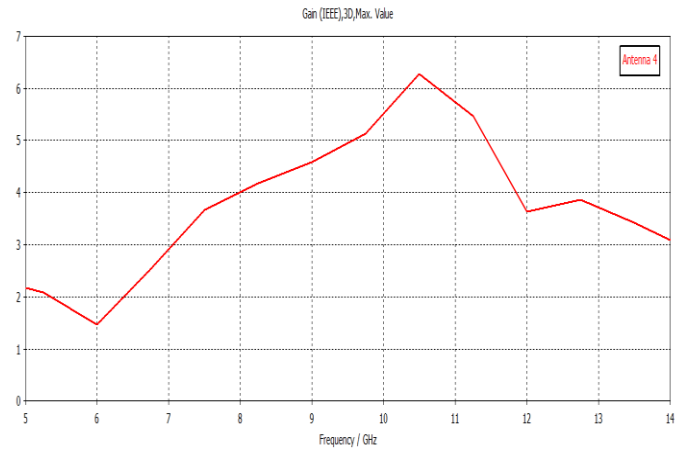


Fig. 14. Variation of simulated Gain in dBi against frequency for designed antenna 4.

IV. CONCLUSION

In this paper, a modified microstrip patch antenna KIFPA structure using the Koch Island fractal geometry for bandwidth enhancement is proposed. The effect of the patch iteration number and the modified ground plane are investigated. The application of Koch fractals to the patch and inserting slots to the simple rectangular partial ground plane are clearly significant in improving the antenna level adaptation performance, enhancing the operating bandwidth of the proposed final antenna. The computer simulations and measured results of the proposed antennas have shown good agreement.

ACKNOWLEDGMENT

The authors would like to gratefully acknowledge Pr. SALAH-BELKHODJA Faouzi from Telecommunications and Digital Signal Processing Laboratory, University of SidiBelabbes for his technical support in providing the experimental data.

REFERENCES

- [1] B. Mandelbrot, "The fractal geometry of nature", Freeman and Company, New York, 1975.
- [2] K Falconer., Fractal Geometry: Mathematical foundations and applications, John Wiley and Sons, Chichester, 1990.
- [3] H.Oraizi and S.Hedayati , "Miniaturization of Microstrip Antennas by the Novel Application of the Giuseppe Peano Fractal Geometries", IEEE Transactions on Antennas and Propagation, Vol. 60, No. 8, pp. 3559-3567, August 2012.
- [4] I Putu Elba Duta Nugraha, I. Surjati, Syah Alam," Miniaturized Minkowski-Island Fractal Microstrip Antenna Fed by Proximity Coupling for Wireless Fidelity Application", TELKOMNIKA, Vol.15, No.3, pp. 1119-1126, September 2017.
- [5] M.Hadji, S. M. Meriah and D. Ziani, " A New Multi Band Microstrip Patch Fractal Antenna for WLAN/WIMAX Applications", Fifth International Conference on Image and Signal Processing and their Applications (ISPA), 2017.
- [6] J. Anguera, E. Martínez, C. Puente, C. Borja, and J. Soler, "Broad-band dual-frequency microstrip patch antenna with modified Sierpinski fractal geometry," IEEE Transactions on Antennas and Propagation, vol. 52, no. 1, pp. 66-73, 2004.
- [7] J. K. Ali, "A new reduced size multiband patch antenna structure based on Minkowski pre-fractal geometry," Journal of Engineering and Applied Sciences, vol. 2, pp. 1120-1124, 2007.

- [8] HuZhangfang, X.We, L.Yuan, HuYinping, Z.Yongxin,"Design of a modified circular-cut multiband fractal antenna", The Journal of China Universities of Posts and Telecommunications, Vol. 23, Issue 6, Pages 68-75, December 2016.
- [9] M.NaghshvarianJahromi, A.Falahati, and Rob. M. Edwards, "Bandwidth and Impedance-Matching Enhancement of Fractal Monopole Antennas Using Compact Grounded Coplanar Wave guide", IEEE Transactions on Antennas and Propagation, Vol. 59, No. 7, pp. 2480-2487, July 2011.
- [10] Y.K. Choukiker, S.K.Behera, "Modified Sierpinski square fractal antenna covering ultra-wide band application with band notch characteristics", IET Microwaves, Antennas & Propagation, Vol. 8, Iss. 7, pp. 506-512, 2014.
- [11] A.Kumar,"Design of Fractal Antenna for Ultra Wide Band Applications", International Journal of Engineering Research and Technology, IJERT, Vol. 4, Issue 07, pp. 541-544, July 2015.
- [12] M. Ali. Dorostkara,b, R. Azima, M. T. Islama, "A Novel Γ -shape Fractal Antenna for Wideband Communications", The 4th International Conference on Electrical Engineering and Informatics (ICEEI), pp. 1285 - 1291.
- [13] A. Altaf, Y. Yang, K.-Yoon Lee,² and Keum Cheol Hwang, "Wideband Circularly Polarized Spidron Fractal Slot Antenna with an Embedded Patch", International Journal of Antennas and Propagation, April 2017.
- [14] H. Oraizi and S. Hedayati, "A Novel Wide Slot Antenna Design using the Giuseppe Peano Fractal Geometry", 20th Iranian Conference on Electrical Engineering, (ICEE2012), May 15-17, Tehran, Iran, 2012.
- [15] N.Sharma and V.Sharma, "A design of Microstrip Patch Antenna using hybrid fractal slot for wideband applications", Ain Shams Engineering Journal, July 2017.
- [16] A. KhannaDinesh, K. SrivastavaJai and P. Saini," Bandwidth enhancement of modified square fractal microstrip patch antenna using gap-coupling", Engineering Science and Technology, an International Journal, Vol18, Issue 2, Pages 286-293, June 2015.
- [17] K.Mandal and P. P. Sarkar, "A Compact Low Profile Wideband U-Shape Antenna With Slotted Circular Ground Plane", AEU - International Journal of Electronics and Communications, Volume 70, Issue 3, Pages 336-340, March 2016.
- [18] R.Kumar and N.Kushwaha" Design And Investigation Of Sectoral Circular Disc Monopole Fractal Antenna And Its Backscattering", Engineering Science and Technology, an International Journal, Volume 20, Issue 1, PP. 18-27, February 2017, Pages 18-27.
- [19] S. Singhal, A. Kumar Singh, "CPW-fed octagonal super-wideband fractal antenna with defected ground structure", IET Microwaves, Antennas & Propagation, Volume: 11, Issue: 3, PP. 370-377, 2017
- [20] CST Microwave Studio Suite 2014.
- [21] M. Rani, R. Ul Haq, D. Kumar Verma," Variants of Koch curve", National Conference on Development of Reliable Information Systems, Techniques and Related Issues (DRISTI), Proceedings published in International Journal of Computer Applications (IJCA), pp.20-25, 2012.
- [22] H.O. Peitgen, H. Jürgens, D. Saupe, "Chaos and Fractals: New Frontiers of Science", Second Edition, pp89-90. Springer-Verlag New York, 2004.
- [23] A. Danideh, R. Sadeghi-Fakhr and H. R. Hassani, "Wideband Co-Planar Microstrip Patch Antenna", Progress In Electromagnetics Research Letters, Vol. 4, pp. 81-89, 2008.
- [24] A. Naghar, F. Falcone, A. Alejos, O. Aghzout and D. Alvarez," A Simple UWB Tapered Monopole Antenna with Dual Wideband-Notched Performance by Using Single SRR-Slot and Single SRR-Shape.d Conductor-Backed Plane" ACES JOURNAL, Vol. 31, No.9,pp 1048-1055, September 2016.
- [25] S. RafathAra And Dr. S. N. Mulgi," Multi Stepped Slotted Partial Ground Plane Dual Notched Ultra Wide Band Rectangular Microstrip Antenna", International Multidisciplinary Research Foundation (IMRF) : Vol. 5, Issue 1, pp 58-64, 2017.

Investigating Saudi Parents' Intention to Adopt Technical Mediation Tools to Regulate Children's Internet Usage

Ala'a Bassam Al-Naim¹, Md Maruf Hasan²

Department of Information Systems
College of Computer Sciences and Information Technology
King Faisal University, Al Ahsa, Saudi Arabia

Abstract—The adverse and harmful effects of Internet on young children have become a global concern. Parents tend to use different strategies to ensure their children's online safety. Many studies have suggested that parental mediation may play a positive role in controlling children's online behavior. The purpose of this study is to identify the factors that shape Saudi parents' intention to regulate their children's online practices using technical mediation tools. An integrated model has been proposed based on famous Information System theories and models to investigate parental intention to adopt technical mediation tools. A questionnaire-based survey is conducted for data collection. Basic descriptive statistical analysis, reliability, and validity assessments were used to analyze the data at the preliminary stage, followed by advanced analysis using Structural Equation Modeling to test the research hypotheses. Research results indicate that effort expectancy, performance expectancy, general computer self-efficacy, perceived severity, and perceived vulnerability are the main predictors of Saudi parent's intention to regulate their children's online behaviors using technical mediation tools.

Keywords—Child and family safety online; parental control and mediation; technology mediation; Unified Theory of Acceptance and Use of Technology (UTAUT); Saudi Arabia

I. INTRODUCTION

Over the last decade, increased use of Internet has been observed worldwide. The Internet became an indispensable utility of our daily life serving multiple purposes such as personal development, education, entertainment, communication, harnessing information for better personal and professional decision-making, etc. Recent research shows that in developed countries it is a home based activity and current generation of young children cannot imagine a world without information and communication technologies (ICT), therefore, they are called digital natives [1]. Although Internet provides a powerful mean to access information and communicates with people to conduct useful activities, it possesses significant risk and threats for children and minors. There are evidences that use of Internet may expose young children to certain risks such as becoming victim of cyber bullying, negative emotional impacts due to unwanted exposure to pornography, violence, explicit language, revealing personal information to sexual predators etc. Since young children lack a sufficient level of maturity to be able to manage these risks, some authors state that this generation is not only to be called "whiz kids" [2], but

also "risk-kids" [3]. Therefore, the Internet impact on children and younger generation has become a global concern resulted in demands for safeguard to protect online privacy when involved with a wide variety of commercial websites and activities. Government in westerns countries is continuously engaged in implementing public policies, framework and legislation to reduce the unnecessary solicitation of personal information from children. One such example is Children's Online Privacy Protection Act (COPPA) implemented by U.S. Federal Trade Commission (FTC) which sets guidelines for online safeguards designed to prevent the collection of personally identifiable information from children unless parental consent is given [4].

However, as the media and communication environment quickly grow it becomes increasingly difficult for governments to formulate legal frameworks and enforce them. Policy makers rely substantially on increasing risk awareness among parents and delegating to them the responsibility for protecting children from online risks. Here the value of parental role becomes critical in view of safe Internet usage and Internet education. Active monitoring of children's online activities by parents can help guard children from the threats of the Internet. Parents implement a range of strategies, favoring active co-use and interaction rules over technical restrictions using filters or monitoring software. Active mediation involves parental guidance and advice through active discussions over online issues and staying nearby or sitting with children when they go online to monitor closely to reduce the likelihood of undesirable and damaging behaviors or attitudes.

However, parents find it increasing difficult to monitor their children actively as well as children don't like intense instructive. This results in development of huge number of technical filtering/monitoring tools and professional computer protection software's. However, studies on parents preferences [5] shows that although more than 75 per cent of parents were concerned about privacy risk and exposure to sexual content risk, only a smaller proportion of parents install filtering software (33 per cent) and monitoring software (23 per cent). Although technical mediation seems less instructive and allows more insight and control over the young Internet users. The main reason for not adopting technology is insufficient technical and Internet usage skills.

It is more challenging when it comes to Saudi parents' context as compare to western world, due to huge cultural difference. The laws implemented in western countries are not acceptable in Saudi Arabia due to tight religion teaching. Furthermore, filtering and monitoring software could not directly implement in Saudi culture due to difference in social/religion norms. Since, what considered normal in western world might conflict with Islamic law. As a result, many sites were blocked at country level. When it comes to parent's level, most parents strictly banned their children to use internet and depriving them to get benefit of rich information content. On the other hand, if they allow their children to use Internet sufficient parental control strategies were not adopted, as a result their children are under constant online risk. There is a pressing need of educating Saudi parents to make them aware of using different technical tool to keep them safe. This main goal of the study is to identify factors that's determines Saudi parental intention to adopt technology mediation for safe Internet usage. Based preliminary on Protection Motivation Theory (PMT) [6], Unified Theory of Acceptance and Use of Technology (UTAUT) [7], and Technology Acceptance Model (TAM) combined with parental perspective an integrated model is proposed. Following main factors were investigated in proposed model to identify parents' acceptance of technology as mediation tool:

- Investigating Saudi patterns of parental mediation on children's online activities.
- Identify parent's awareness on severity of online risk.
- Level of parent's awareness of children internet usage.
- Parental level of internet and technology usage skills, level of comfort and easiness in using those skills.
- Exploring the Saudi parent's perception about technical mediation in providing the required protection for their children.
- Predicting parent's intention to adopt the technology to keep their children safe on Internet.

Rest of paper is organized as follow. Section II briefly explaining background of different parental mediation and highlighting technical mediation is least adopted. The proposed model and research method is presented in Sections III and IV. Research results are discussed in Section V. Finally, Section VI concludes the paper.

II. RELATED WORK

Internet impact on children and younger generation has become a global concern. The impact of parental involvement has been recognized in studies of media and children. In research dealing with television viewing term "mediation" is used as an effort to manage the relationships between the child and mass media [8]. Parents tend to use different strategies and practices to ensure that their children's usage of mass media such as TV, video games, and Internet [9] is under control. Two general mediation strategies have been identified in the context of interaction with mass media Active mediation and Restrictive mediation [10]. Active mediation consists of monitoring children media usage actively by sitting beside

them and talking about media content while child is watching, reading, or listening the medium. Hence, instead of learning directly from environment they learn from their parents. In Restrictive mediation parents set limits for viewing or prohibit the viewing of certain content e.g. in context of TV setting number of hours a child can watch television or prohibit the viewing of a certain program or channel.

Do parents employ similar mediation strategies for the Internet? Some studies [9] stated that same strategies can be applied on internet usage. However, unlike other mass media internet is not shared activity. In addition, it is more interactive as compare to other media, so children are at highest risk. Therefore, parental awareness, and more comprehensive parental mediation strategies are requiring. Among actively monitoring their children and setting rules for internet usage there is a need certain tools to keep children safe online.

Green et al. [11] suggested technical based parental mediation by using parental control software settings for monitoring, filtering, and restricting unwanted contents. However, studies on parents preferences [5] shows only a smaller proportion of parents install filtering software (33 per cent) and monitoring software (23 per cent). Although technical mediation seems less instructive and allows more insight and control over young Internet users. Some studies concluded that using technical tools and settings is the least adopted type of mediation among other types [8], [12]. Besides that, most of parental mediation studies focused on other types of mediation strategies mentioned above while the technical type was partially or completely ignored [9], [13]. The literature did not provide a clear clue of why this type is the least adopted mediation style. Most studies mentioned the most and least used mediation styles with focusing on the most used ones and why they are being used mostly.

As we have explored personally some parental monitoring software and settings on different platforms, we believe in their ability to achieve what the normal styles could not. To fill the gap in current literature, it is important first to predict the reasons of why or why not a parent could think of adopting and accepting technical settings or software to regulate their children Internet usage. It has been observed from technology acceptance studies that the user intention is the main predictor of the usage behavior. These studies addressed the behavioral intention and its aspects in order to know how it can change the amount and tune of usage. The user intention to use a technology can be used to explore the amount of technology adoption or whether a person will use a technology or not [6], [14]. It has been observed from literature [8], [12], [15] that Europe and the US [6], [9], [18] is spending considerable amount of research in this area. When it comes to Saudi Arabia, it seems that Saudi researchers have entered this domain recently by Almoqbel et al. [16]. According to author they are first to study how Saudi parental educational and economic level's impact on their children internet usage habits. The aforementioned study is a social based research with narrow scope about the relationship among parent's specific demographics and mediation behavior. Above all, it seems that parental technical mediation topic was not covered enough by IS research papers. This study is first in Saudi context to explore different factors that can influence Saudi parents'

intention to adopt technical mediation tools to regulate their children Internet usage.

In the next section, an integrated theoretical model is proposed and hypothesis was formulated to find relationship between different factors.

III. RESEARCH MODEL AND HYPOTHESES

Protection Motivation Theory (PMT) [6], Unified Theory of Acceptance and Use of Technology (UTAUT) [7], and Technology Acceptance Model (TAM) are selected as the reference to develop theoretical model for determining parents behavioral intention to adopt the technology as a mediation style. As shown in Fig. 1, proposed model contains six independent constructs, one dependent construct, and five moderators. The variables were either selected based on previous IS studies or self-constructed. Perceived Vulnerability (PV), Perceived Severity (PS), and Self Efficacy (SE) were selected from PMT. Whereas, Effort Expectancy (EE), and Performance Expectancy (PE) adopted from UTAT. Finally, Behavioral Intention was adopted from both UTAUT and TAM as a main indicator for accepting and using the technology. Furthermore, we have extended the model with self-constructed items such as Awareness of children online use [9, 17], Parent’s Internet use [16], [18], Child Age [9], [13], [19], Educational level of Parents [16], Gender of Parents [9], [15], [17] and Parents Age [7], [13], [15]. Details of constructs were presented in Table I.

TABLE I. CONSTRUCTS WITH THEORETICAL DEFINITIONS

Construct	Theoretical Definition
Behavioral Intention (BI)	Person’s willingness and readiness to adopt specific behavior [20, 21].
Perceived Vulnerability (PV)	Person’s evaluation of the possibility of exposure to such threat [22].
Perceived Severity (PS)	Measures the intensity of the consequences to a person or others if the related threat increased [23].
Performance expectancy (PE)	The degree to which a person believes that using a technology will help him to gain the benefits he wants [20].
Effort Expectancy (EE)	The level of easiness related to the use of a technology [22].
Awareness of children online use (AW)	Parents become aware about what their children do online [9, 17]
General computer self-efficacy (GCSE)	Person’s judgment of his ability and skills to any task related to computer [24].

Following hypotheses were constructed to examine proposed model of Saudi parent’s intention to adopt technical mediation tools on children Internet usage.

H1: Perceived vulnerability of online risks for children will significantly predict the parent’s intention

- **H1a:** Parent’s Internet use will positively moderate the relationship between the perceived vulnerability of online risks for children and Saudi parent’s intention.
- **H1b:** Child’s age will negatively moderate the relationship between the perceived vulnerability of online risks for children and Saudi parent’s intention.

- **H1c:** Parent’s educational level will positively moderate the relationship between the perceived vulnerability of online risks for children and Saudi parent’s intention.

H2: Perceived severity of online risks for children will significantly predict the Saudi parent’s intention

H3: Performance expectancy will significantly predict the Saudi parent’s intention

- **H3a:** Gender will positively moderate the relationship between the performance expectancy and parent’s.
- **H3b:** Age will negatively moderate the relationship between the performance expectancy and parent’s.

H4: Effort expectancy will significantly predict the parent’s intention

- **H4a:** Gender will positively moderate the relationship between the effort expectancy and parent’s intention.
- **H4b:** Age will negatively moderate the relationship between the effort expectancy and parent’s intention.

H5: Parent’s awareness of children online usage will significantly predict the parent’s intention

- **H5a:** Parent’s age will negatively moderate the relationship between the awareness of child’s online activities and parent’s intention.

H6: General computer self-efficacy will significantly predict the parent’s intention

IV. RESEARCH METHODOLOGY

A. Instrument Development and Design

In order to test research hypotheses, a survey questionnaire has been designed to collect the data from the targeted sample. The survey consists of 27 questions divided into three parts: First part is on parent’s general information, second part is about child and his/her online usage, and third part contained question used to operationalize research hypothesis. Five point Lickert scale ranging from (1) strongly disagrees to (5) strongly agree and choice based measures are used to collect responses. Questionnaires were reviewed by three professors to check the face and content validity. We translated the questionnaire into Arabic to be appropriate for Saudi people. The translated survey with the original English was reviewed by graduate student from English language department. The goal was checking the validity of translation process. The Arabic version was presented to some parents (1 male and 2 female) to check the clarity of terminologies used.

Table II shows survey items for each construct along with supporting studies.

B. Tools and Techniques

Two quantitative approaches were used to analyze the collected data. First, demographic data is analyzed via descriptive statistics using Statistical Package for Social Sciences (SPSS). Second, WarpPLS 5.0 [25] is used to assess the reliability, validity as well as hypothesis testing. Structural Equation Modeling (SEM) was used for hypothesis testing.

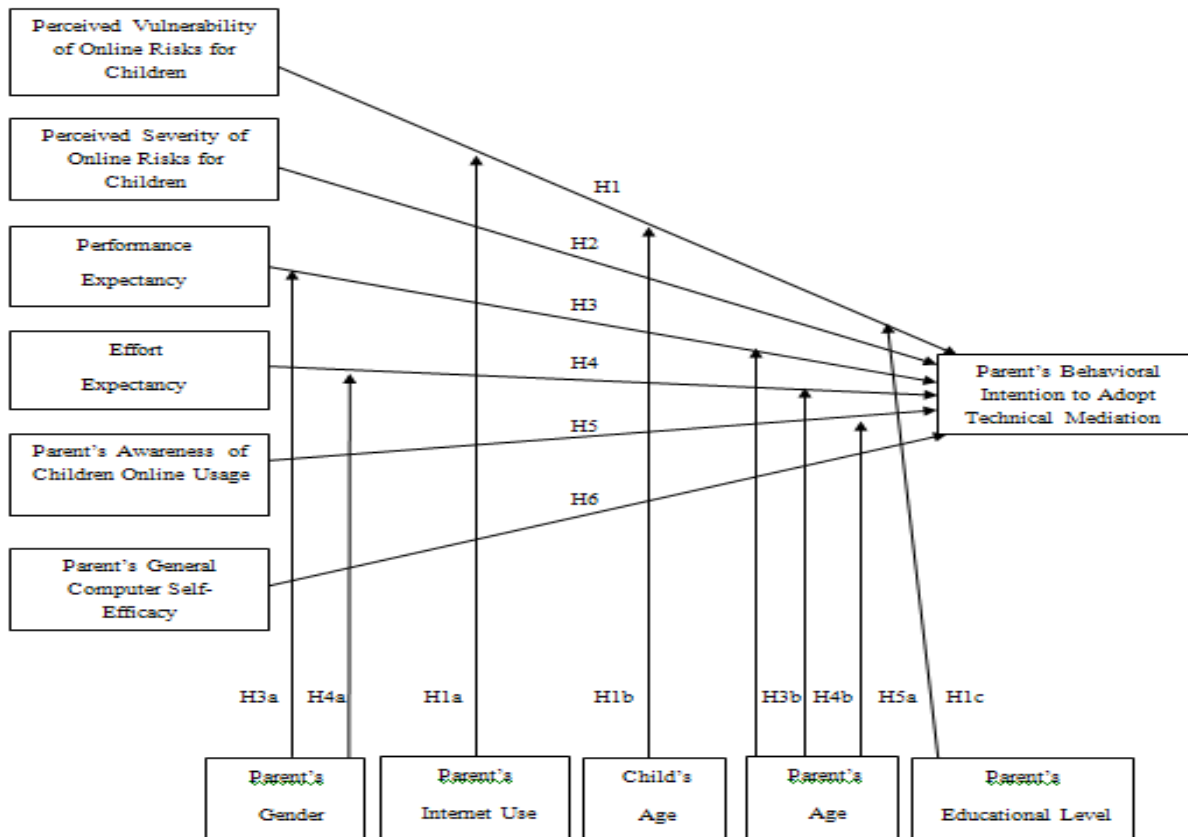


Fig. 1. Research hypothesized model.

TABLE II. CONSTRUCTS AND SURVEY ITEMS

Construct	Survey Items
BI [26]	BI1: I will adopt technical mediation to regulate my child Internet usage.
	BI2: What are the chances in 100 that you will adopt technical mediation to regulate your child Internet usage.
	BI3: To regulate my child's Internet usage, I would adopt technical mediation rather than any other means available.
PV [6]	PV1: My children are exposed to one or more online risks.
	PV2: It is likely that my children will be exposed to one or more online risks.
	PV3: It is possible that my children will be exposed to one or more online risks.
PS [6]	PS1: I believe that my child's exposure to one or more online risks would be a severe problem.
	PS2: I believe that my child's exposure to one or more online risks would be a serious problem.
	PS3: I believe that my child's exposure to one or more online risks would be a significant problem.
PE [27]	PE1: I find using parental monitoring software or parental control settings will be useful in my mediation.
	Using parental monitoring software or PE2: Parental control settings will enable me to accomplish the mediation task more quickly.
	PE3: Using parental monitoring software or parental control settings will enhance more online safety for my children.
EE [27]	EE1: It will be easy for me to become skillful at using parental monitoring software or parental control settings to regulate my child's Internet usage.
	EE2: I find that parental monitoring software or parental control settings will be easy to use for mediating my child's

	Internet usage.
	EE3: Learning to use parental monitoring software or parental control settings to regulate my child's Internet usage will be easy for me.
AW (Self-Constructed)	AW1: I know exactly what does my child do online.
	AW2: I know exactly how much time my child spends online.
GCSE [28]	GCSE1: I believe I have the ability to unpack and set up a new computer.
	GCSE2: I believe I have the ability to install new software applications on a computer.
	GCSE3: I believe I have the ability to use a computer to display or present information in a desired manner.
	GCSE4: I believe I have the ability to identify and correct common operational problems with a computer.
	GCSE5: I believe I have the ability to remove information from a computer that I no longer need.

C. Demographic Characteristics

Data was collected both by means of an online platform and in person by distributing copies. The total number of responses collected were 196 out of 280 copies which shows 70% response rate. Total 55 responses were collected online, hence total number of responses for both paper based and online was 251. Table III shows the general characteristics of respondents. Majority of respondents were females with percentage of (61.4%). Statistics also show that (76.5%) of respondents were university graduates and (59.4%) were using the Internet frequently. The descriptive statistics of the participants' responses are presented in Table IV.

TABLE III. GENERAL CHARACTERISTICS OF THE SAMPLE

Measure	Item	Frequency	Percentage
Parent's gender	Male	97	38.6 %
	Female	154	61.4 %
Parent's age	25-35 years	119	47.4 %
	36-46 years	91	36.3 %
	47-57 years	40	15.9 %
	Over 57 years	1	0.4 %
Educational level	Illiterate	0	0 %
	Elementary	1	0.4 %
	Intermediate	9	3.6 %
	Secondary	33	13.1 %
	University graduate	192	76.5 %
	Postgraduate (PhD or Master)	16	6.4 %
Parent's Internet use	Frequently	149	59.4 %
	Occasionally	92	36.6 %
	Rarely	9	3.6 %
	Never	1	0.4 %
Child's Age	6-7 years	82	32.7 %
	8-9 years	58	23.1 %
	10-11 years	58	23.1 %
	12-13 years	53	21.1 %

V. DATA ANALYSIS AND RESULTS

Two types of quantitative approaches were used to analyze the collected data. First, reliability and validity of the questionnaire is estimated to determine the adequacy of measuring items. Second, SEM is built to test whether proposed hypothesis are supported by data or not.

A. The Assessment of the Measurement Model

Model reliability, Convergent and Discriminant validity [29] is examined to identify the adequacy of measurement model. Construct validity is defined as the degree to which the operational measurement actually reflects the true theoretical meaning of a concept/construct [30]. It was conducted through calculating composite reliability and Average Variance Extracted (AVE) for each latent construct. For a construct, the convergent validity is achieved when it scores a composite reliability above 0.70 and AVE value above 0.50 [31]. Table V shows the results of convergent validity assessment where all constructs reported composite reliability values greater than 0.70 and AVE values greater than 0.50. Discriminant validity was assessed through calculating the square root of AVE for each latent construct and comparing it with the inter-correlations among model constructs. As shown in Table VI, the square root values for each latent construct are exceeded all its correlations with other constructs in the model and so the discriminant validity requirements were met.

TABLE IV. CONVERGENT VALIDITY STATISTICS

Construct	Composite Reliability	AVE
PV	0.877	0.704
PS	0.911	0.774
PE	0.909	0.768
EE	0.903	0.756
AW	0.785	0.647
GCSE	0.860	0.553
BI	0.816	0.597

TABLE V. DISCRIMINANT VALIDITY STATISTICS

	PV	PS	PE	EE	AW	GCSE	BI
PV	0.839						
PS	0.470	0.879					
PE	0.167	0.225	0.876				
EE	0.074	0.215	0.541	0.869			
AW	0.095	0.156	0.088	0.047	0.804		
GCSE	0.050	0.196	0.166	0.315	0.172	0.743	
BI	0.133	0.247	0.513	0.620	0.020	0.290	0.772

Cronbach's coefficient alpha (CA) was calculated to examine the internal consistency for each construct. The results are presented in the Table VII. Generally, all constructs reported appropriate reliability values above 0.70 except behavioral intention (BI) and parent's awareness of child's Internet usage (AW) which scored 0.662 and 0.453 respectively. Therefore, we have used the statistical re-sampling of data to generate simulated data and re-examine the reliability of constructs. Results have shown that all constructs scored composite reliability coefficients above the critical threshold value of 0.70 [32].

TABLE VI. RELIABILITY STATISTICS FOR RESEARCH MODEL

Construct	No. Of Items	Cronbach's Alpha	Composite Reliability
PV	3	0.788	0.877
PS	3	0.854	0.911
PE	3	0.849	0.909
EE	3	0.839	0.903
AW	2	0.453	0.785
GCSE	5	0.795	0.860
BI	3	0.662	0.816

B. The Assessment of the Structural Model

The structural model is evaluated and hypothesis is tested after establishing adequacy of measurement model. SEM is most commonly used multivariate technique [33] for instrument validation and model testing to identify series of relationship constitutes in large-scale model or an entire theory. Following coefficients are calculated to assess SEM: P values (probabilities), path coefficients (Beta coefficients), and R² (explanatory variance power). The results indicate clearly which hypotheses were supported and which of them were not based on P values and Beta coefficients. If the beta coefficient was positive (in the right direction) and the probability value was significant (below one of critical thresholds: 0.05, 0.01, or 0.001), the hypothesis will be supported. In this study, we compared the p values to the three levels of significance

following what has been done by S. Al-Gahtani et al. [32]. R^2 is goodness-of-fit measurement that shows how the behavior of dependent construct is controlled by the behavior of independent construct. It acts as a descriptive or explanatory power of the model to explain the study constructs [34]. As high the R^2 value of the model is as it has high descriptive power.

C. Hypotheses Testing Results

Hypotheses testing results for the hypothesized model are illustrated in Fig. 2 and explained clearly in Table VII. The results of hypotheses testing are reported as follows:

- **H1** (PV→BI) Perceived vulnerability of online risks for children significantly predicts the Saudi parent’s intention to adopt technical mediation tools on children Internet usage. This hypothesis is supported (Beta = 0.133, $P^* = 0.036 < 0.05$).
- **H2** (PS→ BI): Perceived severity of online risks for children significantly predicts the Saudi parent’s

intention to adopt technical mediation tools on children Internet usage. This hypothesis is supported (Beta = 0.247, $P^{***} < 0.001$).

- **H3** (PE→BI): Performance expectancy significantly predicts the Saudi parent’s intention to adopt technical mediation tools on children Internet usage. This hypothesis is supported (Beta = 0.513, $P^{***} < 0.001$).
- **H4** (EE→BI): Effort expectancy significantly predicts the Saudi parent’s intention to adopt technical mediation tools on children Internet usage. This hypothesis is supported (Beta = 0.620, $P^{***} < 0.001$).
- **H5** (AW→BI): Saudi parent’s awareness of children online usage does not predict the Saudi parent’s intention to adopt technical mediation tools on children Internet usage. This hypothesis is not supported (Beta = 0.020, $P = 0.750$).

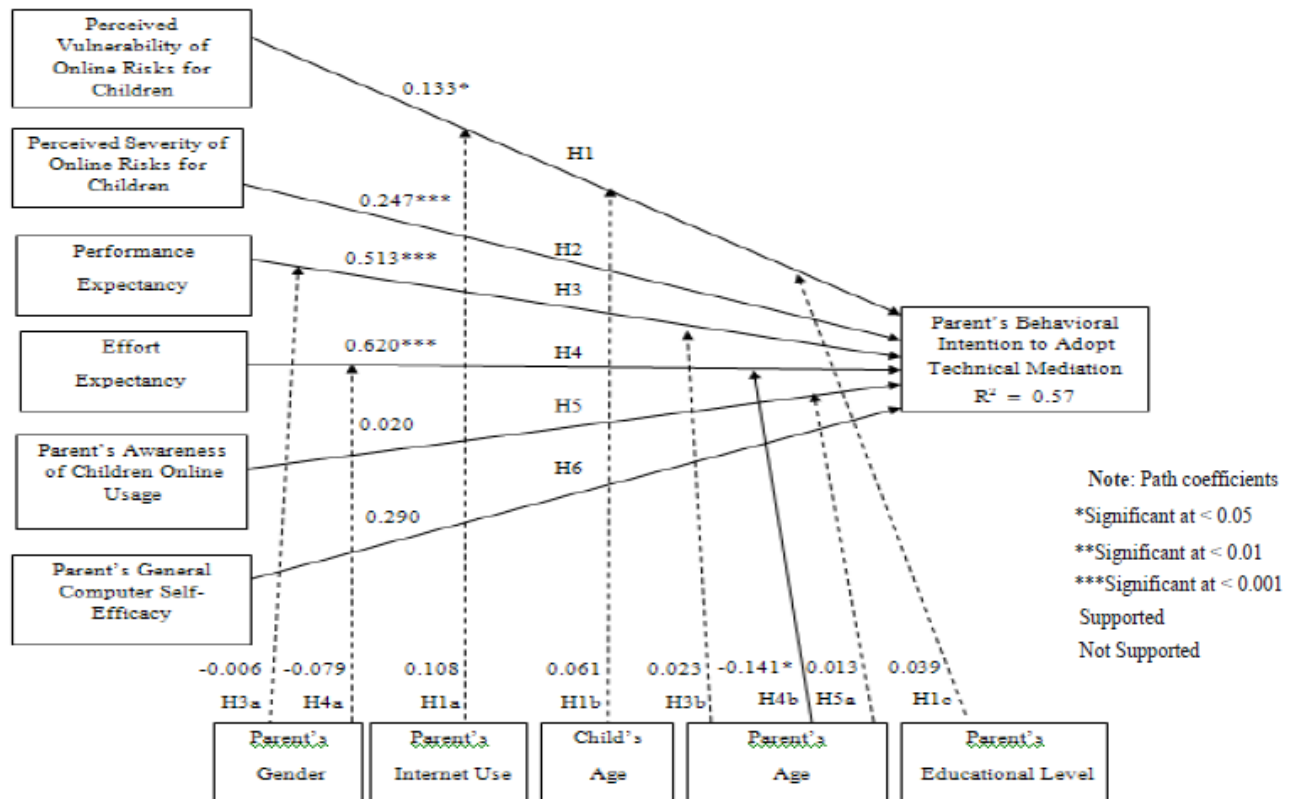


Fig. 2. Research model hypotheses results.

- **H6** (GCSE→BI): Saudi parent’s general computer self-efficacy significantly predicts the Saudi parent’s intention to adopt technical mediation tools on children Internet usage. This hypothesis is supported (Beta = 0.290, $P^{***} < 0.001$).
- **H1a** (P_IU→H1): Saudi parent’s Internet use does not moderate the relationship between the perceived vulnerability of online risks for children and Saudi parent’s intention to adopt technical mediation tools on

children Internet usage positively. This hypothesis is not supported (Beta = 0.108, $P = 0.086$).

- **H1b** (C_Age→H1): Child’s age does not moderate the relationship between the perceived vulnerability of online risks for children and Saudi parent’s intention to adopt technical mediation tools on children Internet usage negatively. This hypothesis is not supported (Beta = 0.061, $P = 0.335$).

- **H1c** (P_Edu →H1): Saudi parent’s educational level does not moderate the relationship between the perceived vulnerability of online risks for children and Saudi parent’s intention to adopt technical mediation tools on children Internet usage positively. This hypothesis is not supported (Beta = 0.039, P = 0.538).
- **H3a** (P_Gende→H3): Saudi parent’s gender does not moderate the relationship between the performance expectancy and Saudi parent’s intention to adopt technical mediation tools on children Internet usage positively. This hypothesis is not supported (Beta = - 0.006, P = 0.920).
- **H3b** (P_Age→H3): Saudi parent’s age does not moderate the relationship between the performance expectancy and Saudi parent’s intention to adopt technical mediation tools on children Internet usage negatively. This hypothesis is not supported (Beta = 0.023, P = 0.713).
- **H4a** (P_Gender→H4): Saudi parent’s gender does not moderate the relationship between the effort expectancy

and Saudi parent’s intention to adopt technical mediation tools on children Internet usage positively. This hypothesis is not supported (Beta = -0.079, P = 0.211)

- **H4b** (P_Age→H4): Saudi parent’s age moderates the relationship between the effort expectancy and Saudi parent’s intention to adopt technical mediation tools on children Internet usage negatively. This hypothesis is supported (Beta = -0.141, P* = 0.026 < 0.05).
- **H5a** (P_Age→H5): Saudi parent’s age does not moderate the relationship between the awareness of child’s online activities and Saudi parent’s intention to adopt technical mediation tools on children Internet usage negatively. This hypothesis is not supported (Beta = 0.013, P = 0.843).

For R squared coefficient, note that each independent variable accounts 57% of the explanatory variance toward BI. In total, the hypothesized model explains 57% of the explanatory variances for BI.

TABLE VII. RESEARCH MODEL HYPOTHESES RESULTS

Hypotheses	Beta value	P value	R ²	Conclusion
H1(PV → BI)	0.133	0.036 P* < 0.05	57%	Supported
H2(PS → BI)	0.247	P***<0.001	57%	Supported
H3(PE→BI)	0.513	P***<0.001	7%	Supported
H4(EE → BI)	0.620	P***<0.001	7%	Supported
H5(AW → BI)	0.020	0.750	57%	Non supported
H6 (GCSE → BI)	0.290	P***<0.001	57%	Supported
H1a (P_IU→H1)	0.108	0.086	N/A	Non supported
H1b(C_Age →H1)	0.061	0.335	N/A	Non supported
H1c (P_Edu →H1)	0.039	0.538	N/A	Non supported
H3a (P_Gender→ H3)	-0.006	0.920	N/A	Non supported
H3b (P_Age → H3)	0.023	0.713	N/A	Non supported
H4a (P_Gender →H4)	-0.079	0.211	N/A	Non supported
H4b (P_Age →H4)	-0.141	0.026 *P < 0.05	N/A	Supported
H5a (P_Age → H5)	0.013	0.843	N/A	No supported

VI. CONCLUSION AND IMPLICATIONS OF THE STUDY

Internet technology is continuously evolving leaving the parents worry regard its hidden threats and dangers on their children. By mediating them appropriately, parents can promote a safe online environment for their children. Different regulation strategies have appeared to apply more control on child’s Internet use. This survey- based study has investigated the factors that shape the intention to use these tools in regulating Internet use of children in Saudi Arabia. Our study found that the main predictors are perceived vulnerability,

perceived severity, performance expectancy, effort expectancy, and general computer self-efficacy. Effort expectancy was found as the most powerful predictor while the perceived vulnerability was the least powerful one. The hypothesis regarding parent’s age and effort expectancy was supported among the other moderating variables. Child and family safety domain is still considered young in Saudi context. This domain needs to be enriched with serious efforts from both public and private sectors. Saudi researchers, programmers, and educators need to work together to raise the Saudi public awareness on how to employ the Internet as an useful and harmless tool. We

envisage that software designed and developed based on this research finding will accommodate effective parental mediation and ensure child and family safety. This study provides the following significant implications:

- Effort expectancy and performance expectancy were found as the most powerful predictors of behavioral intention to use technical mediation. Popular protection software vendors need to understand and shape the different needs of different users. Most of these parental software do not support the Arabic language in its settings, considering the language is important for Saudi parents to make their interaction with software easier since most of them do not talk the English.
- Also, a clear absence of Arabic parental protection software in the market was observed, Arabic developers can benefit from the current research to predict how to attract the public through designing simple fully-functional software; e.g. we found that parents thought about the software abilities to provide their children with the required online protection. This result may inspire the developers to add good protective features.
- The current findings may inspire other researchers to test this model on different contexts and bigger sample sizes.
- The findings could also help those educators or people working in social awareness organizations to design awareness tools to aware the parents of the importance of technology in saving their children online.
- The research is presenting R2 as goodness-of-fit measure which has scored more than 50% among the results.
- Further research could be done by considering and testing other factors like social influence or financial barriers which could increase the model explanatory power.
- The model drives its power from known IS models and theories like PMT and TAM as it borrows some of their constructs.

ACKNOWLEDGMENTS

This research was supported by the Master of Science in Computer Information System (MS-CIS) program at the Colleges of Computer Sciences and Information Technology, at King Faisal University. The authors like to thank Dr. Shaheen Khatoon and Dr. Majed Alshamari for their useful comments and feedback over two semesters as members of the supervisory committee. They are also indebted to Dr. Mohamed Elhassan for his generous help with the model development and verification as well as the data analysis and software support.

REFERENCES

- [1] M. Prensky, "Digital Natives, Digital Immigrants Part II. Do they really think differently? Retrieved 14 May 2008," *On the Horizon*, vol. 9, pp. 1-9, 2001.
- [2] S.-J. Lee and Y.-G. Chae, "Children's Internet use in a family context: Influence on family relationships and parental mediation," *CyberPsychology & Behavior*, vol. 10, pp. 640-644, 2007.
- [3] G. Kuipers, "The social construction of digital danger: debating, defusing and inflating the moral dangers of online humor and pornography in the Netherlands and the United States," *New Media & Society*, vol. 8, pp. 379-400, 2006.
- [4] COPPA. (1998). Children's Online Privacy Protection Act. Available: <http://www.ftc.gov/ogc/coppa1.htm>
- [5] J. Turov and L. Nir, "The Internet and the family: The view of US parents," *Children in the new media landscape*, pp. 331-348, 2000.
- [6] R. E. Crossler, "Protection motivation theory: Understanding determinants to backing up personal data," presented at the System Sciences (HICSS), 2010 43rd Hawaii International Conference on, 2010.
- [7] A. A. Taiwo and A. G. DOWNE, "THE THEORY OF USER ACCEPTANCE AND USE OF TECHNOLOGY (UTAUT): A META-ANALYTIC REVIEW OF EMPIRICAL FINDINGS," *Journal of Theoretical & Applied Information Technology*, vol. 49, 2013.
- [8] R. Hechanova and R. Ortega-Go, "The Good, the Bad and the Ugly: Internet Use, Outcomes and the Role of Regulation in the Philippines," *The Electronic Journal of Information Systems in Developing Countries*, vol. 63, 2014.
- [9] S. N. Hamade, "Parental Awareness and Mediation of Children's Internet Use in Kuwait," presented at the Information Technology-New Generations (ITNG), 2015 12th International Conference on, 2015.
- [10] A. I. Nathanson, "Parent and child perspectives on the presence and meaning of parental television mediation," *Journal of Broadcasting & Electronic Media*, vol. 45, pp. 201-220, 2001.
- [11] L. Green, D. Brady, K. Ólafsson, J. Hartley, and C. Lumby, "Risks and safety for Australian children on the internet," *Cultural Science Journal*, vol. 4, 2011.
- [12] A. Duerager and S. Livingstone, "How can parents support children's internet safety?," *EUKids Online*, 2013.
- [13] M. Valcke, S. Bonte, B. De Wever, and I. Rots, "Internet parenting styles and the impact on Internet use of primary school children," *Computers & Education*, vol. 55, pp. 454-464, 2010.
- [14] Y. Li, "Theories in online information privacy research: A critical review and an integrated framework," *Decision Support Systems*, vol. 54, pp. 471-481, 2012.
- [15] M. Álvarez, A. Torres, E. Rodríguez, S. Padilla, and M. Rodrigo, "Attitudes and parenting dimensions in parents' regulation of Internet use by primary and secondary school children," *Computers & Education*, vol. 67, pp. 69-78, 2013.
- [16] A. Almogbel, M. Begg, and S. H. Wilford, "Analysis of the relationship between Saudi Arabia parents' education and economic level parental control of internet usage," 2015.
- [17] P. Nikken and J. Jansz, "Parental mediation of young children's Internet use," *EU Kids Online*, 2011.
- [18] G. Anderson, D. Ktoridou, N. Eteokleous, and A. Zahariadou, "Exploring parents' and children's awareness on internet threats in relation to internet safety," *Campus-Wide Information Systems*, vol. 29, pp. 133-143, 2012.
- [19] C. Ponte and J. A. Simões, "Asking parents about children's internet use: comparing findings about parental mediation in Portugal and other European countries," in *EU Kids Online-Final Conference*. London, 2009.
- [20] V. Venkatesh and F. D. Davis, "A theoretical extension of the technology acceptance model: Four longitudinal field studies," *Management science*, vol. 46, pp. 186-204, 2000.
- [21] Mariam Al-Khalifa, Shaheen Khatoon, Azhar Mahmood, and I. Fatima, "Factors Influencing Patients' Attitudes to Exchange Electronic Health Information in Saudi Arabia: An Exploratory Study," *International Journal of Advanced Computer Science and Applications*, vol. 7, pp. 197-204, 2016.
- [22] A. Vance, M. Siponen, and S. Pahlila, "Motivating IS security compliance: insights from habit and protection motivation theory," *Information & Management*, vol. 49, pp. 190-198, 2012.
- [23] G. Alnajjar, M. Mahmuddin, and R. Thurasamy, "A conceptual model of mobile commerce acceptance in collectivist cultures," presented at the

- Innovation Management and Technology Research (ICIMTR), 2012 International Conference on, 2012.
- [24] J. L. Claggett and D. L. Goodhue, "Have IS researchers lost bandura's self-efficacy concept? A discussion of the definition and measurement of computer self-efficacy," presented at the System Sciences (HICSS), 2011 44th Hawaii International Conference on, 2011.
- [25] N. Kock, "WarpPLS 5.0 user manual," Laredo, TX: ScriptWarp Systems, 2015.
- [26] F. D. Davis, R. P. Bagozzi, and P. R. Warshaw, "User acceptance of computer technology: a comparison of two theoretical models," *Management science*, vol. 35, pp. 982-1003, 1989.
- [27] N. B. Osman, "Extending the Technology Acceptance Model for Mobile Government Systems," *development*, vol. 5, p. 16, 2013.
- [28] R. D. Johnson, "Gender Differences in E-Learning: Communication, Social Presence," *Innovative Strategies and Approaches for End-User Computing Advancements*, p. 175, 2012.
- [29] D. T. Campbell and D. W. Fiske, "Convergent and discriminant validation by the multitrait-multimethod matrix," *Psychological bulletin*, vol. 56, p. 81, 1959.
- [30] J. Recker, *Scientific research in information systems: a beginner's guide*: Springer Science & Business Media, 2012.
- [31] M. Alshehri, S. Drew, and R. AlGhamdi, "Analysis of Citizens Acceptance for E-government Services: Applying the UTAUT Model," *arXiv preprint arXiv:1304.3157*, 2013.
- [32] S. S. Al-Gahtani, G. S. Hubona, and J. Wang, "Information technology (IT) in Saudi Arabia: Culture and the acceptance and use of IT," *Information & Management*, vol. 44, pp. 681-691, 2007.
- [33] C. Fornell and D. F. Larcker, "Evaluating structural equation models with unobservable variables and measurement error," *Journal of marketing research*, pp. 39-50, 1981.
- [34] N. Kock, "Using WarpPLS in E-Collaboration Studies: Descriptive Statistics, Settings," *Interdisciplinary Applications of Electronic Collaboration Approaches and Technologies*, vol. 62, 2013.

Control of Industrial Systems to Avoid Failures: Application to Electrical System

Yamen EL TOUATI

Department of Computer Science
Faculty of Computing and Information Technology
Northern Border University
Kingdom of Saudi Arabia

Saleh ALTOWAIJRI

Department of Information Systems
Faculty of Computing and Information Technology
Northern Border University
Kingdom of Saudi Arabia

Mohamed AYARI

Department of Information Technology
Faculty of Computing and Information Technology
Northern Border University
Kingdom of Saudi Arabia

Abstract—We resolve the control problem for a class of dynamic hybrid systems (DHS) considering electrical systems as case study. The objective is to guarantee that the plan never reaches unsafe states. We consider a subclass class of DHS called Cumulative Preemptive Event-driven DHS (CPE-DHS). This class is distinguished by the dominance of its discrete aspect characterized by features as cumulative continuous variables combined with actions behavior that may be interrupted and restarted. We utilize a subclass of Rectangular Hybrid Automata (RHA), named Constant Slope RHA (CSRHA), as a solution framework to resolve the control problem. The main contribution is a control Algorithm for the class of systems described above. This algorithm ensures that the system meet the requirement specifications by forcing some events. The forcing action is given in the form of restrictions on the transition guards of the CSRHA. The termination/decidability as well as correctness of the algorithm is given by theorems and formal proofs. This contribution ensures that the system will always be safe states and avoid failure due to the reachability of unsafe states. Our approach can be applied to a large category of industrial systems, especially electrical systems that we consider as case study.

Keywords—Dynamic hybrid systems; supervisory control; hybrid automata; electrical systems; safety

I. INTRODUCTION

Dynamic hybrid systems [1]–[4] (DHS) are systems characterized by the interaction of both discrete and continuous components. A large variety of real-time and embedded systems and many computer automated systems as well as industrial and electrical systems are described by both continuous and discrete aspects. In this paper, we concentrate in a particular class of dynamic hybrid systems where system behavior is captured essentially by preemptive activities which can be produced sequentially or in parallel. Besides, these systems are depicted by an interaction of dominant discrete component with a slight continuous one.

DHS are modeled by a large variety of modeling frameworks. We distinguish essentially several timed and hybrid extensions of finite state automata [5] as well as Petri nets [6], [7]. Petri nets extensions benefit a salient graphical modeling power. However, computations are mostly based on similar

automata extension. On the other hand, there are many extensions of finite state machines, such as time transition systems [8], timed automata [5] and stop watch automata [9]. In these frameworks, time is included in configurations and transitions in the form of constraints and/or speed rate. In order to deal with dynamic hybrid systems, we consider essentially hybrid automata, linear hybrid automata, and rectangular automata [10]. All the previous frameworks capture various aspects of DHS depending on their modelling power which is generally inversely proportional with the decidability of the accessibility problem. In fact, models that cover more classes of systems become more difficult to manage by a computer due to the undecidability problems [11].

In our case, we use a subclass of RHA: the CSRHA to model our systems. This subclass is better managed from decidability side. The control problem, as one of the highly studied problems in literature [12], will be resolved using CSRHA formalism. One of the important problems in the DHS control theory is related to safety verification. This problem states that the controller has to ensure that all the trajectories of the system do not reach any “unsafe” state. In order to guarantee this safety property, the controller may restrict the scope of some controllable events. By taking such decision, the controller avoids that system trajectories interfere with any undesired state induced by uncontrollable events. However, in this paper, we consider that the computational power of the controller is limited to narrowing time intervals on transitions related to controllable events. Technically speaking, this action is similar to modifying guards on transitions associated to controllable events in the CSRHA model.

This paper is organized as follows. The next section provides background of hybrid automata and a description of the CSRHA. In Section 3, we present and solve the supervisory control problem. We note that throughout this paper we use the same case study of an electrical system to illustrate our supervisory control approach.

II. BACKGROUND ON HYBRID AUTOMATA

In the following, we define the retained subclass of RHA: the CSRHA.

A. Constant Piece-wise Rectangular Hybrid Automata

We consider these notations: $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ is a finite set of real valued clocks (variables). $\dot{\mathcal{X}} = \{\dot{x}, x \in \mathcal{X}\}$ denotes the set of first derivative variables of \mathcal{X} . A variable x is considered piece-wise linear variable if $\dot{x} \in \mathbb{R}$. \sim denotes an element of operator's set $\{<, \leq, =, \geq, >, \neq\}$. A rectangular inequality over \mathcal{X} , is an inequality of the form, $x \sim c$, where $c \in \mathbb{R}$, and $x \in \mathcal{X}$. A rectangular predicate over \mathcal{X} is a conjunction of rectangular inequalities over \mathcal{X} . $Rect(\mathcal{X})$ denotes the set of rectangular predicates over \mathcal{X} . A polyhedral inequality over \mathcal{X} is an inequality of the form $c_1x_1 + \dots + c_kx_k \sim c$, where $c, c_1, \dots, c_k \in \mathbb{R}$, and $x_1, \dots, x_k \in \mathcal{X}$. A polyhedral predicate over \mathcal{X} is boolean combination of polyhedral inequalities over \mathcal{X} . $\Psi(\mathcal{X})$ is the set of polyhedral predicates over \mathcal{X} . $\mathbf{v} = (v_1, \dots, v_n)$ denotes an element of \mathbb{R}^n that captures clocks valuation, $v_i \in \mathbb{R}$, of every clock $x_i \in \mathcal{X}$. $v(x_i) = v_i$ corresponds to the value of x_i . We denote by *region*, a subset of \mathbb{R}^n . For a region z and $x_i \in \mathcal{X}$, $z(x_i) = \{v_i | \mathbf{v} \in z\}$. $\psi(\mathbf{v})$ denotes the boolean function which equals **true** if the predicate ψ is satisfied by the input vector \mathbf{v} and **false** if not. We denote by $\llbracket \psi \rrbracket$, the region composed by the set of vectors $\mathbf{v} \in \mathbb{R}^n$, where the predicate ψ is **true** when we substitute each x_i by its corresponding v_i . $\llbracket \psi \rrbracket(x_i)$ denotes the interval of values captured by v_i , $\forall \mathbf{v} \in \llbracket \psi \rrbracket$.

Definition 1: In [13]–[15] A constant piece-wise linear hybrid automata (CSRHA) is a tuple $\mathcal{A} = (\mathcal{X}, \mathcal{Q}, \mathcal{T} \cup \{e_0\}, inv, dyn, guard, assign, l_0)$ where:

- \mathcal{X} , is a finite set of variables.
- \mathcal{Q} , is a finite set of locations.
- \mathcal{T} , is a finite set of transitions. A transition $e = (l, l') \in \mathcal{T}$, leads the system from the source location $l \in \mathcal{Q}$, to the end location $l' \in \mathcal{Q}$. The entry transition of the initial state l_0 is denoted by e_0 .
- *inv*: $\mathcal{Q} \rightarrow \Psi(\mathcal{X})$ is the location *invariant*, it associates a predicate to each location.
- *dyn*: $\mathcal{Q} \times \mathcal{X} \rightarrow \mathbb{R}$, is a function describing the evolution of variables. This evolution is usually of the form $l, \dot{x} = k$, $k \in \mathbb{R}$ or simply $\dot{x} = k$ in the location l . $\dot{\mathcal{X}}(l)$ denotes the evolution of all variables in the location l .
- *guard*: $\mathcal{T} \rightarrow \Psi(\mathcal{X})$ is the guard function. It associates a predicate, C_e to each transition, e . The guard, C_e should equals true to allow the execution of the transition e .
- *assign*, is the initialization function. It associates a relation, $assign_e$ to each transition e defining the clocks to be reset.
- $l_0 \in \mathcal{Q}$, is the initial location. □

The semantic of a constant piece-wise linear hybrid automata (CSRHA) is given by the following definition:

Definition 2: The semantic of a CSRHA $\mathcal{A} = (\mathcal{X}, \mathcal{Q}, \mathcal{T} \cup \{e_0\}, inv, dyn, guard, assign, l_0)$ is defined by a timed transition system $S_{\mathcal{A}} = (\mathcal{Q}, q_0, \rightarrow)$ with

- $\mathcal{Q} = \mathcal{Q} \times \mathbb{R}^n$ with $n = |\mathcal{X}|$.
- $q_0 = (l_0, init)$ is the initial state.
- $\rightarrow \in (\mathcal{Q} \times (\mathcal{T} \cup \mathbb{R}_+) \times \mathcal{Q})$ is defined by:
 - $(l, v) \xrightarrow{a} (l', v')$ (**jump transition**) if $\exists e = (l, l') \in \mathcal{A}$ s.t.

$$\begin{cases} a = e \\ guard(e)(v) = true \\ v' = assign_e(v) \\ inv(l')(v') = true \end{cases}$$
 - $(l, v) \xrightarrow{\epsilon(t)} (l', v')$ (**flow transition**) if

$$\begin{cases} l = l' \\ v' = v + t * \dot{\mathcal{X}}(l) \\ inv(l')(v') = true \end{cases} \quad \square$$

A *run* of CSRHA \mathcal{A} is a path in $S_{\mathcal{A}}$ started from q_0 . $\llbracket \mathcal{A} \rrbracket$ denotes the set of all runs of \mathcal{A} . We note $(l, v) \xrightarrow{\epsilon(t)} (l', v') \xrightarrow{a} (l'', v'')$ is equivalent to $(l, v) \xrightarrow{a} (l'', v'')$. A state (l_i, v_i) is considered as *reachable*, if $\exists (l_0, v_0) \xrightarrow{\epsilon(t_0)} (l_1, v_1) \xrightarrow{\epsilon(t_1)} (l_2, v_2) \xrightarrow{\epsilon(t_2)} \dots \xrightarrow{\epsilon(t_i)} (l_i, v_i)$ where $(l_0, v_0) = q_0$. A run $(l_0, v_0) \xrightarrow{\epsilon(t_0)} (l_1, v_1) \xrightarrow{\epsilon(t_1)} (l_2, v_2) \xrightarrow{\epsilon(t_2)} \dots \xrightarrow{\epsilon(t_i)} (l_i, v_i) \dots$ starting from $q_0 = (l_0, v_0)$ is a *timed trace*, denoted as $w = (a_0 = e_0, \delta_0) \rightarrow (a_1, \delta_1) \rightarrow (a_2, \delta_2) \rightarrow \dots (a_2, \delta_2) \dots$, where w is a sequence of pairs (a_i, δ_i) , with $a_i \in \mathcal{T} \cup \{e_0\}$ a transition, and $\delta_{i+1} \in \mathcal{R}_+$ is the delay between the two successive events a_i and a_{i+1} , where $\delta_0 = 0$, and $\forall i \geq 1, \delta_i = \epsilon(t_i) - \epsilon(t_{i-1})$.

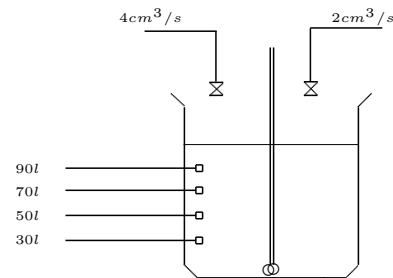


Fig. 1. Electrical system for blending chemical solution.

Example 2.1: Consider the electrical system for mixing chemical solution given in Fig. 1. Filling action is composed of two stages. Firstly, a tray is replenished by a chemical solution with a rate of $2cm^3/s$. We assume that initially the tray is filled by $10dm^3$ of a neutral liquid. This phase is accomplished when the current content of the tray is bounded by 30 and 50 liters. The next phase should be fulfilled before a deadline of 18s, elapses in order to avoid the risk of obtaining improper solution. An authorization at a random time prompts the second stage which has a deadline of 16s once started.

When the next stage is activated, a second chemical solution is replenished with the rate of $4cm^3/s$. The filling process is accomplished when the total content of the tray is bounded by $70dm^3$ and $90dm^3$. The CSRHA modeling of this electrical system is illustrated in Fig. 2.

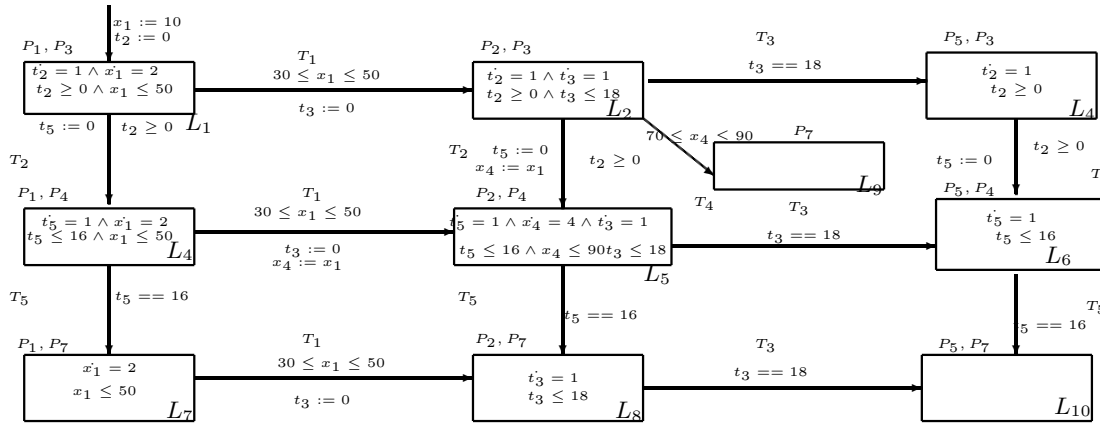


Fig. 2. The CSRHA of the electrical system.

III. CONTROL OF CPE-DHS

In the following, we describe our contribution to resolve the control problem. Our solution define a derived space where all trajectories satisfy the requested specifications to avoid system failure. Thus, all unsafe locations will be inaccessible. The safety specification is considered as the set of forbidden locations. The control action acts by reducing transition guard intervals. By nature, some events are not eligible for narrowing their time occurrence scope. Such events are considered uncontrollable from the controller perspective. An event is controllable if the controller has the power to reduce its occurrence time slot. In general, event connected to forbidden locations are uncontrollable, otherwise it becomes trivial to define the control solution. Moreover, the restriction action on the time intervals should be minimalist.

A. Specification of the Control Problem

The inputs are the set of unsafe locations and the partition of events as controllable/uncontrollable. The main steps that we propose to resolve the control problem are as follows:

Steps:

- 1) Mark all unsafe locations considering the safety specification.
- 2) Mark all transitions as controllable and uncontrollable considering the input events partition.
- 3) Perform a computation of the desired space adopted by the controller in all the locations to ensure that the system is not accessing forbidden locations.
- 4) Reassign the restricted guards of transition related to controllable events and update any necessary location invariant to force that the system remains in safe states.

B. Control Algorithm

Let $\mathcal{A} = (L, l_0, X, \Sigma, E, inv, Dif)$ the CSRHA model of the system to be controlled. \mathcal{A}^d represents the output (controlled) CSRHA. We consider these notations:

- L^F represents the set of forbidden locations (given by the safety specification).

- E^F represents the set of CSRHA transitions where the output location is a forbidden.

$$E^F = \{e \in E | e = (l, \delta, \alpha, Aff, \rho, l'), l' \in L^F\}$$

- $e_{l,l'}$ represents a transition $e = (l, \delta, \alpha, Aff, \rho, l')$ where the source location is l and the destination location is l' .
- E_l represents the set of transitions having l as source location.

$$E_l = \{e \in E | e = (l, \delta, \alpha, Aff, \rho, l'), l' \in L\}$$

- $E_l^F = E_l \cap E^F$ represents the forbidden transitions having l as source location.
- $\overline{E_l^F} = E_l - E^F$ represents the non forbidden transitions having l as source location.

- In Algorithm III.1, we consider that $|E^l| = |\overline{E_l^F}| \cup |E_l^F| = m$ such as $|\overline{E_l^F}| = k$, $|E_l^F| = m - k$. E^l is parted into $\overline{E_l^F}$ and E_l^F as follows:

$$\begin{aligned} \circ \overline{E_l^F} &= \{e_1 = (l, \delta_1, \alpha_1, Aff_1, \rho_1, l_1), e_2 = (l, \delta_2, \alpha_2, Aff_2, \rho_2, l_2), \dots, e_k = (l, \delta_k, \alpha_k, Aff_k, \rho_k, l_k)\} \\ \circ E_l^F &= \{e_{k+1} = (l, \delta_{k+1}, \alpha_{k+1}, Aff_{k+1}, \rho_{k+1}, l_{k+1}), e_{k+2} = (l, \delta_{k+2}, \alpha_{k+2}, Aff_{k+2}, \rho_{k+2}, l_{k+2}), \dots, e_m = (l, \delta_m, \alpha_m, Aff_m, \rho_m, l_m)\} \end{aligned}$$

- $L^R \in (2^L)^L$ represents the set of reachable locations from l in $-\mathcal{A}^l$. In other words, a location $l' \in L^R(l)$ if it exists a run from l' to l . Formally, $L^R(l) = \{l \in L, \exists k \in \mathbb{N}, (l', v') \xrightarrow{t_1, e_1} (l_1, v_1) \xrightarrow{t_2, e_2} (l_2, v_2) \dots \xrightarrow{t_k, e_k} (l_k, v_k), l_k = l\}$. This corresponds to the closure of the set $\{l\}$ under the relation $\{(p, q) : \text{there is a transition } e = (p, \delta, \alpha, Aff, \rho, q) \in E, q \in L^R(l)\}$.

- $\overline{L^R(l)}$ represents $L - L^R(l)$.

Algorithm III.1 Control Algorithm

¹ $-\mathcal{A}$ is the reversed automata of \mathcal{A} ([16]).

```

1: function Control( $\mathcal{A}, M^F$ ): $\mathcal{A}^d$ 
2: initialize the output CSRHA by the entry CSRHA.  $\mathcal{A}^d := \mathcal{A}$ 


---


3: function initialize()
4: calculate the set  $E^F$  :
5: for all  $e_{l,l'} \in E$  with  $l' \in L^F$  do
6:    $E^F := E^F \cup \{e_{l,l'}\}$ 
7: end for
8: calculate  $L^R(l)$ 
9: initialize  $L^R(l) := \{l\}$ 
10: while  $\exists e = (l', \delta, \alpha, Aff, \rho, l'') \in E$ , with  $l'' \in L^R(l)$  and  $l' \notin L^R(l)$  do
11:    $L^R(l) := L^R(l) \cup \{l'\}$ .
12: end while
13: for all location  $l \in L \setminus L^F$  do
14:   calculate  $E_l^F$  and  $\overline{E}_l^F$ :
15:   for all  $e_{l,l'} \in E^F, l' \in L$  do
16:      $E_l^F := E_l^F \cup \{e_{l,l'}\}$ 
17:   end for
18:   calculate  $\overline{E}_l^F := E_l - E_l^F$ 
19:   if  $\overline{E}_l^F = \emptyset$  then
20:      $L^F := L^F \cup \{l\}$ 
21:   end if
22: end for
23: if  $L^F$  is modified then
24:   re-invoke initialize() to reconsider the new forbidden locations
25: end if
26: if  $L^F = L$  then
27:   exit (there is no solution)
28: end if
29: end function


---


30: for all location  $l$  where  $E_l^F \neq \emptyset$  do
31:   recalculate the guard predicates of all the transitions included into the set  $\overline{E}_l^F$  :
32:   for all  $e_i \in \overline{E}_l^F, i \in \llbracket k+1, m \rrbracket$  do
33:     calculate the new guard  $\delta_i^n$  regarding  $\delta_i$  and guards of transitions in  $E_l^F$  :2

$$\delta_i^n := \delta_i \wedge \neg \delta_1 \wedge \neg \delta_2 \dots \wedge \neg \delta_k$$

34:   end for
35: end for
36: do a forward analysis, started at the initial location. We note by  $S_l^{forward}$  the reachable space3 calculated by forward analysis at location  $l$ .
37: for all location  $l$  where  $E_l^F \neq \emptyset$  do
38:   do a backward analysis started at location  $l$  considering  $\delta_{k+1}^n \vee \delta_{k+2}^n \vee \dots \vee \delta_m^n$  as initial entry space. We note by  $S_{l,l'}^{backward}$  the space calculated by backward analysis (from location  $l'$ ) in the location  $l \in L_l^R$ 
39: end for
40: for all  $l' \in L^R(l)$  where  $E_{l'}^F \neq \emptyset$  do
41:   calculate the final space of backward analysis at loca-

```

²Our goal is to reduce the state space in order to avoid the possibility of occurrence of prohibited events.

³The reachable space at a given location is a polyhedron with dimension $|X|$ defined the inequalities system $A.XR \leq b$, with $A \in \mathcal{M}_{a,|X|}(\mathbb{R})$ a matrix with a lines and $|X|$ columns, and $X \in \mathbb{R}^n$ the vector of CSRHA variables.

tion l' :

$$S_{l'}^{backward} := \bigwedge_{l \in E_l^F} S_{l,l'}^{backward}$$

```

42: end for
43: for all location  $l_i$  do
44:   calculate the desired space  $S_l^d$  at location  $l$ 

$$S_l^d = S_l^{backward} \wedge S_l^{forward}$$

45:   calculate the new location invariant  $l$  given by

$$inv^d(l) := inv(l) \wedge S_l^d$$

46:   for all transition  $e_{l,l'} \in E_l^F$  do
47:     redefine the guards :  $\delta_{l,l'}^d := \delta_{l,l'} \wedge S_l^d$ 
48:   end for
49: end for
50: end function


---



```

The CSRHA modelling a CPE-DHS system is the input of the Algorithm III.1. Algorithm III.1 produces the output as an updated CSRHA where forbidden states can never be reached. The control algorithm computes the new transition guards and the new location invariants.

Theorem 1: The Algorithm III.1 terminates if the entry CSRHA has no loop.

Proof: 1 The Algorithm III.1 terminates if the computation of reachable space (both backward and forward) terminates. This analysis use discrete and continuous predecessor and successor operators which perform certain geometric calculus on regions [14]. Software like PHAVER [17] and SpaceEx [18], [19] implement such region operations, using polyhedral libraries, to accomplish the reachable space computation. We note that these analysis terminate if the CSRHA is acyclic. Nevertheless, for more general forms, the accessibility problem is known as undecidable [14], [20]. ■

In the following, we present some particular and interesting cases where this problem is decidable.

Theorem 2: The Algorithm III.1 terminates if the input CSRHA satisfies the following proprieties:

- 1) All derivative variables in the locations are **non negative** or **null**.
- 2) Guards and invariants are defined by **single non negative** constraints.
- 3) Assignments are of the form $x' := x$ or $x' = c$.

Proof: This is ensured due to the decidability of accessibility problems in that case [21]. ■

Furthermore, we can ensure the algorithm decidability for these interesting classes of CSRHA:

- 1) CSRHA where each loop contains at least one initialization of all clocks [22].
- 2) CSRHA where each loop contains at most one transition guard in the form of “dangerous” test [22].
- 3) CSRHA where the dynamic changing (the derivative value) of a variable between two locations is accompanied by resetting the variable assignment at the transition between the two locations [16].

Theorem 3: The automaton \mathcal{A}^d obtained by applying Algorithm III.1 ensures that all reachable spaces respect the safety specification while being maximal permissive.

Proof: Consider the CSRHA $\mathcal{A} = (L, l_0, X, \Sigma, E, inv, Dif)$.

Part 1: We demonstrate (by contradiction) that the reachable space meets the safety specification.

Suppose that $\exists l \in L^F$ such as it exists a run in \mathcal{A}^d from initial state:

$$(l_0, v_0) \xrightarrow{t_0, e_0} (l_1, v_1) \dots (l_a, v_a) \xrightarrow{t_a, e_a} (l, v)$$

We have $l \in L^F \implies e_a \in E^F$. Suppose that $e_a = (l_a, \delta_a, \alpha_a, Aff_a, \rho_a, l)$. According to the TTS of \mathcal{A}^d , we have $inv(l)(v_a) = true$ and $\delta(v_a) = true$. However, according to Algorithm III.1, the calculation of $S_l^{backward}$ conclude that $inv^d(l) = inv(l) \wedge \neg\delta_1 \wedge \neg\delta_2 \dots \wedge \neg\delta_k, \forall e_i \in E_{l_a}^F, i \in \llbracket 1, k \rrbracket$. According to the construction of the set $E_{l_a}^F$ in the Algorithm, we have $e_a \in E_{l_a}^F$. Thus, $\exists j \in \llbracket 1, k \rrbracket$ such as $e_a = e_j$. This implies that $inv(l_j^a)(v) = false$, which contradicts the starting assumption.

Part 2: We demonstrate (by contradiction) that the reachable space at \mathcal{A}^d is maximal permissive.

To do this, let us suppose that there is a location $(l, v) \in \mathcal{Q}_A$ such that $(l, v) \notin \mathcal{Q}_{A^d}$ and $l \notin L^F$. Also suppose that (l, v) do not lead to forbidden locations by the specification. As $(l, v) \notin \mathcal{Q}_{A^d}$, we obtain $inv^d(l)(v) = false$. Similarly $(l, v) \in \mathcal{Q}_A \implies inv^d(l)(v) = true$.

The fact that (l, v) does not lead to unauthorized locations, means that there is no run from (l, v) leading to a state (l', v') with $l' \in L^F$.

Let $l^f \in L^F$ a location such that $l \in L^R(l_f)$. Since there is no run from (l, v) leading to forbidden location, thus, (l_f, v_f) is not reachable since (l, v) , and that, for any $v_f \in \mathbb{R}^{|X|}$. Similarly, (l, v) is not reachable from (l_f, v_f) at the reverse automaton $-\mathcal{A}$ (or by backward analysis). Let $S_{l, l_f}^{backward}$ the obtained space at l by backward analysis from (l_f, v_f) . Thus, we have $S_{l, l_f}^{backward}(v) = false$.

According to Algorithm III.1, S_l^d start by the initial space $\neg\delta_1 \wedge \neg\delta_2 \dots \wedge \neg\delta_k \forall e_i \in E_{l_f}^F, i \in \llbracket 1, k \rrbracket, k = |E_{l_f}^F|$. Thus, $S_l^d(v) = true$. Moreover, according to the calculation formula of location invariant, we have $inv^d(l)(v) = true$.

$\implies (l, v) \in \mathcal{Q}^d$. Thus, any location leading exclusively to locations respecting the specification is in the reachable space of \mathcal{A}^d . Consequently, \mathcal{A}^d is maximal permissive. ■

Example 3.1: We reconsider the CSRHA of the electrical system illustrated in Fig. 2. According to the safety specification, we consider the following unsafe locations: $\mathcal{SF} = \{l_7, l_8, l_{10}, l_4, l_6\}$. The results related to the reachable space computation by forward and backward analysis are performed by PHAVer [17] and SpaceEx [18], [19] software. The intersection between backward and forward spaces is illustrated in Table I. The results meets with the safety specification. Thus, the controller defines a derived CSRHA where invariant locations and transition guards are truncated by the new obtained polyhedral equations in each location. This derived

TABLE I. INTERSECTION SPACE

l_5	$E_5^c = (-x_4 - 4t_2 + 4t_3 \geq -130) \wedge (-x_4 \geq -90) \wedge (-x_4 + 4t_3 \geq -50) \wedge (-x_4 + 2t_2 + 4t_3 > -42) \wedge (x_4 - 4t_3 > -2) \wedge (-t_5 > -16) \wedge (t_3 \geq 0) \wedge (t_2 \geq 0) \wedge (x_4 - 4t_5 > 6) \wedge (x_4 - 2t_2 + 2t_3 \geq 10) \wedge (x_4 \geq 30) \wedge (7x_4 + 18t_2 - 28t_3 \geq 210) \wedge (-t_3 > -18)$
l_4	$E_4^c = (x_1 - 2t_2 - 2t_5 == 10) \wedge (-x_1 \geq -50) \wedge (-x_1 + 2t_2 > -42) \wedge (t_2 \geq 0) \wedge (x_1 - 2t_2 \geq 10)$
l_1	$E_1^c = (x_1 - 2t_2 == 10) \wedge (-x_1 \geq -50) \wedge (x_1 \geq 10)$
l_9	$E_9^c = (-x_4 - 4t_2 + 4t_3 \geq -130) \wedge (-x_4 \geq -90) \wedge (-x_4 + 4t_3 \geq -50) \wedge (-x_4 + 2t_2 + 4t_3 > -42) \wedge (-t_3 > -18) \wedge (-t_5 > -16) \wedge (t_2 \geq 0) \wedge (x_4 \geq 70) \wedge (7x_4 + 18t_2 - 28t_3 \geq 210)$
l_2	$E_2^c = (x_1 - 2t_2 + 2t_3 == 10) \wedge (-x_1 \geq -50) \wedge (-x_1 + 2t_2 \geq -10) \wedge (3x_1 - 4t_2 > 18) \wedge (x_1 \geq 30)$

automaton is maximal permissive and describes all possible trajectories that obey to the requirements.

Table I illustrates the intersection space, obtained by PHAVer and SpaceEx. This allows capturing the maximal polyhedron that meet with requirements. For example, the updated location invariant of l_4 is given by $I_4^c = I_4 \wedge (x_1 - 2t_2 - 2t_5 == 10) \wedge (-x_1 \geq -50) \wedge (-x_1 + 2t_2 > -42) \wedge (t_2 \geq 0) \wedge (x_1 - 2t_2 \geq 10)$. Besides, the updated guard of transition of $T_{4,5}$ is $g_{4,5}^c = g_{4,5} \wedge t_5 < 16 \wedge (x_1 - 2t_2 - 2t_5 == 10) \wedge (-x_1 \geq -50) \wedge (-x_1 + 2t_2 > -42) \wedge (t_2 \geq 0) \wedge (x_1 - 2t_2 \geq 10)$. Similarly, all guards and invariants will be updated according to the results given by the intersection space. Furthermore, we omit any outgoing transition from a forbidden location (since it becomes unreachable).

IV. CONCLUSION

In this paper, our main contribution is to solve the problem of supervisory control of the particular class of dynamic hybrid systems (DHS) called Cumulative Preemptive Event-driven DHS (CPE-DHS) by narrowing guards and invariants of transitions relative to controllable events in a way that forbidden states remain inaccessible. Our proposed solution can be applied in a systematic way to any system that fits with our requirements. Then we applied this approach to an electrical system as case study. Generally speaking, the control problem is known to be undecidable for this class of complex systems. Nevertheless, in quest of decidability, we propose some restrictions that makes the problem decidable. In our future directions, we will focus on the supervisor generation while considering uncontrollable variables.

ACKNOWLEDGMENT

The authors gratefully acknowledge the approval and the support of this research from the Deanship of Scientific Research study by the grant no 4665-CIT-2016-1-6-F K.S.A, Northern Border University, Arar, KSA.

REFERENCES

- [1] A. van der Schaft and H. Schumacher, *An introduction to hybrid dynamical systems Lecture Notes in Control and Information Sciences*. London: Springer-Verlag London Ltd., 2000, vol. 251.
- [2] J. Lygeros, K. Johansson, S. Simic, J. Zhang, and S. Sastry, "Dynamical properties of hybrid automata," *IEEE Transactions on Automatic Control*, vol. 48, no. 1, pp. 2–17, 2003.

- [3] R. Sanfelice., "Robust hybrid control systems," Ph.D. dissertation, University of California, Santa Barbara, 2007.
- [4] M. Kosmykov, "Hybrid dynamics in large-scale logistics networks," Ph.D. dissertation, University of Bremen, 2011.
- [5] R. Alur, "Timed automata," *Theoretical Computer Science*, vol. 126, pp. 183–235, 1999.
- [6] B. Berthomieu and M. Diaz, "Modeling and verification of time dependent systems using time Petri nets," *IEEE transactions on software Engineering*, vol. 17, no. 3, pp. 259–273, 1991.
- [7] R. David and H. Alla, *Discrete, Continuous, and Hybrid Petri Nets*, 1st ed. Springer, Heidelberg, 2004.
- [8] T. A. Henzinger, Z. Manna, and A. Pnueli, "Timed transition systems," in *Real-Time: Theory in Practice*, ser. Lecture Notes in Computer Science, J. Bakker, C. Huizing, W. Roever, and G. Rozenberg, Eds. Springer Berlin Heidelberg, 1992, vol. 600, pp. 226–251. [Online]. Available: <http://dx.doi.org/10.1007/BFb0031995>
- [9] F. Casse and K. Larsen, "The impressive power of stopwatches," in *In Proc. of CONCUR 2000: Concurrency Theory*. Springer, 1999, pp. 138–152.
- [10] T. A. Henzinger, P. W. Kopke, A. Puri, and P. Varaiya, "What's decidable about hybrid automata?" in *STOC '95: Proceedings of the twenty-seventh annual ACM symposium on Theory of computing*. New York, NY, USA: ACM, 1995, pp. 373–382, a revised version appeared in *Journal of Computer and System Sciences*, vol. 57, p. 94124, 1998.
- [11] E. Asarin, V. P. Mysore, A. Pnueli, and G. Schneider, "Low dimensional hybrid systems - decidable, undecidable, don't know," *Inf. Comput.*, vol. 211, pp. 138–159, Feb. 2012. [Online]. Available: <http://dx.doi.org/10.1016/j.ic.2011.11.006>
- [12] R. Alur, T. Henzinger, and P.-H. Ho, "Automatic symbolic verification of embedded systems," *IEEE Transactions on Software Engineering*, vol. 22, no. 3, pp. 181–201, 1996, (A preliminary version appeared in the Proceedings of the 14th Annual Real-Time Systems Symposium (RTSS), IEEE Computer Society Press, 1993, pp. 2-11.).
- [13] T. A. Henzinger, X. Nicollin, J. Sifakis, and S. Yovine, "Symbolic model checking for real-time systems," *Information and Computation*, vol. 111, pp. 394–406, 1994. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.26.7422>
- [14] R. Alur, C. Courcoubetis, N. Halbwachs, T. A. Henzinger, P.-H. Ho, X. Nicollin, A. Olivero, J. Sifakis, and S. Yovine, "The algorithmic analysis of hybrid systems," *Theoretical Computer Science*, vol. 138, no. 1, pp. 3–34, 1995, (A preliminary version appeared in the Proceedings of the 11th International Conference on Analysis and Optimization of Systems: Discrete-Event Systems (ICAOS), Lecture Notes in Control and Information Sciences 199, Springer-Verlag, 1994, pp. 331-351.). [Online]. Available: citeseer.ist.psu.edu/alur95algorithmic.html
- [15] Y. E. Touati, N. B. Hadj-Alouane, and M. Yeddes, "Modeling and control of constant speed dynamic hybrid systems using extended time petri networks," in *Control and Decision Conference (CCDC), 2012 24th Chinese*, May 2012, pp. 634–641.
- [16] T. A. Henzinger and V. Rusu, "Reachability verification for hybrid automata," in *HSCC 98: Hybrid Systems Computation and Control, Lecture Notes in Computer Science 1386*. Springer-Verlag, 1998, pp. 190–204.
- [17] G. Frehse, "Compositional verification of hybrid systems using simulation relations," Ph.D. dissertation, Radboud Universiteit Nijmegen, 10 2005.
- [18] G. F. Alexandre Donzé, "Modular, hierarchical models of control systems in spaceex," in *Proc. European Control Conf. (ECC'13)*, Zurich, Switzerland, 2013.
- [19] G. Frehse, R. Kateja, and C. Le Guernic, "Flowpipe approximation and clustering in space-time," in *Proceedings of the 16th International Conference on Hybrid Systems: Computation and Control*, ser. HSCC '13. New York, NY, USA: ACM, 2013, pp. 203–212. [Online]. Available: <http://doi.acm.org/10.1145/2461328.2461361>
- [20] R. Alur, C. Courcoubetis, T. Henzinger, and P. Ho, "Hybrid automata: An algorithmic approach to the specification and verification of hybrid systems," in *Hybrid Systems*, ser. Lecture Notes in Computer Science, R. Grossman, A. Nerode, A. Ravn, and H. Rischel, Eds. Springer Berlin, Heidelberg, 1993, vol. 736, pp. 209–229, 10.1007/3-540-57318-6-30. [Online]. Available: <http://dx.doi.org/10.1007/3-540-57318-6-30>
- [21] Y. Kesten, A. Pnueli, J. Sifakis, and S. Yovine, "Integration graphs: A class of decidable hybrid systems," in *Hybrid Systems, volume 736 of Lecture Notes in Computer Science*. Springer-Verlag, 1993, pp. 179–208, (appeared in In Proceedings of Workshop on Theory of Hybrid Systems , Lyngby, Denmark, June 1992).
- [22] Y. Kesten, A. Pnueli, J. Sifakis, and S. Yovine, "Decidable integration graphs," in *Information and Computation, volume 150(2)*, 1999, pp. 209–243.

Mobility Management using the IP Protocol

Imtiaz A. Halepoto¹, Adnan Manzoor², Nazar H. Phulpoto², Sohail A. Memon³, Muzamil Hussain⁴

¹Department of Computer Systems Engineering, QUEST Nawabshah, Pakistan

²Department of Information Technology, QUEST Nawabshah, Pakistan

³Department of Mathematics, SALU Khairpur, Pakistan

⁴Department of Computer Science, QUEST Nawabshah, Pakistan

Abstract—Time critical applications, such as VoIP and video conferencing require Internet connectivity all of the time for better performance. Moreover, in case of vehicular networks, it is very common for mobile devices to move from one network to another. In such scenarios the sudden changes in the network connectivity may cause problems, which affects the data transmission rate. The movement of a mobile node from one network to another is also a challenge for the routers to maintain the routing information as well as to forward the data to the corresponding node. In all of the aforementioned scenarios, the switching between the networks with minimum latency improves the performance, i.e. in terms of mobility and availability of the network. The Mobile IP protocol serves the purpose of seamless handover of mobile devices from one network to another. A mobile node maintains its permanent IP address using the Mobile IP protocol while moving to a foreign network. When a mobile node establishes the connection with the foreign network the data packets transmitted from the home network are redirected to the foreign network. The Mobile IP protocol establishes a tunnel between the home network and the foreign network. The process of tunneling continues until the mobile node moves back to the home network or when the foreign network advertises the new IP address of the mobile node. With the increasing number of wireless devices the mobility is the key challenge. The devices with multiple interfaces such as mobile phone which uses 4G as well as WiFi, the urge for the availability of the Internet is also high. This paper provides a deep discussion about the Mobile IP protocol and its implementation. A network scenario is proposed with the configuration of the Mobile IP. According to the obtained results of the simulations, the Mobile IP protocol increases the availability of the network connection as well as it achieves the larger throughput when compared with the scenario without using the Mobile IP.

Keywords—Mobility; Mobile IP; 4G; foreign network; permanent IP

I. INTRODUCTION

The main protocol that is most widely used over the Internet is TCP [1]. However, TCP needs the services of the network layer protocols. The famous network layer protocol is the IP protocol. At the time of IP proposal, it was intended for the stationary networks with a constant IP address assigned to each host. With the recent advances in the technology like ubiquitous computing and pervasive computing the mobility of devices while maintaining the Internet connection have become the normal trend. As the IP is mainly for stationary devices so in case of mobility [12] when a device wanders far, it disconnect with the home network. For example, when a device is outside wireless range of the service provider. In case of cellular networks for example 4G [11], the technique like link layer and roaming minimize the handover [2]. However, many of of companies still use the mobile IP. In many cases

some of the handover management [16] schemes are used to reduce the transmission latency when a device travels from the home network to a foreign network. However, the latency in the handover management schemes is not according the expectations and also raise many questions regarding the availability of the network.

In order to cope with mobility and the availability of the Internet (or network connection) some of the protocols are proposed at the transport layer, such as the stream control transmission protocol (SCTP). With SCTP, it is possible to remain connected with multiple networks at the same time. SCTP increases the Internet availability as well as provides features to enhance the mobility. In order to facilitate the mobility, SCTP uses the multihoming technique [3], which provides the ability to a device to remain connected with multiple networks (i.e. home and foreign agents). However, there are many issues in SCTP such as congestion control, flow control and scheduling of data transmission [4]–[7]. In Fig. 1, a network is shown with two multihomed devices. Device A with two interface and device B with three interfaces connected by using the SCTP protocol. The other solution to the mobility and the availability is the use of Mobile IP protocol [8]–[10]. The assortment of Mobile IP with the other mobility techniques is beneficial. With Mobile IP it is allowed to devices to use the permanent IP address while moving to a foreign network. This protocol maintains the Internet connectivity and improves the transmission rate in the scenarios of voice over IP (VoIP), virtual private networks, video streaming and in wireless sensor networks.

The paper provides the details of the implementation of Mobile IP protocol in NS2 [13]. The simulation are performed on network where nodes as stationary and mobile. The result are collected while using the Mobile IP protocol and are compared with the same scenario that does not uses the Mobile IP protocol. A network is proposed to observe and analyze the importance of Mobile IP with the following scenarios:

- The movement of a mobile host from one network another network.
- Transmission of data from a wired device to a wireless mobile node.
- The movement of two mobile nodes moving to a foreign network while receiving data from a fixed node.
- the movement of two mobile nodes to a foreign network while sending data each other.

In the rest of the paper, the details of the Mobile IP are

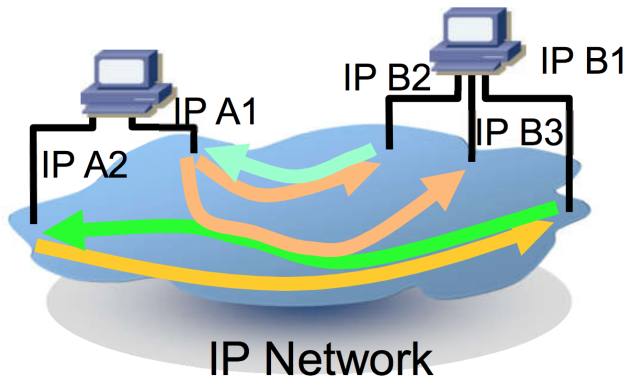


Fig. 1. Sctp multihoming.

presented in Section II. The creation of network scenario and implementation details of mobile IP are described in Section III. The discussion on the observed results are presented in Section IV. The conclusion and recommendations for future are presented in Section V.

II. MOBILITY AT NETWORK LAYER

The Internet protocol that provides seamless connectivity while maintaining the IP address is called the Mobile IP protocol. It is standardized in IETF RFC 5944. The trend of wireless and mobile communication is increasing, the device with multiple interfaces, such as mobile phone are also large in number. Many of the service providers of 4G has incorporated new techniques for the mobility management of the mobile phones. The implementation of mobility management techniques in cellular technology uses the link layer methods, however such methods face many problems specifically in security and trustworthy. Mobile IP is an alternative to such link layer techniques. It provides the location independent routing over the Internet. Each device using the Mobile IP maintains a permanent IP address. When a device travels to a foreign network, it creates a tunnel which redirects the packets from belong the mobile host in home network to the mobile host in foreign network.

While using the Mobile IP, a host maintain two kinds of addresses. First, the home address, it assigned to all of the devices connected to the home network. Second, the care of address, it is the address of a mobile host in the foreign network. Two kinds of agents are also used for the communication purpose. These agents serve as the routers and help in data forwarding. They also maintain a data structure to maintain the IP information for data forwarding. One agent is called the home agent (HA), it stores the permanent IP address when a mobile host is in the home network. Other is called the foreign agent (FA), it advertises the care of address of a mobile host when in the foreign network. The working of these agents is as a three step process, these are agent discovery, agent advertisement and agent solicitation. The key functionality of these agents is as follows:

- Agent discovery is a method used by the mobile host to discover the available agents in the area.
- Agent advertisement is done by the agents through

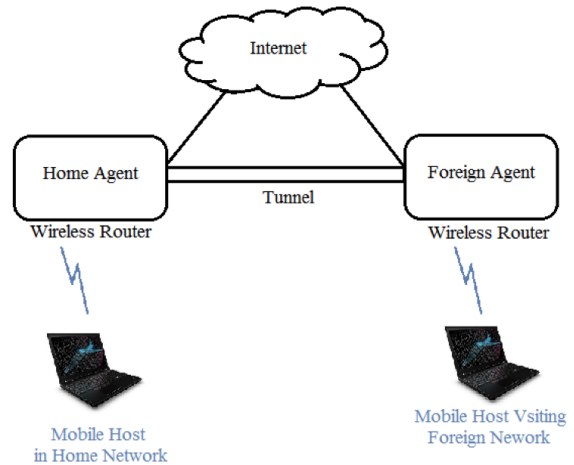


Fig. 2. The concept of Mobile IP.

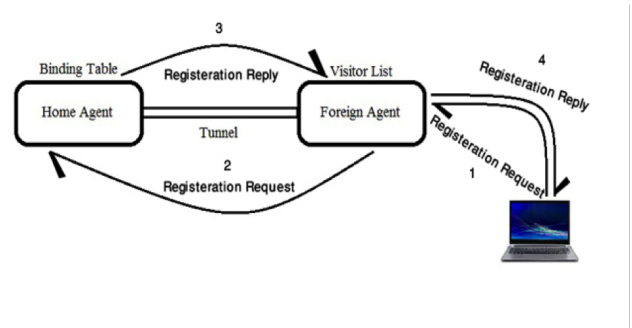


Fig. 3. Registration process in Mobile IP.

the ICMP (Internet message control protocol) to show their presence and availability.

- Agent solicitation messages are used by the mobile hosts to ask for the advertisement from the agents. Once the mobile receives the advertisement, it registers itself with the foreign agent, which forwards information of care of address to the home agent.

A mobile host using Mobile IP is also allowed to register with more than one foreign agents. When the mobile node returns to the home network it de-register itself from other networks. Due to the wireless communication, the Mobile IP face challenges in providing authentication. At the time of registration, the authentication is very critical, Mobile IP uses MD5 (128bits) algorithm to minimize the issue. This algorithm also helps in minimizing the denial of service attack, where an attacker sends a large number of bogus registration requests (see, Fig. 2 and 3).

III. MOBILE IP IMPLEMENTATION AND CONFIGURATION

The simulation scenario consist of wired and wireless nodes as shown in Fig. 4. Topology implementation is taken from the Marc Greis in [14], however the node movements are configured to simulate the two mobile hosts to send and receive data to each other. Number of wired nodes are 2. Number

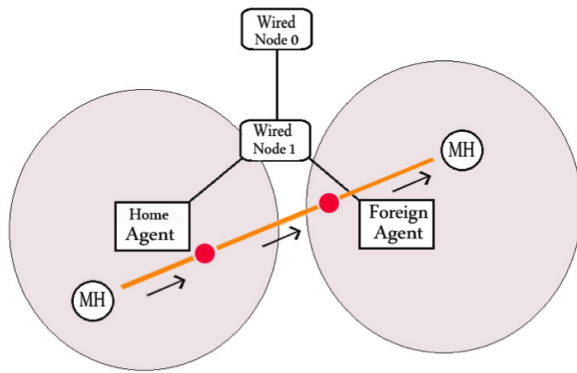


Fig. 4. Network topology.

of mobile nodes are 2. Number of base station nodes are 2. Channel type used is wireless. Propagation model is configured to two-way. Antenna type is omni-antenna. Link Layer type is LL. Packet length is set to 50. Routing Protocol used at the network layer is DSDV (Destination sequenced distance vector). Dimensions of topology are $x=y=670$. The total simulation time is 250seconds. The simulations are carried out in NS2 and results are collected on the basis of average of 20 experiments for each of the simulation scenario. All other parameters of the network layer are set to default values. On the transport layer, TCP protocol is configured. TCP agents are defined as the source and the destination. For example, in first scenario the source is wired node and the destination is mobile host.

After the topology is designed, a simple scenario of two simple wired nodes w_0 & w_1 , which are then connected to each other using 5mbps link. Then two base station are added, which are named as Home Agent (HA) and Foreign Agent (FA), which are connected to the wired node w_1 using 5mbps link. A mobile host is defined to move in between HA and FA. In first movement the mobile host stays remain in the range of Home Network and after 100 seconds it starts travel to the foreign agent. Then after 200 seconds it travels back to the Home Network. From the operating system side, Ubuntu 14.04 is used on the virtual machine. AWK scripting is used to calculate the throughput [15]. Three kinds of simulations are performed by changing the movement of mobile nodes:

- 1) The mobile host moves from the HA to the FA and then back to the area of HA. The data transmitter is w_0 and the receiver is the mobile host.
- 2) In this scenario, two mobile hosts are configured. One is in the HA and one is in the FA. The sender is w_0 and two mobile hosts are the receiver.
- 3) In this scenario, two mobile hosts send and receive the data to each other as shown in Fig. 4.

In all of the simulations, Mobile IP simulation is compared with the same scenario without the Mobile IP.

IV. RESULTS AND DISCUSSIONS

A. One Mobile Host

Two scenarios using one with Mobile IP protocol and other without using Mobile IP protocol are evaluated. In both

scenarios, the start time is 0 and the end time 250seconds. Mobile host from 0s to 100s remains in the range of Home Network. After 100s Mobile Host travels towards the Foreign Network and stays there until 200s. After 200s it travels back to the Home Network.

Fig. 5 clearly shows that from time about 0s to 110s in both cases the throughput is about 0.7Mbps. After 110s the graph suddenly goes down because mobile host in both scenarios are disconnected because mobile host is neither in the range of the home network nor in the range of the foreign network. After some time about 160s Mobile Host using Mobile IP is connected and registered to the foreign network by that data transmission resumes and throughput goes up again about 0.7Mbps. One other hand in the same scenario but mobile host without mobile IP support remains disconnected. When the time slot in graph reaches to 210s throughput goes down due to the coordinates of mobile host, i.e. outside home agent and foreign agent. It remains disconnected until 240s where the mobile host in both scenarios are back to the network and in the graph it is represented with a throughput value of 0.7Mbps. In Fig. 4, from 110s to 240s, it clear that the mobile host remains disconnected when no Mobile IP protocol is used. It is represented by the red line. However, with the use of Mobile IP protocol the mobile host resumes the connection at around time 160s to time 210s while maintaining the connection with the foreign agent. The throughput with and without using the Mobile IP in this scenario are 59.80kbps and 42.81kbps.

B. Two Mobile Hosts

In this scenario two Mobile Hosts are used. The mobile hosts belong to different networks. One mobile host uses home agent HA. Other uses the home agent FA. The simulation is repeated two times, i.e. with and without using the Mobile IP protocol.

Fig. 6 graph shows the performance of Mobile IP when two mobile hosts are used. Mobile node with and without Mobile IP part of HA. In both cases MH uses home access. Throughput rise about 0.7Mbps. After 110s throughput goes 0Mbps. Reason is that mobile node in both conditions are out of reach of home agent and foreign agent. MH configured to Mobile IP regaining throughput at about 155s. MH configured to Mobile IP have throughput about 0.7Mbps. MH without Mobile IP is down because of unknown to foreign network. After 210s both are returning to home. Unavailable throughput for interval of time due to out of range to HA & FA. Reaching back at home location the both lines up to throughput is available on the time slot 240s goes to 0.7Mbps. It is shown in Fig. 6 that with Mobile IP, the mobile host connects to the foreign agent at approximately time 160s and remains connected up to time 210s. Similarly, in Fig. 7, the seconds mobile host resumes connection from approximately 145s to 210s as represented by red line.

C. Mobile Sender and Mobile Receiver

In third scenario, the role of sending and receiving is assigned to the mobile nodes. The mobile node 1 and 2 are configured to send and receive data from each other. The mobile node 1 starts from the HA and mobile node 2 starts from the FA as shown in Fig. 4.

TABLE I. PARAMETERS AND VALUES

Scenario	Average Throughput (kbps)	Packets Sent	Packets Received	End-to-end Delay
1.				
Mobile IP	59.80	15290	14950	106795.95
Without Mobile IP	42.81	11601	10703	76346.42
2.				
Mobile IP	55.99	13868	10572	104115.27
Without Mobile IP	42.29	10927	13541	80811.15
3.				
Mobile IP	64.78	16037	15721	112724.19
Without Mobile IP	36.74	9511	9166	58487.06
4.				
Mobile IP(HA-FA)	98.46	24724	24573	89640.36
Mobile IP(FA-HA)	110.01	27576	27455	59423.14

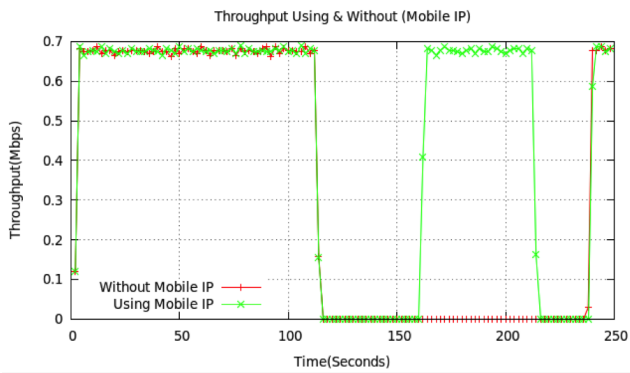


Fig. 5. Experiment 1: One mobile node.

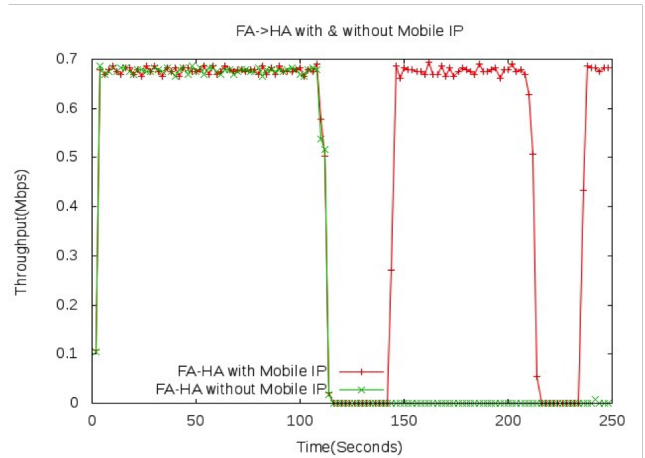


Fig. 7. Experiment 2.2: Two mobile node.

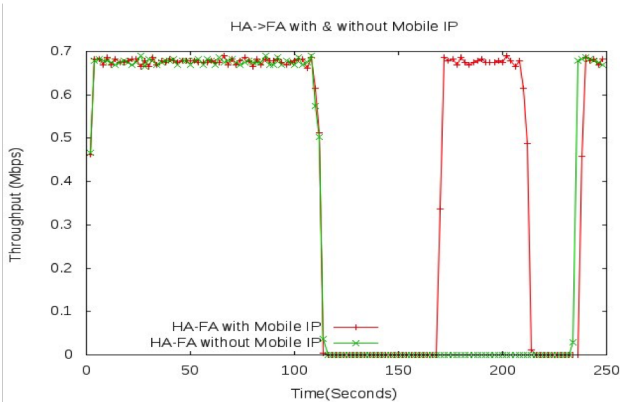


Fig. 6. Experiment 2.1: Two mobile node.

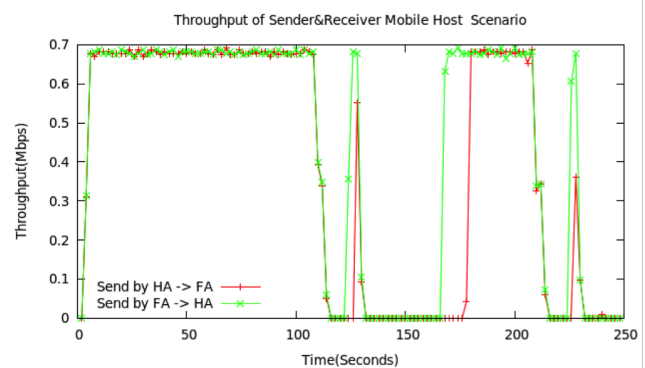


Fig. 8. Experiment 3: Two sending receiving mobile nodes.

Fig. 8 resultant graph showing that from time 0s the both mobile hosts have available throughput is about 0.7Mbps until the time 110s and throughput goes 0Mbps because both mobile goes out of reach of home agent and mobile agent. And both agent again available for throughput for few second and disconnect. Mobile Host of FA is again regaining throughput at about 145s and mobile host of HA is remain unconnected and will throughput available at 160s. Both hosts connectivity

is maintained with 0.7Mbps until 210s after that gain mobile agents are back to the home location remain unconnected until they are in range and throughput for both is available on the time slot 240s at about 0.7Mbps. The overall result of this experiments shows that the increased mobility while using the Mobile IP. The packets sent and received and the throughput

of both mobile node 1 and 2 are depicted in Table I.

V. CONCLUSION AND RECOMMENDATIONS FOR FURTHER STUDY

With the rapid growth in the mobile applications for video streaming and conferencing, the most important is the availability of the Internet connection. The availability of the network connection is related with the mobility, as the mobile devices move from one place to other and join the foreign networks. The transport layer also provide the solutions to the availability and mobility with the help of multi-interface devices. Where, a device maintains the Internet connection through a number of network interface cards. One of the simplest solution to the mobility and the availability of the Internet is the use of Mobile IP protocol at the network layer. Mobile IP uses tunneling concept, that allows a mobile device to use the same IP address in a remote or foreign network. Tunneling redirects the packets of mobile host to the foreign network. This paper provides a detailed study of Mobile IP protocol. The protocol is implemented and tested against the traditional network that do not support mobility. The evaluation shows that the Mobile IP increases the network availability and mobility. By that, it also improves the data transmission rate. In future, the concept of Mobile IP with other mobility management schemes such as handover should be considered. The evaluation of Mobile IP concept with the IPv6 would be good research.

REFERENCES

- [1] Allman, Mark, Vern Paxson, and Ethan Blanton. TCP congestion control. No. RFC 5681. 2009.
- [2] Borella, Michael S. "System and method for control of packet data serving node selection in a mobile internet protocol network." U.S. Patent No. 7,346,684. 18 Mar. 2008.
- [3] Ishakian, Vatche, Joseph Akinwumi, Flavio Esposito, and Ibrahim Matta. "On supporting mobility and multihoming in recursive internet architectures." *Computer Communications* 35, no. 13 (2012): 1561-1573
- [4] Halepoto, Imtiaz Ali. "Scheduling and flow control in CMT-SCTP." HKU Theses Online (HKUTO) (2014).
- [5] Halepoto, Imtiaz A., Francis CM Lau, and Zhixiong Niu. "Scheduling over dissimilar paths using CMT-SCTP." *Ubiquitous and Future Networks (ICUFN), 2015 Seventh International Conference on.* IEEE, 2015.
- [6] BHANGWAR, Noor H., Imtiaz A. HALEPOTO, Intesab H. SADHAYO, Suhail KHOKHAR, and Asif A. LAGHARI. "On Routing Protocols for High Performance." *Studies in Informatics and Control* 26, no. 4: 441-448, 2017.
- [7] Halepoto, Imtiaz A., Francis CM Lau, and Zhixiong Niu. "Management of buffer space for the concurrent multipath transfer over dissimilar paths." *Digital Information, Networking, and Wireless Communications (DINWC), 2015 Third International Conference on.* IEEE, 2015.
- [8] Perkins, Charles E. "IP mobility support for IPv4, revised." (2010).
- [9] Jung, J-W., et al. "Performance evaluation of two layered mobility management using mobile IP and session initiation protocol." *Global Telecommunications Conference, 2003. GLOBECOM'03.* IEEE. Vol. 3. IEEE, 2003.
- [10] Jnsson, Ulf, et al. "MIPMANET: mobile IP for mobile ad hoc networks." *Proceedings of the 1st ACM international symposium on Mobile ad hoc networking & computing.* IEEE Press, 2000.
- [11] Frattasi, Simone, et al. "Defining 4G technology from the users perspective." *IEEE network* 20.1 (2006): 35-41.
- [12] Van de Groenendaal, Johan, and Amitava Chakraborty. "Mobility management in wireless networks." U.S. Patent No. 7,634,252. 15 Dec. 2009.
- [13] Issariyakul, Teerawat, and Ekram Hossain. *Introduction to network simulator NS2.* Springer Science & Business Media, 2011.
- [14] Greis, Marc. "Marc Greis tutorial for the ucblbn/vint network simulator ns." (2004).
- [15] Aho, Alfred V., Brian W. Kernighan, and Peter J. Weinberger. *The AWK programming language.* Addison-Wesley Longman Publishing Co., Inc., 1987.
- [16] Tsang, Ken CK, Cho-Li Wang, and Francis CM Lau. "Handoff performance comparison of mobile IP, fast handoff and mSCTP in mobile wireless networks." *Parallel Architectures, Algorithms, and Networks, 2008. I-SPAN 2008. International Symposium on.* IEEE, 2008

Experimental Results on Agent-Based Indoor Localization using WiFi Signaling

Stefania Monica, Federico Bergenti

Dipartimento di Scienze Matematiche, Fisiche e Informatiche

Università degli Studi di Parma

Parco Area delle Scienze 53/A, 43124 Parma, Italy

Abstract—This paper discusses experimental results on the possibility of accurately estimating the position of smart devices in known indoor environments using agent technology. Discussed localization approaches are based on WiFi signaling, which can be considered as an ubiquitous technology in the large majority of indoor environments. The use of WiFi signaling ensures that no specific infrastructures nor special on-board sensors are required to support localization. Localization is performed using range estimates from the fixed access points of the WiFi network, which are assumed to have known positions. The performance of two range-based localization algorithms are discussed. The first, called Two-Stage Maximum-Likelihood algorithm, is well-known in the literature, while the second is a recent optimization-based algorithm that uses particle swarm techniques. Results discussed in the last part of the paper show that a proper processing of WiFi-based range estimates allows obtaining accurate position estimates, especially if the optimization-based algorithm is used.

Keywords—WiFi-based localization; indoor localization; particle swarm optimization; agent technology

I. INTRODUCTION

In recent years, mobile devices and smart appliances have assumed a relevant impact on everyday life and their number has rapidly increased. Such devices can be commonly used in a variety of contexts, including, e.g., fitness activity and sport, tourism and outdoor navigation, monitoring of environmental parameters, social networks and social games. One of the weaknesses observed in available smart appliances is that they do not yet offer support for accurate indoor localization. While outdoor localization is effectively achieved using various assisted technologies, such as the *Global Positioning System (GPS)*, accurate indoor localization is still an open issue (e.g., [1], [2]). Recently, the use of specific technologies to support indoor localization has been investigated. Among such technologies, it is worth recalling *Ultra Wide Band (UWB)*, which seems to be very promising (e.g., [3], [4]). The large bandwidth and the high time resolution of UWB signals reduce the impact of phenomena which typically interfere with wireless communications in indoor environments (e.g., non-line-of-sight propagation, multi-path, and multiple access interference) [5]. The main drawback of this technology is that it requires a dedicated infrastructure which is not widely available in common indoor scenarios, and which is also still expensive today. In order to overcome such a limitation, we focus in this paper on indoor localization approaches based on the use of WiFi technology, which can be easily found in almost all indoor environments.

In the experimental scenario discussed in this paper, WiFi-

based localization is addressed by means of agent technology. Presented experimental results are obtained by using an add-on module for the *Java Agent and DEvelopment framework (JADE)* [6]. JADE has been developed during the last 20 years and, in 1998, four of the major manufactures of mobile appliances of the time started a joint research initiative [7] to bring agent technology to what we used to call *Java-enabled phones* at the time. The results of such an initiative eventually became the base of JADE for Android [8]. Since then, *nomadic agents* have been considered one of the most promising applications of agent technology and JADE, together with its companion language JADEL [9]–[11], is a consolidated tool in this field which has been used for many applications (e.g., [12]–[14]). The significant opportunities that the synergic combination of agents and smart appliances offer have already been investigated, e.g., in the scope of the *Agent-based Multi-User Social Environment (AMUSE)* project [15]. AMUSE is a major evolution of JADE and it consists of an open-source platform built on top of JADE that addresses specific issues of online social games, which is an application domain where the use of agents is particularly promising. AMUSE has already been fruitfully used to experiment mixed-reality games, in which the possibility of interacting with the physical world becomes crucial [16]. In addition, other examples of games that can be developed using AMUSE include social games in indoor areas with high concentration of potential users, like halls of shopping malls, waiting areas of airports and train stations, interactive museums and exhibits, and covered markets in historic towns. Such areas typically offer dedicated WiFi coverage by means of *Access Points (APs)* spread in the environment, which can be used to support effective localization by means of dedicated techniques as discussed in the following sections. Note that discussed techniques only requires that the WiFi receiver is active on smart appliances and they do not necessitate that appliances are connected to one of the WiFi networks of the area.

The JADE module used to obtain experimental results discussed in this paper can be used to develop agents capable of sensing their positions with respect to a fixed reference frame in known indoor environments [17]. Localization is performed by properly processing estimates of the distances between the smart appliance where the agent runs and the APs of the WiFi network, whose positions are assumed to be known. In such scenarios, estimated distances among the smart appliance and APs can be used to feed a *localization algorithm* which is in charge of providing the agent with an estimate of its position. Experimental results discussed in this paper show that the accuracy of the two discussed localization algorithms

are comparable. Actually, both can achieve an accuracy in the order of 1 m, and even better accuracy can be obtained by properly averaging range estimates. Note that the paper focuses on presenting two of the localization algorithms available in the JADE add-on module, and a detailed description of the architecture of the module is left for a future paper.

This paper is organized as follows. Section II introduces the problem of localization and discusses relevant related work. Section III shows how WiFi signaling can be used to obtain information on the position of smart appliances and introduces useful notations. Section IV illustrates the two discussed localization algorithms. Section V introduces metrics to evaluate the performance of discussed approaches and shows illustrative experimental results. Section VI concludes the paper.

II. RELEVANT APPROACHES FOR INDOOR LOCALIZATION

Many localization algorithms have been proposed in the literature to provide location information to nodes of a network (e.g., [18]). In particular, *range-based localization methods* rely on the knowledge of inter-node distances or angle information and they can be classified into *active* and *passive* [3]. In active methods, all nodes are equipped with sensors and with an electronic device which sends information to a positioning system. Passive localization, instead, is based on the fact that wireless communications strongly depend on the surrounding environment. Relying on the scattering caused by small targets during signal propagation and/or on the variance of a measured signals, changes in the received signals can be used to detect and locate targets and for tracking purposes [19].

In this paper, we focus on active range-based algorithms, which typically involves two steps. First, proper parameters related to signals traveling between the smart appliance and some nodes of the network with known positions are estimated. Such parameters can be the *Time of Flight (ToF)*, the *Angle of Arrival (AoA)*, or the *Received Signal Strength (RSS)* [20]. Then, the parameters evaluated in the first step are used to estimate the position of the smart appliance using a proper localization algorithm. Concerning the first step, note that range-estimation techniques based on AoA rely on the measurements of angles between nodes, which are usually taken by means of antenna arrays that require dedicated hardware. Moreover, the installation cost of antenna arrays can be high and the number of signal paths in indoor environments can be large for the presence of obstacles, which make accurate angle estimation challenging in practical scenarios. Time-based range-estimation techniques, instead, rely on measurements of the ToF of signals traveling between nodes, and, therefore, they require high time resolution in the processing of considered signals and they are not particularly well suited when WiFi technology is used. For the reasons mentioned above, we focus on localization approaches based on RSS, which is typically used when performing localization with WiFi signaling. Under the assumption that the energy of transmitted signals is known, the availability of a relation between the received power of a signal traveling between two nodes and the distance between the nodes can be used to estimate the distance using RSS measurements at the receiver node. In detail, in order to use RSS measurements to estimate the distance that separates a receiving node from a transmitting node, we rely on the *Friis transmission equation*, according to which the received power

$P(\rho)$ at distance ρ can be expressed as [5]

$$P(\rho) = P_0 - 10\beta \log_{10} \frac{\rho}{\rho_0}, \quad (1)$$

where, P_0 is the known power at reference distance ρ_0 and β is a parameter related to the details of the transmission [21]. An estimate of the received power $P(\rho)$ yields the value of the distance ρ by inverting (1).

In the remaining of this paper, we assume that range estimates are obtained from measurements of the RSS of signals traveling between nodes using (1). The adopted add-on module for JADE uses the functionality of the smart appliance where it is installed to obtain the RSS of signals used by WiFi APs to support network scanning. Measured RSS is used to estimate the distances between the smart appliance where the module is installed, which is normally called *Target Node (TN)*, and the APs of the network, which are assumed to be in known and static positions. Each communication between the TN and one of the AP allows obtaining an estimate of the distance between them together with other valuable information, such as the *Basic Service Set Identification (BSSID)* of the responding AP. Assuming that the position of each AP is known, each mapped BSSID can be associated with the coordinates of the corresponding AP and, hence, each distance estimate can also be related to the coordinates of the corresponding AP. The possibility of associating the position of an AP with each distance estimate between a TN and that AP is crucial to support localization, as discussed in the following.

In order to provide agents with information on their current positions, the add-on module for JADE integrates localization algorithms that use estimated distances from APs to estimate the position of the TN and feed it to interested agents. In detail, once range estimates are acquired, they are used to feed one of the available localization algorithms, which computes an estimate of the position of the TN. Finally, interested agents are informed of the current position of the TN, which is used as the position of the agents that are host in the smart appliance where the add-on module is installed. The current version of the JADE module includes different localization algorithm and, in this paper, we compare the performance of two of them:

- The *Two-Stage Maximum-Likelihood (TSML)* algorithm, which is well-known in the literature because it can attain the Cramer-Rao lower bound, which is a lower bound for the variance of an estimator [22].
- A recent optimization-based algorithm [23], which relies on the use of *Particle Swarm Optimization (PSO)* to obtain effective localization.

Preliminary results on the accuracy of discussed algorithms were presented in [24], where only the TSML algorithm was considered to address a localization scenario different from the scenario discussed in the following sections. Experimental results discussed in this paper show that the accuracy of the two localization algorithms are comparable, and that they can both achieve an accuracy in the order of 1 m. Such accuracy can be further improved, especially for the second algorithm, by using proper processing of estimated distances.

III. WIFI-BASED RANGE ESTIMATES

In this section, we introduce relevant notation that will be used to describe the localization scenario and localization algorithms. Throughout the paper, we denote the number of available APs as M . Note that the coordinates of APs are assumed to be known and static, and using a proper coordinate system, they are denoted as

$$\underline{s}_i = (x_i, y_i, z_i)^T. \quad (2)$$

The unknown position of the TN is denoted as

$$\underline{u} = (\bar{x}, \bar{y}, \bar{z})^T. \quad (3)$$

For the sake of simplicity, in the following we assume that there exist a unique TN whose position needs to be estimated, but it is worth noting that discussed approaches can be used also with a large number of TN, provided they are all equipped with the JADE add-on module. Using the introduced notation, the distance between the TN and the i -th AP is

$$\rho_i = \|\underline{u} - \underline{s}_i\|,$$

where $\|\underline{x}\|$ denotes the Euclidean norm of vector \underline{x} . The position of the TN defined in (3) can be found by intersecting the spheres with radii $\{\rho_i\}_{i=1}^M$ centered in $\{\underline{s}_i\}_{i=1}^M$, i.e., by solving the following non-linear system of equations:

$$\begin{cases} (x - x_1)^2 + (y - y_1)^2 + (z - z_1)^2 = \rho_1^2 \\ \vdots \\ (x - x_M)^2 + (y - y_M)^2 + (z - z_M)^2 = \rho_M^2. \end{cases} \quad (4)$$

In order to guarantee that the system of equations (4) has a unique solution, the number of available APs needs to be at least equal to 4.

Unfortunately, when performing localization in realistic scenarios, the exact values of distances $\{\rho_i\}_{i=1}^M$ between the TN and APs are unknown. As a matter of fact, their knowledge together with the system of equations (4) would also imply that the true TN position \underline{u} would be known. For this reason, it is necessary to rely on range estimates acquired by a proper processing of the RSS of WiFi signals. As soon as range estimates from the M APs are available, a proper localization algorithm can use them to estimate the position of the TN. It is assumed that this acquisition and processing procedure is iterated $L > 1$ times, thus leading to L position estimates at different instants. We denote the estimated distance between the i -th AP and the TN at the j -th iteration as $\hat{\rho}_{i,j}$.

Before describing the localization algorithms used to obtain discussed results, let us make additional comments on range estimates. Instead of using a single range estimate from each APs, it is possible to acquire $K > 1$ range estimates from each AP and then use the average of such range estimates to derive an estimate of the TN position. In detail, let us define

$$\hat{\rho}_{i,j}^K = \frac{1}{K} \sum_{h=j}^{j+K-1} \hat{\rho}_{i,h} \quad (5)$$

which represent the averaged range estimates from the i -th AP obtained by averaging $\hat{\rho}_{i,j}$ over K consecutive acquisitions. Averages $\hat{\rho}_{i,j}^K$ can be used instead of single estimates $\hat{\rho}_{i,j}$ to alleviate problems related to acquisition noise and multipath.

From now on, we make an additional assumption which allows simplifying the localization algorithm, i.e., we assume that the height \bar{z} of the considered TN is known. Even though this may seem a strong assumption, we remark that, in considered scenarios, users are holding their smart appliances, i.e., the TNs, in their hands or in their pockets. Hence, even if the true height is not accurately known, it can be reasonably approximated to, e.g., $\bar{z} = 1$ m. Errors in the order of a few centimeters on the value of \bar{z} do not have a strong impact on the accuracy of discussed algorithms because range estimates error are typically in the order of 10 cm [25]. Additionally, this assumption has the advantage of simplifying the localization algorithm as if the considered scenario was a bidimensional one, namely, as if the coordinates of the i -th AP were (x_i, y_i, \bar{z}) for $i \in \{1, \dots, M\}$. In other words, the third coordinate can be neglected and the mathematical description of the discussed localization algorithms is simplified. Let us define the difference between the height \bar{z} of the TN and the height z_i of the i -th AP as

$$h_i = \|\bar{z} - z_i\|. \quad (6)$$

Given this definition, it is possible to evaluate the projections of the distances $\{\rho_i\}_{i=1}^M$ between the TN and the i -th AP on the plane $z = \bar{z}$ where the TN lies. According to the Pythagorean theorem, under the assumption that the true values of distances $\{\rho_i\}_{i=1}^M$ are known, the projections of such distances can be written as

$$r_i = \sqrt{\rho_i^2 - h_i^2}. \quad (7)$$

An illustrative geometric representation of the relations among ρ_i , h_i , and r_i is shown in Fig. 1, assuming that the coordinates of the considered AP, expressed in meters, are $(0, 0, 3)$ and those of the TN, expressed in meters, are $(1, 1, 1)$.

Under the assumption that the height \bar{z} of the TN is known, the abscissa and the ordinate of the TN can be found by intersecting the circumferences centered in $\{(x_i, y_i)\}_{i=1}^M$ with radii $\{r_i\}_{i=1}^M$, i.e., by solving the following non-linear system of equations:

$$\begin{cases} (x - x_1)^2 + (y - y_1)^2 = r_1^2 \\ \vdots \\ (x - x_M)^2 + (y - y_M)^2 = r_M^2. \end{cases} \quad (8)$$

As previously observed, the values of $\{\rho_i\}_{i=1}^M$ are unknown and, consequently, the values of $\{r_i\}_{i=1}^M$ are also unknown. Hence, it is necessary to rely on range estimates obtained from the RSS of WiFi signals. Let us define the projections of distances $\hat{\rho}_{i,j}$ from the i -th AP at the j -th iteration on the plane $z = \bar{z}$ as

$$\hat{r}_{i,j} = \sqrt{\hat{\rho}_{i,j}^2 - h_i^2}. \quad (9)$$

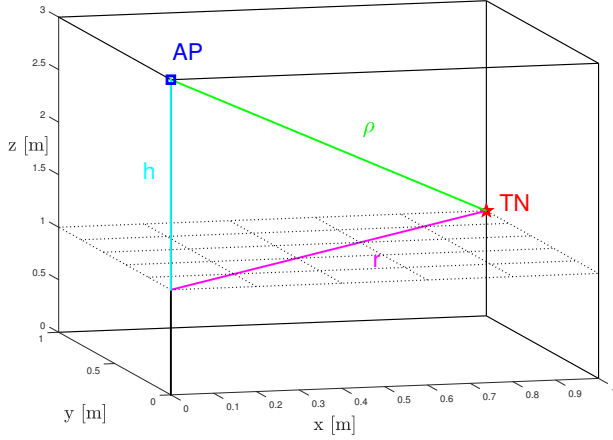


Fig. 1. The distance ρ between an AP (blue square) and a TN (red star) is shown (green line), together with the height h of AP with respect to the TN (cyan line). The projection r of the distance ρ on the plane on which the TN lies is also shown (magenta line).

The position estimate of the TN at the j -th iteration will be denoted as

$$\hat{\underline{u}}_j = (\hat{x}_j, \hat{y}_j, \bar{z}). \quad (10)$$

Similarly, from the definition of averaged range estimates $\hat{\rho}_{i,j}^K$ introduced in (5), it is also possible to define the projections of the averaged range estimates introduced in (5) on the plane $z = \bar{z}$ as

$$\hat{r}_{i,j}^K = \sqrt{(\hat{\rho}_{i,j}^K)^2 - h_i^2}. \quad (11)$$

Finally, the position estimates obtained by feeding the localization algorithm with averaged range estimates $\{\hat{r}_{i,j}^K\}_{i=1}^M$ are denoted as

$$\hat{\underline{v}}_j^K = (\hat{x}_j^K, \hat{y}_j^K, \bar{z}). \quad (12)$$

We remark that the j -th position estimate $\hat{\underline{v}}_j^K$ relies on the M averaged range estimates $\{\hat{\rho}_{i,j}^K\}_{i=1}^M$, and, therefore, on the $M \cdot K$ range estimates $\{\hat{\rho}_{i,h}^K\}_{i=1}^M$ with $h \in \{j, \dots, j+K-1\}$. For this reason, the first position estimate $\hat{\underline{v}}_1^K$ can be obtained only after an initial phase during which K range estimates from each of the M APs are acquired. Once $\hat{\underline{v}}_1^K$ is evaluated, the position estimates that follow can be determined as soon as a new M -tuple of range estimates from the M APs is acquired. We remark that using the notation introduced for averaged range estimates it is also possible to consider the case without range averaging by setting $K = 1$ in (5).

IV. TWO RELEVANT LOCALIZATION ALGORITHMS

Various range-based localization algorithms have been proposed in the literature (e.g., [26]) and they can all be integrated with the adopted JADE add-on module. This section describes the two algorithms used to obtain experimental results shown in next section, namely the TSML algorithm and the PSO-based algorithm. The starting point for the considered localization algorithms is (8), where the exact distances $\{r_i\}_{i=1}^M$ are replaced by their estimates $\hat{r}_{i,j}$. For the sake of simplicity, in the descriptions of localization algorithms we neglect the

subscript j , which counts iterations, and we denote as $\{\hat{r}_i\}_{i=1}^M$ a generic set of range estimates at a given iteration. Using this notation, the non-linear system of equations (8) is replaced by

$$\begin{cases} (\hat{x} - x_1)^2 + (\hat{y} - y_1)^2 = \hat{r}_1^2 \\ \vdots \\ (\hat{x} - x_M)^2 + (\hat{y} - y_M)^2 = \hat{r}_M^2. \end{cases} \quad (13)$$

The system of equations (13) shows the equations of the M circumferences lying on the plane $z = \bar{z}$, centered in $\{(x_i, y_i)\}_{i=1}^M$ with radii $\{\hat{r}_{i,j}\}_{i=1}^M$. If the radii of such circumferences were the exact distances $\{r_i\}_{i=1}^M$, they would all intersect in the same point, which would correspond to the exact TN position. Instead, the radii of circumferences in (13) are range estimates, and they are affected by errors. Hence, circumferences do not intersect in a unique point and, therefore, proper localization algorithms are needed. In the following, we denote the solution of (13) as $\hat{\underline{u}} = (\hat{x}, \hat{y})$.

A. Two-Stage Maximum-Likelihood Algorithm

In order to simplify the description of the algorithm, let us define the following quantity, related to the solution $\hat{\underline{u}}$ of (13)

$$\hat{n} = \hat{x}^2 + \hat{y}^2. \quad (14)$$

Then, the system of equations (13) can be reformulated in matrix notation as

$$\underline{G}_1 \hat{\underline{\omega}} = \hat{\underline{h}}_1, \quad (15)$$

where

$$\underline{G}_1 = \begin{pmatrix} -2x_1 & -2y_1 & 1 \\ \vdots & \vdots & \vdots \\ -2x_M & -2y_M & 1 \end{pmatrix} \quad (16)$$

$$\hat{\underline{\omega}} = \begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{n} \end{pmatrix} \quad \hat{\underline{h}}_1 = \begin{pmatrix} \hat{r}_{1,j}^2 - a_1^2 \\ \vdots \\ \hat{r}_{M,j}^2 - a_M^2 \end{pmatrix}$$

and

$$a_i = \sqrt{x_i^2 + y_i^2} \quad (17)$$

We remark that (15) is not a linear system since the third element of vector $\hat{\underline{\omega}}$ depends on the first two elements according to (14). Neglecting this dependence, the solution $\hat{\underline{\omega}}$ of (15) is determined through a *Maximum-Likelihood (ML)* approach as

$$\hat{\underline{\omega}} = (\underline{G}_1^T \underline{W}_1 \underline{G}_1)^{-1} \underline{G}_1^T \underline{W}_1 \hat{\underline{h}}_1, \quad (18)$$

where \underline{W}_1 is a positive definite matrix [27]. For the sake of simplicity, in the implementation used to obtain discussed experimental results, matrix \underline{W}_1 is equal to the identity matrix.

Once the solution $\hat{\underline{\omega}}$ of (15) is evaluated, the dependence of \hat{n} on \hat{x} and \hat{y} can be taken into account by considering the following system of equations:

$$\underline{G}_2 \hat{\underline{\phi}} = \hat{\underline{h}}_2, \quad (19)$$

where $\hat{\phi} = (\hat{x}^2, \hat{y}^2)^\top$ and

$$\underline{G}_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} \quad \hat{h}_2 = \begin{pmatrix} \hat{\omega}_1^2 \\ \hat{\omega}_2^2 \\ \hat{\omega}_3 \end{pmatrix}. \quad (20)$$

Let us remark that $\hat{\omega}_j$ in (20) denotes the j -th component of $\hat{\omega}$. The solution of the rectangular system (19) can be determined, using a ML approach, as

$$\hat{\phi} = (\underline{G}_2^\top \underline{W}_2 \underline{G}_2)^{-1} \underline{G}_2^\top \underline{W}_2 \hat{h}_2, \quad (21)$$

where \underline{W}_2 is a positive definite matrix. In the implementation used to obtain discussed experimental results, matrix \underline{W}_2 is set equal to the identity matrix, as done for \underline{W}_1 . Given the solution $\hat{\phi}$, the estimated abscissa and ordinate at a generic iteration of the considered localization algorithm can be expressed as

$$\tilde{u} = (\hat{x}, \hat{y}) = \underline{U} \left(\sqrt{\hat{\phi}_1}, \sqrt{\hat{\phi}_2} \right)^\top, \quad (22)$$

where $\underline{U} = \text{diag}(\text{sign}(\hat{\omega}))$, and $\hat{\phi} = (\hat{\phi}_1, \hat{\phi}_2)$. Finally, reintroducing the subscript j to denote the corresponding iteration, the estimated position of the TN at the j -th iteration can be written as

$$\hat{u}_j = (\hat{x}_j, \hat{y}_j, \bar{z}). \quad (23)$$

The same localization algorithm can be applied also to the averaged range estimates $\hat{r}_{i,j}^K$ defined in (11). In this case, the initial system of equations can be obtained from (13) by substituting range estimates $\hat{r}_{i,j}$ with the averaged range estimates $\hat{r}_{i,j}^K$, namely,

$$\begin{cases} (\hat{x} - x_1)^2 + (\hat{y} - y_1)^2 = (\hat{r}_{1,j}^K)^2 \\ \vdots \\ (\hat{x} - x_M)^2 + (\hat{y} - y_M)^2 = (\hat{r}_{M,j}^K)^2. \end{cases} \quad (24)$$

The M equations in (24) represent the M circumferences lying on the plane $z = \bar{z}$, centered in $\{(x_i, y_i)\}_{i=1}^M$, with radii equal to the averaged range estimates $\{\hat{r}_{i,j}^K\}_{i=1}^M$. By applying the TSML algorithm to the system of equations in (24), it is possible to obtain position estimates evaluated using averaged distances and, in the following, such position estimates are denoted as

$$\hat{u}_j^K = (\hat{x}_j^K, \hat{y}_j^K, \bar{z}). \quad (25)$$

B. The PSO-Based Algorithm

Observe that the system of equations (13) can be re-written in matrix notation as

$$\underline{1} \tilde{u}^\top \tilde{u} + \underline{A} \tilde{u} = \hat{k}, \quad (26)$$

where $\underline{1}$ is the vector with M elements equal to 1, \hat{k} is a vector whose i -th element is $\hat{r}_i^2 - (x_i^2 + y_i^2)$, and \underline{A} is the

following $M \times 2$ matrix

$$\underline{A} = \begin{pmatrix} -2x_1 & -2y_1 \\ -2x_2 & -2y_2 \\ \vdots & \vdots \\ -2x_M & -2y_M \end{pmatrix}. \quad (27)$$

Also observe that the solution of the system of equations (26), and, hence, of the system of equations (13), can be reinterpreted as the solution of a related optimization problem. In detail, the solution $\tilde{u} = (\hat{x}, \hat{y})$ can be found as the solution of the following minimization problem:

$$\tilde{u} = \arg \min_{\underline{u}} F(\underline{u}) \quad (28)$$

where the fitness function $F(\underline{u})$ is defined as

$$F(\underline{u}) = \|\hat{k} - (\underline{1} \underline{u}^\top \underline{u} + \underline{A} \underline{u})\|. \quad (29)$$

In order to solve the minimization problem (28), thus finding estimates for the abscissa and the ordinate of the TN, we proposed to use the PSO algorithm [28]. The PSO algorithm was first introduced in [29] and it considers the set of potential solutions of an optimization problem as a swarm of S particles which move through a *search space* according to proper rules. In detail, it is assumed that, at every instant, each particle is associated with a position in the search space and with a velocity. Positions and velocities of particles are iteratively updated according to proper rules, which are meant to move all particles towards the solution of the minimization problem, namely towards the position which minimizes the fitness function (29). The rules that are normally adopted to update positions and velocities are inspired by the rules which govern the behaviors of birds in swarms.

The use of PSO to support the localization of a TN works as follows. At initialization, the positions of the particles are randomly initialized in the search space, which, in our context, corresponds to the plane $z = \bar{z}$ where the abscissa and the ordinate of the TN are supposed to be situated. The initial positions are denoted as $\underline{x}^{(i)}(0)$, where $i \in \{1, \dots, S\}$ is the index of a generic particle and S is the number of particles. Similarly, the velocities of all the particles are randomly initialized and they are denoted as $\underline{w}^{(i)}(0)$, where $i \in \{1, \dots, S\}$ is the index of a generic particle. After this initialization phase, positions and velocities of all particles are updated at each iteration $t > 0$ to simulate interactions among individuals [30]. The position of the i -th particle at the t -th iteration is denoted as $\underline{x}^{(i)}(t)$ and its velocity is denoted as $\underline{w}^{(i)}(t)$. At each iteration, the velocity of the i -th particle is updated according to a specific rule expressed as the sum of three addends. In detail, the velocity of the i -th particle at the $(t+1)$ -th iteration is [31]

$$\begin{aligned} \underline{w}^{(i)}(t+1) = & \omega(t) \underline{w}^{(i)}(t) + c_1 R_1(t) (\underline{y}^{(i)}(t) - \underline{x}^{(i)}(t)) \\ & + c_2 R_2(t) (\underline{u} - \underline{x}^{(i)}(t)), \end{aligned} \quad (30)$$

where the following quantities [32] are used:

- $\underline{y}(t)$ is the best position globally reached so far;
- $\underline{y}^{(i)}(t)$ is the best position reached so far by the i -th particle;

- $\omega(t)$ is called *inertial factor*;
- c_1 is called *cognition* parameter;
- c_2 is called *social* parameter; and
- $R_1(t)$ and $R_2(t)$ are independent random variables uniformly distributed in $(0, 1)$.

The first addend in (30) is related to the velocity of the i -th particle at previous iteration t , which is weighed according to the inertial factor $\omega(t)$. In the second addend, the cognition parameter is a positive real parameter and the best position reached so far by the i -th particle can be expressed as

$$\underline{y}^{(i)}(t) = \arg \min_{\underline{z} \in X^{(i)}} F(\underline{z}), \quad (31)$$

where $X^{(i)} = \{\underline{x}^{(i)}(0), \dots, \underline{x}^{(i)}(t)\}$. Hence, this second addend is meant to move each particle towards its best position reached so far. Finally, in the third addend, the social parameter is a positive real parameter and the best position reached so far by any particle in the swarm can be expressed as

$$\underline{y}(t) = \arg \min_{\underline{z} \in Y^{(i)}} F(\underline{z}), \quad (32)$$

where $Y^{(i)} = \{\underline{y}^{(1)}(t), \dots, \underline{y}^{(S)}(t)\}$. The third addend is meant to move each particle towards the *global best position*, namely the position which corresponds to the lowest value of the fitness function among all those reached by any particle in the swarm [30]. The velocities computed with (30) are used to update the positions of particles at each iteration. Such updates are performed by adding the velocities evaluated in (30) to the previous positions of each particle, namely

$$\underline{x}^{(i)}(t+1) = \underline{x}^{(i)}(t) + \underline{w}^{(i)}(t+1). \quad (33)$$

The PSO algorithm is iterated until a termination condition is met. One of the possible termination conditions for the PSO algorithm is the reach of a maximum number of iterations. Once the execution of the algorithm is terminated, the solution corresponds to the position of the particle in the global best position, namely the position of the particle with the lowest value of the fitness function. This solution correspond to the estimated abscissa and ordinate of the TN.

The PSO algorithm outlined previously is used to solve the localization problem formulated in (28) and, hence, to solve the localization problem described in (13). The same algorithm can be applied also when the system of equations (24), where averaged range estimates $\hat{r}_{i,j}^K$ defined in (11) appear, is used. In this case, (24) can be re-written in matrix notation as

$$\underline{1} \underline{\tilde{u}}^T \underline{\tilde{u}} + \underline{A} \underline{\tilde{u}} = \underline{\hat{k}}_j \quad (34)$$

where $\underline{1}$ is the vector with M elements equal to 1, $\underline{\hat{k}}_j$ is a vector whose i -th element is $(\hat{r}_{i,j}^K)^2 - (x_i^2 + y_i^2)$, and \underline{A} is defined in (27). By applying the PSO algorithm to (34), it is possible to obtain position estimates evaluated using averaged distances and, in the following, such position estimates are denoted as

$$\hat{\underline{v}}_j^K = (\hat{x}_j^K, \hat{y}_j^K, \bar{z}). \quad (35)$$

Experimental results shown in the remaining of this paper are obtained with a population of $S = 40$ particles. The inertial factor is set to $\omega(t) = 0.5$ and the values of c_1 and c_2 are both set to 2, so that the average values of $c_1 R_1(t)$ and of $c_2 R_2(t)$ correspond to 1. The termination condition for the PSO algorithm corresponds to the reach of 50 iterations. These values proved to be effective for localization purposes [33]. Illustrative experimental results about the performance of the PSO-based algorithm are shown in next section.

V. PERFORMANCE EVALUATION

In the experimental campaign described in this section we consider three values for K , namely, $K = 1$ (i.e., no averaging); $K = 10$; and $K = 100$. The performance of discussed localization approaches is evaluated in terms of the distances between the true TN position and its estimates. In order to evaluate the performance of the discussed localization algorithms, let us define the distance error as

$$d_j = \|\hat{\underline{u}}_j - \underline{u}\|. \quad (36)$$

Observe that, since we assume that the height of the TN is known, the third component of $\hat{\underline{u}}_j$ is equal to the third component of the vector \underline{u} which represents the true TN position. Therefore, (36) represents the projection of the distance error on the plane $z = \bar{z}$. The definition of the distance error (36) allows introducing the maximum value of the distance errors, which can be denoted as

$$d_{\max} = \max_{j \in \{1, \dots, L\}} d_j. \quad (37)$$

Let us also introduce the average value of the distance error, which can be expressed as

$$d_{\text{avg}} = \frac{1}{L} \sum_{j=1}^L d_j. \quad (38)$$

Finally, the standard deviation of the distance error is

$$\sigma_d = \sqrt{\frac{1}{L} \sum_{j=1}^L (d_j - d_{\text{avg}})^2}. \quad (39)$$

Analogous values relative to position estimates $\hat{\underline{v}}_j^K$ obtained using averaged range estimates can be defined. In detail, let us define

$$\delta_j^K = \|\hat{\underline{v}}_j^K - \underline{u}\| \quad (40)$$

which represents the distance error on the plane $z = \bar{z}$ between the true TN position and its estimate in the j -th iteration obtained using averaged range estimates over K consecutive range acquisitions. The definition of δ_j^K allows computing the maximum value of the distance error as

$$\delta_{\max}^K = \max_{j \in \{1, \dots, L\}} \delta_j^K, \quad (41)$$

and the average value of the distance error as

$$\delta_{\text{avg}}^K = \frac{1}{L} \sum_{j=1}^L \delta_j^K. \quad (42)$$

Finally, the standard deviation of the distance error is

$$\sigma_{\delta}^K = \sqrt{\frac{1}{L} \sum_{j=1}^L (\delta_j^K - \delta_{\text{avg}}^K)^2}. \quad (43)$$

Observe that the values in (37), (38), and (39) can be equally defined using the more general notation in (41), (42), and (43), respectively, for $K = 1$.

In order to assess the accuracy of the localization algorithms previously described, we performed an experimental campaign in an illustrative indoor scenario which consists of a square room whose sides are 4 m long. The considered scenario is shown in Fig. 2, where $M = 3$ APs are shown (blue squares). Observe that $M = 3$ is the minimum number of APs which allows the application of described localization algorithms. The coordinates of the APs are denoted as $\{\text{AP}_i\}_{i=1}^3$ and they are positioned in the room in such a way that, in a proper coordinate system, they can be expressed, in meters, as

$$\begin{aligned} \text{AP}_1 &= (0, 0, 3)^\top \\ \text{AP}_2 &= (0, 4, 3)^\top \\ \text{AP}_3 &= (4, 4, 3)^\top. \end{aligned} \quad (44)$$

In Fig. 2, three different TN positions are also shown (red stars) and their position in the same coordinate system can be expressed in meters as

$$\begin{aligned} \underline{u}_1 &= (1, 1, 1)^\top \\ \underline{u}_2 &= (1, 2, 1)^\top \\ \underline{u}_3 &= (2, 2, 1)^\top. \end{aligned} \quad (45)$$

Using the described configuration of fixed APs, the three different TN positions are estimated. Results of such position estimates are discussed in the remaining of this section, using both localization algorithms introduced in previous section, on the basis of the distance error discussed above. Let us remark that, even if in the considered scenario all APs are placed at the same height (i.e., 3 m), and all TNs are placed at the same height (i.e., 1 m), both proposed localization algorithms are general and they do not require that APs share the same height. Moreover, different heights for the TNs could also be considered, provided that they are known, so that the values of h_i defined in (6) can be computed.

In the remaining of this section, relevant comments on the localization of the three TNs are presented. In all scenarios, the number of iterations is set equal to $L = 100$. Hence, the average value and the standard deviations of distance errors is based on 100 position estimates. In the following figures and tables, in order to distinguish between position estimates obtained using the two algorithms, we add superscript (T) and (P) to denote position estimates derived using the TSML and the PSO-based algorithm, respectively.

A. First Scenario

We start by considering the TN denoted as TN_1 in Fig. 2, whose coordinates are $\underline{u}_1 = (1, 1, 1)^\top$. Since the coordinates of the APs are also known, the true distances $\{\rho_i\}_{i=1}^3$ between

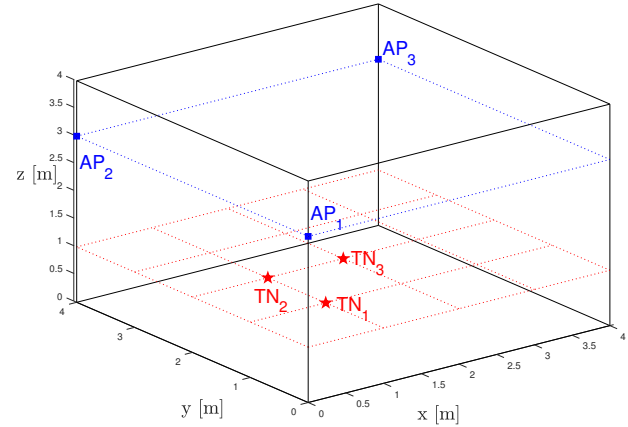


Fig. 2. The positions of the three considered APs (blue squares) and three different TN positions (red stars) are shown.

the i -th AP and the TN can be computed. In detail, from (44) and (45), the values of $\{\rho_i\}_{i=1}^3$ can be computed as

$$\rho_1 \simeq 2.45 \text{ m} \quad \rho_2 \simeq 3.74 \text{ m} \quad \rho_3 \simeq 4.69 \text{ m}.$$

The projections $\{r_i\}_{i=1}^3$ of the range estimates $\{\rho_i\}_{i=1}^3$ on the plane $z = 1$ m can be also computed according to the Pythagorean theorem. Simple algebraic manipulations show that the values of $\{r_i\}_{i=1}^3$ are

$$r_1 \simeq 1.41 \text{ m} \quad r_2 \simeq 3.16 \text{ m} \quad r_3 \simeq 4.24 \text{ m}.$$

In order to apply the discussed algorithms, we acquire range estimates from each AP to have: 100 position estimates obtained without range averaging; 100 position estimates obtained by averaging over $K = 10$ consecutive range estimates; and 100 position estimates obtained by averaging over $K = 100$ consecutive range estimates.

Fig. 3 shows the true position of TN_1 on the plane $z = 1$ (red star). In the same figure, the projections of the position estimates on the plane $z = 1$ (black circles) obtained without averaging are also shown, together with the projections of the position estimates (magenta crosses) obtained by averaging over 10 consecutive range estimates from each AP, and the projections of the position estimates (green triangles) obtained by averaging over 100 consecutive range estimates from each AP. From Fig. 3 it can be observed that, as expected, the position estimates are closer to each other as K increases. This is in agreement with the intuitive idea that averaging range estimates over large values of K leads, on average, to more accurate values of $\hat{r}_{i,j}^K$. The more accurate are the range estimates used to feed the localization algorithm, the more accurate are the obtained position estimates. Fig. 4 shows the *Cumulative Distribution Functions (CDFs)* of the distance errors without range averaging (black line), with $K = 10$ (dashed magenta line), and with $K = 100$ (dash-dotted green line). As intuitively expected from Fig. 3, the larger is K , the steepest is the graph of the CDF.

Fig. 5 shows the true position of TN_1 on the plane $z = 1$ (red star), together with: the projections of position estimates

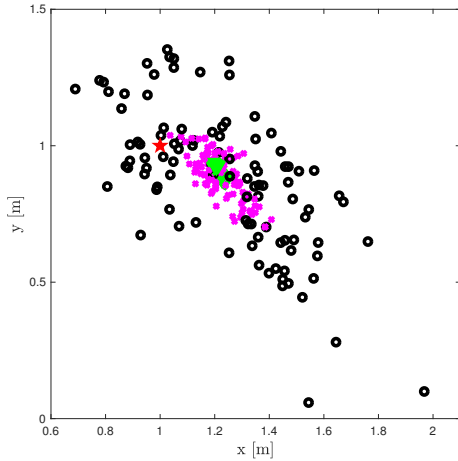


Fig. 3. Projections of the position estimates of TN_1 in Fig. 2 obtained using the TSML algorithm on the plane $z = 1$: (i) without range averaging (black circles); (ii) with $K = 10$ (magenta crosses); and (iii) with $K = 100$ (green triangles). The true TN position (red star) is also shown.

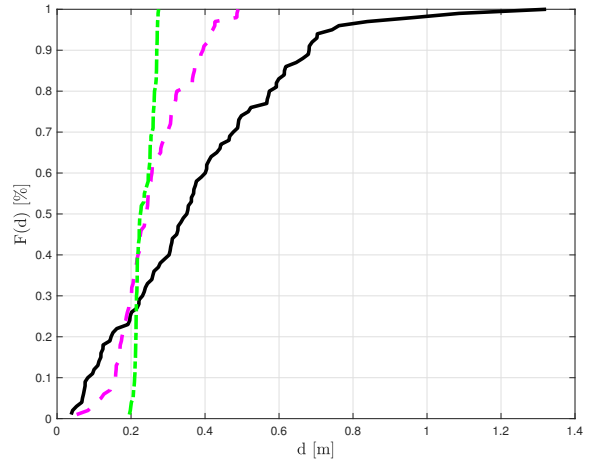


Fig. 4. The cumulative distribution function of distance errors from Fig. 3: (i) without range averaging (black line); (ii) with $K = 10$ (dashed magenta line), and (iii) with $K = 100$ (dash-dotted green line), relative to the position estimates of TN_1 obtained using the TSML algorithm.

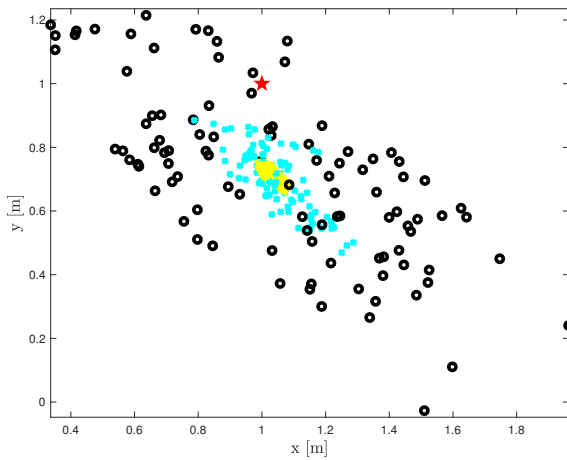


Fig. 5. Projections of the position estimates of TN_1 in Fig. 2 obtained using the PSO algorithm on the plane $z = 1$: (i) without range averaging (black circles); (ii) with $K = 10$ (cyan crosses); and (iii) with $K = 100$ (yellow triangles). The true TN position (red star) is also shown.

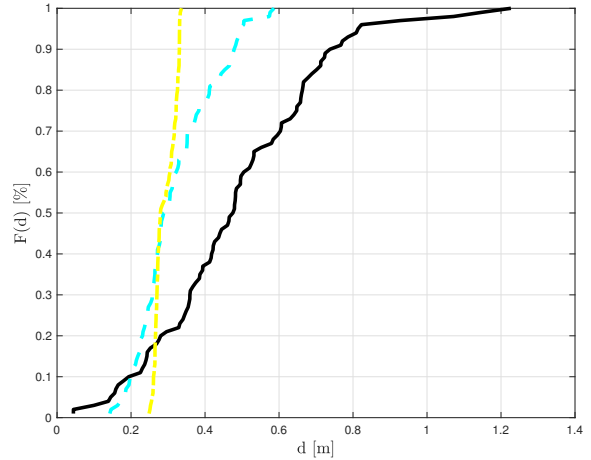


Fig. 6. The cumulative distribution function of distance errors from Fig. 5: (i) without range averaging (black line); (ii) with $K = 10$ (dashed cyan line), and (iii) with $K = 100$ (dash-dotted yellow line), relative to the position estimates of TN_1 obtained using the PSO algorithm.

on the plane $z = 1$ (black circles) obtained without range averaging; the projections of position estimates (cyan crosses) obtained by averaging over 10 consecutive range estimates from each AP; and the projections of position estimates (yellow triangles) obtained by averaging 100 consecutive range estimates from each AP. As observed when using the TSML algorithm, the position estimates are closer to the TN position as K increases. Hence, also when considering the PSO-based algorithm, more accurate range estimates lead to more accurate position estimates. Fig. 6 shows the CDFs of distance errors without range averaging (black line), with $K = 10$ (dashed cyan line), and with $K = 100$ (dash-dotted yellow line). As expected, larger values of K correspond to steeper CDFs.

Table I shows the values of the maximum distances and of the average distances between the considered TN and its estimates, and the values of standard deviations of distance

errors for $K = 1$, $K = 10$, and $K = 100$. It can be observed that when localization is performed using the TSML algorithm,

TABLE I. VALUES OF THE MAXIMUM DISTANCE ERROR (FIRST ROW), OF THE AVERAGE DISTANCE ERROR (SECOND ROW), AND OF THE STANDARD DEVIATION OF THE DISTANCE ERROR (THIRD ROW) ARE SHOWN, FOR DIFFERENT VALUES OF K , RELATIVE TO POSITION ESTIMATES OF TN_1 OBTAINED WITH THE TSML ALGORITHM AND WITH THE PSO-BASED ALGORITHM, RESPECTIVELY

TN_1	TSML			PSO			
	K	1	10	100	1	10	100
δ_{\max}^K [m]		1.32	0.49	0.27	1.23	0.59	0.34
δ_{avg}^K [m]		0.38	0.26	0.24	0.49	0.32	0.29
σ_{δ}^K [m]		0.24	0.09	0.02	0.23	0.10	0.03

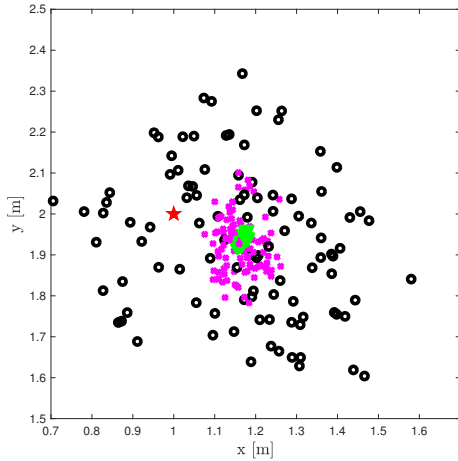


Fig. 7. Projections of the position estimates of TN_2 in Fig. 2 obtained using the TSML algorithm on the plane $z = 1$: (i) without range averaging (black circles); (ii) with $K = 10$ (magenta crosses); and (iii) with $K = 100$ (green triangles). The true TN position (red star) is also shown.

the values of the maximum distance between the considered TN and its estimates decrease as K increases. In detail, without range averaging the values of the maximum distance between the considered TN and its estimates is equal to 1.32 m and it decreases to 0.49 m when $K = 10$ and to 0.27 m when $K = 100$. Analogous considerations hold when analyzing results relative to the PSO-based algorithm. In this case, the value of the maximum distance between the considered TN and its estimates without range averaging equals 1.23 m, and it is slightly smaller than that evaluated when considering the TSML algorithm. The values of the maximum distance between the considered TN and its estimates are equal to 0.59 m and to 0.34 m when $K = 10$ and $K = 100$, respectively, and they are slightly higher than those obtained by applying TSML algorithm. From the fourth row of Table I it can be observed that also the values of the average distance between the considered TN and its estimates decrease as K increases, starting from 0.38 m when no range averaging is considered to 0.24 m when $K = 100$ is considered. The values of the average distance δ_{avg}^K between the considered TN and its estimates when using the PSO-based algorithm are comparable to those obtained when using the TSML algorithm, and they are equal to 0.49 m, when range averages are not performed, and to 0.32 m and 0.29 m when $K = 10$ and $K = 100$, respectively. Finally, Table I shows that the values of the standard deviations of distance errors σ_{δ}^K also decrease when K increases. This results was expected also from Fig. 4 and Fig. 6, from which it is evident that increasing K not only reduces the distances between the position estimates and the TN, but it also reduces the distances between different position estimates because it alleviates the influence of acquisition errors. For the same choice of K , the values of the standard deviation evaluated using the TSML algorithm and using the PSO-based algorithm are similar and their order of magnitude is 0.2 m without range averaging and 0.02 m when $K = 100$.

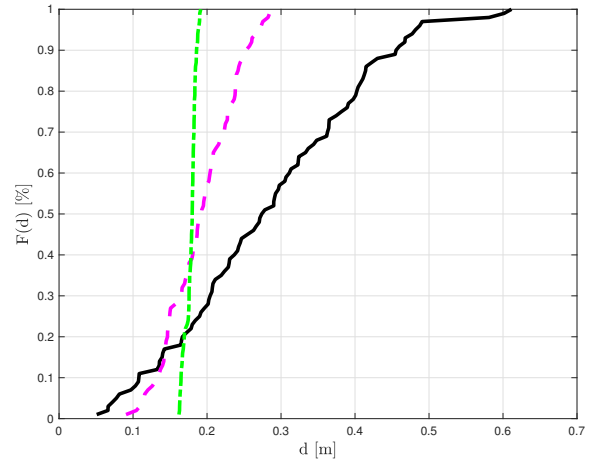


Fig. 8. The cumulative distribution function of distance errors from Fig. 7: (i) without range averaging (black line); (ii) with $K = 10$ (dashed magenta line), and (iii) with $K = 100$ (dash-dotted green line), relative to the position estimates of TN_2 obtained using the TSML algorithm.

B. Second Scenario

We now consider the TN positioned in the point denoted as TN_2 in Fig. 2, whose coordinates are denoted as \underline{u}_2 in (45). In this case, the true distances $\{\rho_i\}_{i=1}^3$ between the i -th AP and the TN are

$$\rho_1 \simeq 3 \text{ m} \quad \rho_2 \simeq 3 \text{ m} \quad \rho_3 \simeq 4.12 \text{ m}.$$

The projections of the range estimates on the plane $z = 1$ m can be evaluated, according to the Pythagorean theorem, as

$$r_1 \simeq 2.23 \text{ m} \quad r_2 \simeq 2.23 \text{ m} \quad r_3 \simeq 3.60 \text{ m}.$$

In order to estimate the position of TN_2 , range estimates from the APs are taken in order to have: 100 position estimates obtained without range averaging; 100 position estimates obtained by averaging over $K = 10$ consecutive range estimates; and 100 position estimates obtained by averaging over $K = 100$ consecutive range estimates.

In Fig. 7, the true position of TN_2 on the plane $z = 1$ (red star) is shown. Fig. 7 also shows the projections of the position estimates on the plane $z = 1$ (black circles) obtained without range averaging. Moreover, the projections of the position estimates (magenta crosses) obtained by averaging over 10 consecutive range estimates from each AP, and the projections of the position estimates (green triangles) obtained by averaging over 100 consecutive range estimates from each AP are shown. Fig. 7 shows that the distance among the TN position and its estimates decreases as K increases. This is motivated by the fact that large values of K lead to more precise averaged range estimates, which allow more accurate position estimates to be derived.

Fig. 8 shows the CDFs of the distance errors without range averaging (black line), with $K = 10$ (dashed magenta line), and with $K = 100$ (dash-dotted green line). As when considering TN_1 , the steepness of the CDF increases as K also increases.

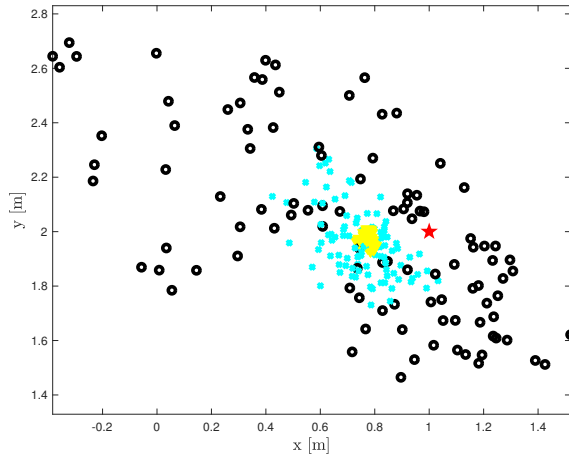


Fig. 9. Projections of the position estimates of TN_2 in Fig. 2 obtained using the PSO algorithm on the plane $z = 1$: (i) without range averaging (black circles); (ii) with $K = 10$ (cyan crosses); and (iii) with $K = 100$ (yellow triangles). The true TN position (red star) is also shown.

TABLE II. VALUES OF THE MAXIMUM DISTANCE ERROR (FIRST ROW), OF THE AVERAGE DISTANCE ERROR (SECOND ROW), AND OF THE STANDARD DEVIATION OF THE DISTANCE ERROR (THIRD ROW) ARE SHOWN, FOR DIFFERENT VALUES OF K , RELATIVE TO POSITION ESTIMATES OF TN_2 OBTAINED WITH THE TSML ALGORITHM AND WITH THE PSO-BASED ALGORITHM, RESPECTIVELY

TN_2	TSML			PSO		
	1	10	100	1	10	100
δ_{\max}^K [m]	0.61	0.29	0.19	1.53	0.59	0.27
δ_{avg}^K [m]	0.28	0.19	0.18	0.54	0.28	0.23
σ_{δ}^K [m]	0.13	0.05	0.008	0.35	0.10	0.01

Fig. 9 shows the true position of TN_2 on the plane $z = 1$ (red star), together with, the projections of the position estimates (black circles) obtained without range averaging, the projections of the position estimates (cyan crosses) obtained with $K = 10$, and the projections of the position estimates (yellow triangles) obtained with $K = 100$. As observed when using the TSML algorithm, the accuracy of the position estimates increases as K increases. Hence, also when considering the PSO-based algorithm, more accurate range estimates lead to more accurate position estimates.

Fig. 10 shows the CDFs of the distance errors without range averaging (black line), with $K = 10$ (dashed cyan line), with $K = 100$ (dash-dotted yellow line). As expected from Fig. 9, the larger is K , the steepest is the graph of the CDF.

Table II shows the values of the maximum distances and of the average distances between the considered TN and its estimates, and the values of the standard deviations of the distance errors for $K = 1$, $K = 10$, and $K = 100$. As when considering the first scenario, all such values decrease as the number of range averages K increases. Table II also shows that with no range estimate averaging, the maximum distance and the average distance between the considered TN and its estimates when using the PSO-based algorithm correspond to the double of the analogous values obtained with the TSML

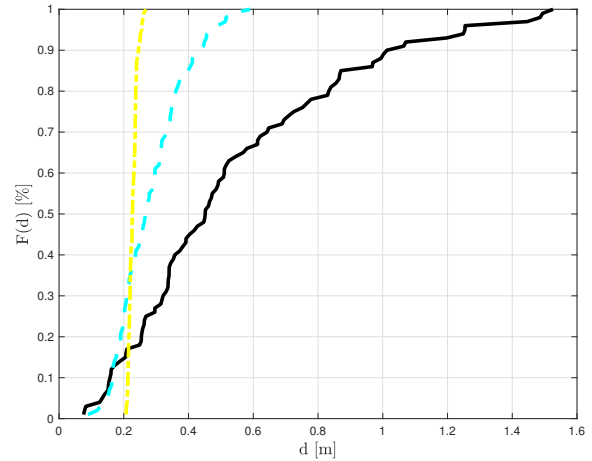


Fig. 10. The cumulative distribution function of distance errors from Fig. 9: (i) without range averaging (black line); (ii) with $K = 10$ (dashed cyan line), and (iii) with $K = 100$ (dash-dotted yellow line), relative to the position estimates of TN_2 obtained using the PSO algorithm.

algorithm. When $K = 10$ and $K = 100$, values obtained with the TSML algorithm are more similar to each other, even though the values obtained with the latter are slightly higher than those obtained with the former. Concerning the standard deviations of the distance error, it can be observed that, as in the first scenario, the values corresponding to $K = 100$ are one order of magnitude lower than those obtained when $K = 1$.

C. Third Scenario

Finally, let us now consider the TN positioned in the middle of the room, denoted as TN_3 in Fig. 2, whose coordinates are denoted as \underline{u}_3 in (45). In this case, the true distances $\{\rho_i\}_{i=1}^3$ between the i -th AP and the TN, expressed in meters, are

$$\rho_1 = \rho_2 = \rho_3 \simeq 3.46 \text{ m.}$$

Since the values of $\{h_i\}_{i=1}^3$ are 2 m, the projections of the distances on the plane $z = 1$ m can be computed as

$$r_1 = r_2 = r_3 \simeq 2.83 \text{ m.}$$

As in previous scenarios, range estimates from each of the three APs are acquired in order to have: 100 position estimates obtained without range averaging; 100 position estimates obtained by averaging over $K = 10$ consecutive range estimates; and 100 position estimates obtained by averaging over $K = 100$ consecutive range estimates.

Fig. 11 shows the projection of the true position of TN_3 on the plane $z = 1$ (red star). In the same figure, the projections of the position estimates (black circles) obtained without range averaging, the projections of the position estimates (magenta crosses) obtained by averaging over 10 consecutive range estimates from each AP, and the projections of the position estimates (cyan triangles) obtained by averaging over 100 consecutive range estimates from each AP, are also shown. As in the previous cases, Fig. 11 shows that the distance between position estimates and the TN decreases as K increases.

TABLE III. VALUES OF THE MAXIMUM DISTANCE ERROR (FIRST ROW), OF THE AVERAGE DISTANCE ERROR (SECOND ROW), AND OF THE STANDARD DEVIATION OF THE DISTANCE ERROR (THIRD ROW) ARE SHOWN, FOR DIFFERENT VALUES OF K , RELATIVE TO POSITION ESTIMATES OF TN_3 OBTAINED WITH THE TSML ALGORITHM AND WITH THE PSO-BASED ALGORITHM, RESPECTIVELY

TN_3	TSML			PSO		
	1	10	100	1	10	100
δ_{\max}^K [m]	0.64	0.31	0.16	0.88	0.47	0.33
δ_{avg}^K [m]	0.31	0.18	0.15	0.43	0.34	0.32
σ_{δ}^K [m]	0.15	0.05	0.006	0.23	0.05	0.004

Fig. 12 shows the CDFs of the distance errors without range averaging (black line), with $K = 10$ (dashed magenta line), and with $K = 100$ (dash-dotted green line), which is steeper as K increases.

Fig. 13 shows the true position of TN_3 on the plane $z = 1$ (red star), together with the projections of the position estimates on the plane $z = 1$ (black circles) obtained without range averaging, the projections of the position estimates (cyan crosses) obtained by averaging over 10 consecutive range estimates from each AP, and the projections of the position estimates (yellow triangles) obtained by averaging over 100 consecutive range estimates from each AP. As in previous scenarios, the position estimates are closer to each other and to the true TN position as K increases. Hence, also when considering the PSO-based algorithm, more accurate range estimates lead to more accurate position estimates.

Fig. 14 shows the CDFs of the distance errors without range averaging (black line), with $K = 10$ (dashed cyan line), and with $K = 100$ (dash-dotted yellow line). Once again, the larger is K , the steepest is the graph of the CDF.

Table III shows the values of the maximum distances and of the average distances between the considered TN and its estimates, and the values of the standard deviations of the distance errors for $K = 1$, $K = 10$, and $K = 100$.

As when considering previous scenarios, all such values decrease as the number of range averages K increases. Moreover, for the same choice of K values obtained with the TSML algorithm are of the same order of magnitude than those obtained with the PSO-based algorithm, even though the values obtained with the latter are slightly higher than those obtained with the former. Concerning the standard deviations of the distance error, it can be observed that, as in previous scenarios, the values corresponding to $K = 100$ are one order of magnitude lower than those obtained when $K = 1$. It is worth observing that, according to results shown in Tables I, II and III, the performance of both algorithms are similar in all scenarios.

VI. CONCLUSION

This paper presented an experimental evaluation of two approaches to indoor localization which both use ordinary WiFi signaling with no dedicated localization infrastructure. In both discussed approaches, agents acquire range estimates from the APs of the WiFi infrastructure, and they use such

estimates to obtain real-time information on the position of the smart appliances which hosts them.

The results obtained in the presented experimental campaign are meant to give a quantitative assessment on the performance of WiFi-based indoor localization, and they show that the level of accuracy of WiFi-based localization can be increased by a proper pre-processing of acquired range estimates. Obtained results show that agents can reach a localization accuracy of less than 1 m, thus making the proposed approach adequate for many application scenarios. In particular, such an accuracy is sufficient to support social games in large environments like shopping malls, waiting areas of airports and train stations, and covered markets in historical areas of towns. It is worth noting that presented results are valid under the assumption that the smart appliance does not move, or that it moves slowly with respect to range acquisition rate. Such an assumption does not necessarily hold for social games, and further investigation on dynamic scenarios is in progress.

REFERENCES

- [1] Z. Farid, R. Nordin, and M. Ismail, "Recent advances in wireless indoor localization techniques and system," *Journal of Computer Networks and Communications*, vol. 2013, 2013.
- [2] N. Bulusu, J. Heidemann, and D. Estrin, "GPS-less low cost outdoor localization for very small devices," *IEEE Personal Communications*, vol. 7, no. 5, pp. 28–34, October 2000.
- [3] H. Liu, H. Darabi, P. Banerjee, and J. Liu, "Survey of wireless indoor positioning techniques and systems," *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, vol. 37, no. 6, pp. 1067–1080, November 2007.
- [4] S. Monica and F. Bergenti, "A comparison of accurate indoor localization of static targets via WiFi and UWB ranging," in *Trends in Practical Applications of Scalable Multi-Agent Systems, the PAAMS Collection, 14th International Conference on Practical Applications of Agents and Multi-Agent Systems (PAAMS 2016), Special Session on Agents and Mobile Devices (AM)*, Sevilla, Spain, June 2016, pp. 111–123.
- [5] S. Gezici and H. V. Poor, "Position estimation via Ultra-Wide-Band signals," *Proceedings of the IEEE*, vol. 97, no. 2, pp. 386–403, February 2009.
- [6] F. Bellifemine, F. Bergenti, G. Caire, and A. Poggi, "JADE – A Java Agent DEvelopment framework," in *Multi-Agent Programming*. Springer, 2005, pp. 125–147.
- [7] G. Adorni, F. Bergenti, A. Poggi, and G. Rimassa, "Enabling FIPA agents on small devices," in *5th International Workshop on Cooperative Information Agents (CIA 2001)*, 2001, pp. 248–257.
- [8] F. Bergenti, G. Caire, and D. Gotta, "Agents on the move: JADE for Android devices," in *Proceedings of the 15th Workshop From Objects to Agents*, ser. CEUR Workshop Proceedings, vol. 1260. RWTH Aachen, 2014.
- [9] F. Bergenti, "An introduction to the JADEL programming language," in *Procs. of the IEEE 26th Int'l Conference on Tools with Artificial Intelligence (ICTAI)*, 2014, pp. 974–978.
- [10] F. Bergenti, E. Iotti, S. Monica, and A. Poggi, "Interaction protocols in the JADEL programming language," in *International Workshop on Programming based on Actors, Agents, and Decentralized Control (AGERE) at the ACM SIGPLAN conference on Systems, Programming, Languages, and Applications: Software for Humanity (SPLASH)*, 2016.
- [11] F. Bergenti, E. Iotti, S. Monica, and A. Poggi, "Agent-oriented model-driven development for JADE with the JADEL programming language," *Computer Languages, Systems and Structures*, vol. 50, pp. 142–158, 2017.
- [12] A. Poggi and F. Bergenti, "Developing smart emergency applications with multi-agent systems," *Int. J. E-Health Med. Commun.*, vol. 1, no. 4, pp. 1–13, 2010.

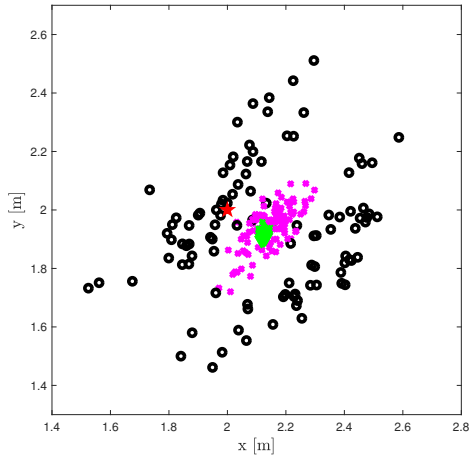


Fig. 11. Projections of the position estimates of TN_3 in Fig. 2 obtained using the TSML algorithm on the plane $z = 1$: (i) without range averaging (black circles); (ii) with $K = 10$ (magenta crosses); and (iii) with $K = 100$ (green triangles). The true TN position (red star) is also shown.

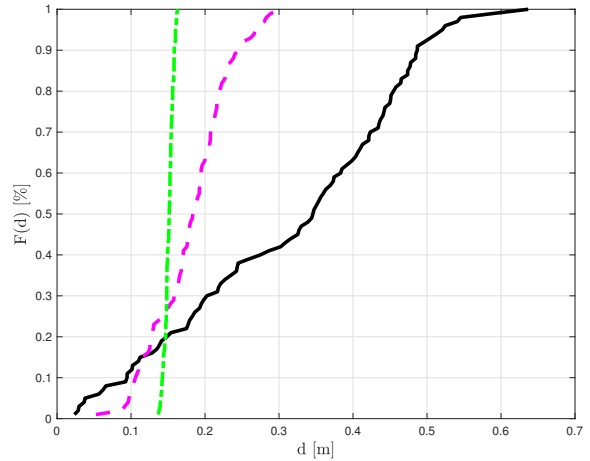


Fig. 12. The cumulative distribution function of distance errors from Fig. 11: (i) without range averaging (black line); (ii) with $K = 10$ (dashed magenta line), and (iii) with $K = 100$ (dash-dotted green line), relative to the position estimates of TN_3 obtained using the TSML algorithm.

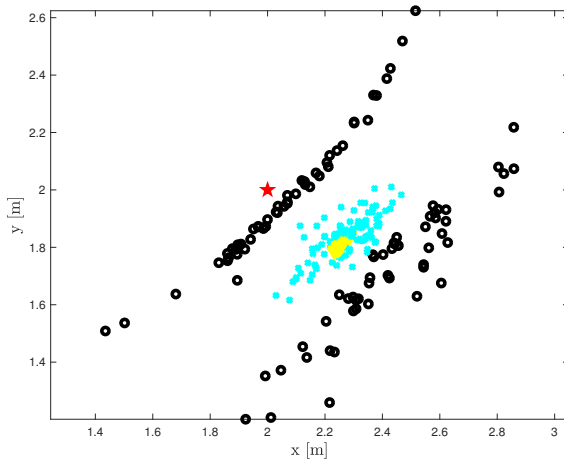


Fig. 13. Projections of the position estimates of TN_3 in Fig. 2 obtained using the PSO algorithm on the plane $z = 1$: (i) without range averaging (black circles); (ii) with $K = 10$ (cyan crosses); and (iii) with $K = 100$ (yellow triangles). The true TN position (red star) is also shown.

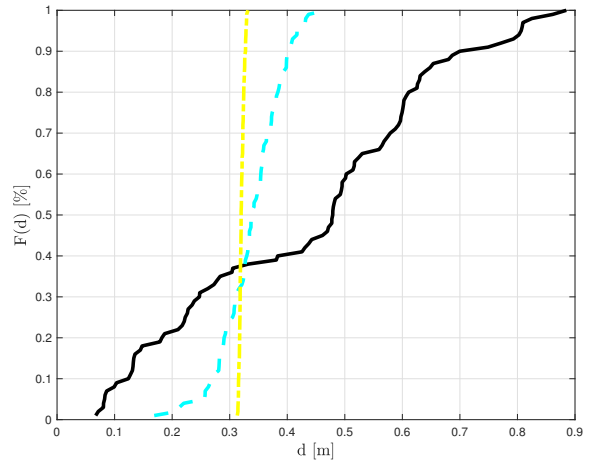


Fig. 14. The cumulative distribution function of distance errors from Fig. 13: (i) without range averaging (black line); (ii) with $K = 10$ (dashed cyan line), and (iii) with $K = 100$ (dash-dotted yellow line), relative to the position estimates of TN_3 obtained using the PSO algorithm.

[13] F. Bergenti, E. Franchi, and A. Poggi, "Agent-based social networks for enterprise collaboration," in *20th IEEE International Workshop on Enabling Technologies: Infrastructure for Collaborative Enterprises*, 2011, pp. 25–28.

[14] F. Bergenti, G. Caire, and D. Gotta, "Large-scale network and service management with WANTS," in *Industrial Agents: Emerging Applications of Software Agents in Industry*. Elsevier, 2015, pp. 231–246.

[15] F. Bergenti, G. Caire, and D. Gotta, "An overview of the AMUSE social gaming platform," in *Proceedings of the 14th Workshop From Objects to Agents*, ser. CEUR Workshop Proceedings, vol. 1099. RWTH Aachen, 2013.

[16] F. Bergenti and S. Monica, "Location-Aware Social Gaming with AMUSE," in *Trends in Practical Applications of Scalable Multi-Agent Systems, the PAAMS Collection, 14th International Conference on Practical Applications of Agents and Multi-Agent Systems (PAAMS 2016)*, ser. LNCS, vol. 9662, 2016, pp. 36–47.

[17] S. Monica and F. Bergenti, "Location-aware JADE agents in indoor scenarios," in *Proceedings of the 16th Workshop From Objects to*

Agents, ser. CEUR Workshop Proceedings, vol. 1382. RWTH Aachen, 2015, pp. 103–108.

[18] T. He, C. Huang, B. M. Blum, J. A. Stankovic, and T. Abdelzaher, "Range-free localization schemes for large scale sensor networks," in *Proceedings of the 9th Annual International Conference on Mobile Computing and Networking (MobiCom)*, San Diego, CA, September 2003, pp. 81–95.

[19] D. Dardari and R. D'Errico, "Passive ultrawide bandwidth RFID," in *Proceedings of the IEEE Global Telecommunications Conference (GLOBECOM '08)*, New Orleans, LA, December 2008, pp. 1–6.

[20] Z. Sahinoglu, S. Gezici, and I. Guvenc, *Ultra-wideband positioning systems: Theoretical limits, ranging algorithms and protocols*. Cambridge, U.K.: Cambridge University Press, 2008.

[21] N. Patwari, J. N. Ash, S. Kyperountas, A. O. Hero, R. L. Moses, and N. S. Correal, "Locating the nodes," *IEEE Signal Processing Magazine*, vol. 22, no. 4, pp. 54–69, July 2005.

[22] Y. Chan and K. C. Ho, "A simple and efficient estimator for hyperbolic location," *IEEE Transactions on Signal Processing*, vol. 42, no. 8, pp.

- 1905–1915, August 1994.
- [23] S. Monica and G. Ferrari, “A swarm-based approach to real-time 3D indoor localization: Experimental performance analysis,” *Applied Soft Computing*, vol. 43, pp. 489–497, June 2016.
- [24] S. Monica and F. Bergenti, “An experimental evaluation of agent-based indoor localization,” in *The Science and Information Computing Conference, technically sponsored by IEEE*. IEEE Press, July 2017, pp. 638–646.
- [25] S. Monica and F. Bergenti, “Experimental evaluation of agent-based localization of smart appliances,” in *Proceedings of the European Conference on Multi-Agent Systems (EUMAS 2016)*, ser. LNAI, vol. 10207, 2017.
- [26] S. Monica and G. Ferrari, “Accurate indoor localization with UWB wireless sensor networks,” in *Proceedings of the 23rd IEEE International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE 2014)*. IEEE Press, 2014, pp. 287–289.
- [27] K. C. Ho, X. Lu, and L. Kovavisaruch, “Source localization using TDOA and FDOA measurements in the presence of receiver location errors: Analysis and solution,” *IEEE Transactions on Signal Processing*, vol. 55, no. 2, pp. 684–696, February 2007.
- [28] S. Monica and G. Ferrari, “Swarm intelligent approaches to auto-localization of nodes in static UWB networks,” *Applied Soft Computing*, vol. 25, pp. 426–434, December 2014.
- [29] J. Kennedy and R. Eberhart, “Particle swarm optimization,” in *Proceedings of the IEEE International Conference on Neural Networks (ICNN)*, Perth, Australia, November 1995, pp. 1942–1948.
- [30] R. Poli, J. Kennedy, and T. Blackwell, “Particle swarm optimization,” *Swarm Intelligence Journal*, vol. 1, no. 1, pp. 33–57, June 2007.
- [31] Y. Shi and R. Eberhart, “A modified particle swarm optimizer,” in *Proceedings of the IEEE International Conference on Evolutionary Computation (ICEC)*, Washington, DC, July 1999, pp. 69–73.
- [32] R. Eberhart and J. Kennedy, “A new optimizer using particles swarm theory,” in *Proceedings of the 6th International Symposium on Micro Machine and Human Science (MHS)*, Nagoya, Japan, October 1995, pp. 39–43.
- [33] S. Monica and G. Ferrari, “Particle swarm optimization for auto-localization of nodes in wireless sensor networks,” in *Proceedings of the 11th International Conference on Adaptive and Natural Computing Algorithms (ICANNGA '13)*, ser. LNCS, vol. 7824, 2013, pp. 456–465.

Divide and Conquer Approach for Solving Security and Usability Conflict in User Authentication

Shah Zaman Nizamani
Department of IT
Quaid-e-Awam University
Nawabshah, Pakistan

Waqas Ali Sahito
Department of IT
Quaid-e-Awam University
Nawabshah, Pakistan

Shafique Awan
Department of Computer Science
Benazir Bhutto Shaheed University
Liyari, Pakistan

Abstract—Knowledge based authentication schemes are divided into textual password schemes and graphical password schemes. Textual password schemes are easy to use but have well known security issues, such as weak against online security attacks. Graphical password schemes are generally weak against shoulder surfing attacks. Usability is another issue with most of the graphical password schemes. For improving security of knowledge-based authentication schemes complex password entry procedures are used, which improve security but weakens useability of the authentication schemes. In order to resolve this security and usability conflict, a user authentication scheme is proposed, which contains one registration and two login screens called easy and secure login screens. Easy login screen provides easy and quick way of authentication while secure login screen is resilient to different online security attacks. A user has to decide based upon the authentication environment, which login screen to be used for authentication. For secure environment, where chances of security attacks are less easy login screen is recommended. For insecure environments where chances of security attacks are high, secure login screen is recommended for authentication. In the proposed scheme, image based passwords can also be set along with alphanumeric passwords. Results suggest that proposed scheme improves security against offline and online attacks.

Keywords—Authentication; alphanumeric passwords; security; passwords memorability

I. INTRODUCTION

Textual password scheme is easy to use because it has very simple password entry procedure. However, this scheme is weak in security because passwords can be recorded or observed from the login screen. Textual passwords can also be guessed through dictionary attacks because users mostly use dictionary words in their passwords [1]. In many applications some restrictions are enforced to set strong passwords such as minimum length of passwords. These restrictions does not fully resolve the issue of weak textual passwords because users still use dictionary words after applying the restrictions [2]. Complex or strong alphanumeric passwords are difficult to remember [3], therefore such passwords are not widely used. Strong textual passwords are difficult to guess from offline guessability attacks but they can be theft by observability and recordability attacks [4]. Another issue with textual passwords is that users generally set similar passwords in different accounts [5]. Due to this approach strong alphanumeric passwords can be guessed through offline guessability attacks after hacking a password from one user account [6].

Graphical passwords are used to solve security and memorability issues of textual passwords, but this technique has their

own set of problems, specially shoulder surfing and useability related issues [7]. Usability wise an authentication scheme is required to be easy to use, easy to learn and users' satisfaction need to be high with respect to performance and design. While, with respect to security an authentication scheme needs to provide enough resilience against different security attacks. Graphical password schemes lie in the range from secure and less useable to highly usable and less secure. It is because of conflicting nature between security and usability in user authentication schemes.

Different researchers improve the security of graphical password schemes by adding some logic in password entry techniques such as persuasive technique [8]. However, different usability or memorability issues arise due to inclusion of such logic because users have to complete multiple authentication steps or they need to provide large amount of information for authentication.

Security and usability parameters does not efficiently fit into one solution due to their conflicting nature [9]. Therefore, in user authentication schemes either security or usability sacrifices. Researchers generally give more importance to security because it is the most essential feature for an authentication scheme. In this research, both parameters are balanced by two login screens. First screen provides quick and easy way of password entry but it has some security weaknesses against online security attacks. While, other screen is resilient to online security attacks but it requires comparably more time for password entry. Users have option to authenticate with any of the login screen by using same password.

II. RELATED WORK

Passwords in knowledge based authentication schemes are alphanumeric or graphical. Alphanumeric or textual passwords are widely used for authentication but it has security and memorability issues. In order to overcome the issues of textual passwords, graphical passwords are proposed. Graphical password schemes are divided into pure recall based, cued recall based and recognition based schemes [10]. All the categories of graphical password schemes are discussed here in detail.

1) *Pure recall based schemes*: In this category of graphical password schemes, the passwords consist of some lines. Jermyn *et al.* [11] proposed a pure recall based graphical password scheme known as DAS (Draw-A-Secret). In this scheme, users draw some lines inside 2D grid-based login screen and the lines are considered to be the passwords of the users. In this category, passwords can be quickly inserted

but this category of authentication schemes has some security issues. For example, passwords can be easily viewed from login screen and dictionary attacks can also be applied.

Dunphy and Yan [12] proposed modified version of DAS scheme known as BDAS (Background DAS). In this scheme, background image is used inside 2D grid-based login screen. Background image helps the users to set complex passwords [13]. BDAS scheme is weak against shoulder surfing attack because passwords can be easily viewed from the login screen. Android unlock scheme [14] is widely used pure recall based graphical password scheme. In this scheme, nine points are shown in the login screen and the users have to create a password by connecting the points. In Android unlock scheme passwords are easy to enter but this scheme is weak against shoulder surfing attack.

2) *Cued recall based schemes*: In cued recall based graphical password schemes, passwords consist of some points inside a login screen. Blonder [15] proposed first cued recall based graphical password scheme in 1996. In this scheme, a password is created by selecting some predefined locations inside a picture. For authentication, a user needs to click on the locations which were selected as a password at the time of account registration. Password points are easy to select in Blonder's scheme but it has some security issues such as low password space and weak against shoulder surfing and guessability attack.

Wiedenbeck *et al.* [16] uses the idea of Blonder's scheme and proposed an authentication scheme known as Passpoint. In this scheme, users have no restriction of selecting password points inside the predefined locations. For authentication, a user just needs to click on the points which were selected at the time of password registration. This scheme requires short amount of time for authentication [13] but it is weak against shoulder surfing attack and guessability attacks [17].

3) *Recognition based scheme*: Images are used as password elements in recognition based schemes. For authentication, users have to correctly select their password images. Dhamija *et al.* [18] proposed a recognition based graphical password scheme known as Deja Vu, in which Abstract art images are used for password selection. A password is entered by clicking on the password images. Advantage of abstract art images is that they are difficult to guess by the attackers but such images are difficult to memorize. This scheme is also weak against shoulder surfing attacks [13].

CHC [19] is another recognition based graphical password scheme in which large number of icons are shown to the users for password selection. Users are authenticated when they correctly click on the logical triangles formed by the password icons. This scheme is resilient to many security attacks such as shoulder surfing and spyware attacks but it requires large amount of time for authentication. Therefore, usability is the main issue of CHC scheme.

Davis *et al.* [20] proposed another recognition based graphical password scheme known as story scheme. In this scheme, images of different categories are used for password selection. Idea of this scheme is that, users may create stories from password images and the stories will help in memorization of password images. This scheme has password memorability

advantage over CHC scheme [19] but this scheme is weak against shoulder surfing attack.

III. TOWARDS SOLUTION

Large number of knowledge based authentication schemes are proposed but easy to use authentication schemes are widely used such as traditional textual password scheme and Android unlock scheme. Relatively secure but difficult to use authentication schemes are not used. For example CHC scheme [19] provides a secure mechanism for authentication but it is not used for authentication due to difficult mechanism of password insertion.

Security and usability have conflicting nature in the field of user authentication, as a result one solution is difficult to design which equally resolves both the conflicting requirements of the authentication process. Therefore, in the proposed scheme users have been provided two options for authentication, one is secure and other is easy to use. Depending upon the login environment, the users can authenticate by any of the login options. Chances of password hacking increases when a user authenticates inside an office or public network, in such environments secure login approach is recommended. While, easy login approach is recommended for private networks such as home. When a user mistakenly authenticate through easy login approach in the insecure environment then same level of password security is achieved which is present in traditional textual password scheme.

IV. PROPOSED SCHEME

In the proposed scheme, one registration and two login screens are designed as shown in Fig. 1. Multiple login screens are designed to solve security and usability conflict in user authentication. In the first login screen (called easy login screen) users just need to type the password elements similar to textual password scheme. This login screen provides easy to use password entry procedure for better usability. Second login screen (called secure login screen) is designed for better security. In the secure login screen, 50% elements (alphanumeric characters and images) are presented. For password entry, users need to enter the count of visible password elements inside the password field. This login screen is resilient to many online security attacks such as keylogger attacks. This login screen provides security advantages but it requires relatively more time for password entry than easy login screen. Therefore, easy login screen is recommended for authentication in secure login environments and the secure login screen is recommended for insecure login environments.

The proposed scheme contains both textual and graphical elements for password selection. A password may be consist of alphanumeric characters, images or combination of both. Users need to remember single password for both the login screens.

A. Registration Activity

In the registration activity, a user inserts profile and authentication information for account creation. Registration screen of the proposed scheme presents twenty four images for password selection as shown in Fig. 2. All the images are selected from the categories of fruits, electronics, birds,

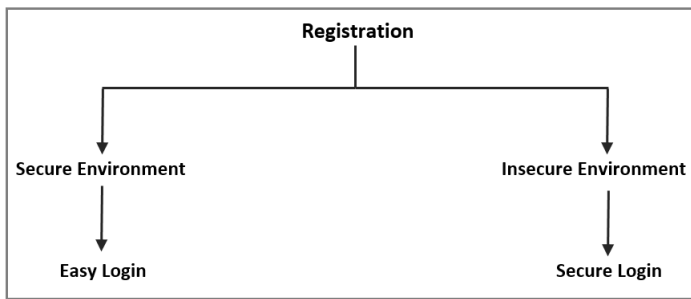


Fig. 1. Proposed solution.

emoji and animals. Four images are selected from each of the category, based upon familiarity among the users.

1) *Password entry*: Alphanumeric characters of a password are selected by typing the keys of the characters and a password image is selected by typing the shortcut key, which is the combination of “control” and “alter” keys along with two initial characters of an image. For example, if password is “abc” and image of “horse”, then the password is selected by typing “abc” and pressing “ctrl+alt+h” keys altogether inside the password field. In the database, some Unicode symbols are saved against the password images. For example, in current scenario the image of horse may be represented in the database by Unicode symbol β . The Unicode symbols for the images are not fixed they can be changed in every deployment of the scheme. For improving security against dictionary attacks, it is better to use different Unicode symbols in each deployment of the proposed scheme.

B. Login Activity

In the proposed scheme, login activity can be completed by any of the two login screens. First screen is called “easy login screen” and second is called “secure login screen”. Password entry process is different in both the login screens, which is explained here.

1) *Easy login screen*: This login screen is almost similar to registration screen as shown in Fig. 3, only fields for inserting profile information are not presented. Password is entered in the easy login screen, similar to the registration screen. A user can insert alphanumeric characters by pressing the keys related with the alphanumeric characters and the images are selected by typing the shortcut keys of the images. After inserting username and password, a user just needs to click on login button for authentication. The user will be authenticated once username and password matches with the stored authentication information.

Easy login screen is designed for quick and easy authentication. Login process of easy login screen is similar to textual password scheme, therefore learnability is not an issue with this login screen. However, security attacks such as keylogger and spyware attacks may work in the easy login screen due to exact insertion of password elements inside the password field. These security weaknesses are intentional because this login screen is designed for better usability. Easy login screen has same security issues as in textual password scheme against online security attacks as same password entry procedure is

used. However, this login screen is better in offline guessability attacks because image based passwords can also be selected.

2) *Secure login screen*: This login screen is designed to resist online secure attacks. All of the attacks are resisted by indirectly gathering passwords from the users. In the secure login screen randomly 50% elements are presented. For authentication a temporary number is inserted into the password field and the number depends upon the visibility of the elements in the secure login screen. The password number changes in every login session, therefore this login screen resist online security attacks.

a) *Password entry*: In the secure login screen as shown in Fig. 4 and 5, randomly 59 out of 118 (50%) elements are presented. The 59 elements (47 alphanumeric characters and 12 images) are randomly selected but they are shown in natural order. A user has to count the password elements currently visible inside the secure login screen. The count of visible password elements is then inserted into the password field for authentication. For example, if password of a user is “abcde” and an image of “horse” then based upon the login screens as shown in Fig. 4 and 5, the user has to insert “4” in the password field because character “a” and the image of “horse” is not visible in the login screen. If a password contains same element multiple times then the element will be counted more than once. For example, if the password is “abccde” and login screen is same as shown in Fig. 4, then the user has to enter decimal number “5” inside the password field because the character “c” is presented two times within the password.

Login process in secure login screen consists of three steps. In each step different arrangements of alphanumeric characters and images are shown. A user has to count and enter visible password elements for all the three steps. Each step appears by clicking on the tab “Step” as shown at the top of Fig. 4 and 5. Multiple steps are created for reducing the chances of blind guessing attack. The login steps can be increased for improving security against blind guessing attack but it will require more time for authentication.

b) *Password matching*: In secure login, password numbers are compared for authentication instead of actual password elements. The server generates a password number which is compared with the user’s provided password number. When both the numbers become equal then the server allows sign-in. Steps for authentication in secure login are listed below.

- (i) Username and a password number are received by an authentication server.
- (ii) If username does not present in the database then the server will close authentication process.
- (iii) If username present in the database then the server fetches and decrypt the password based upon the provided username.
- (iv) Server gets the session variable which stores all the visible elements for the current login session.
- (v) Server counts all the decrypted password elements inside the session variable.
- (vi) Server compares the count generated in step V with the password number given by the user.
- (vii) A user is authenticated when both the user’s provided number and the system generated count is equal in all the three steps.

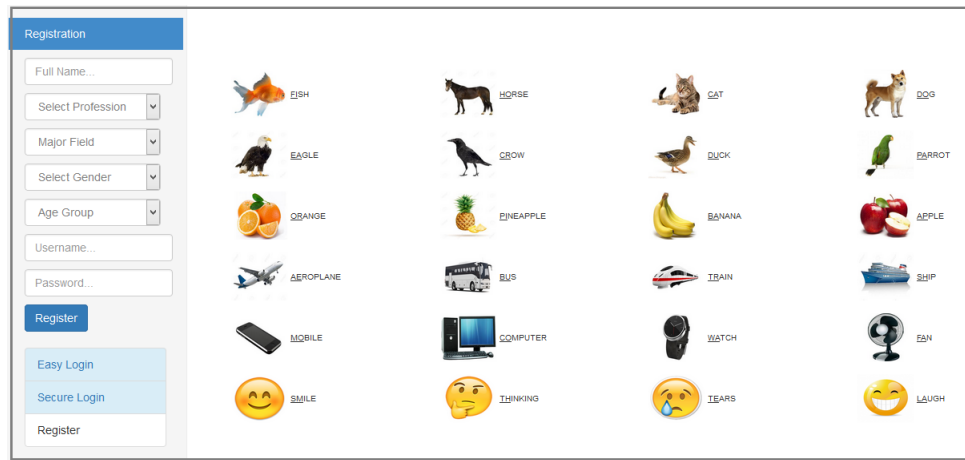


Fig. 2. Registration screen.

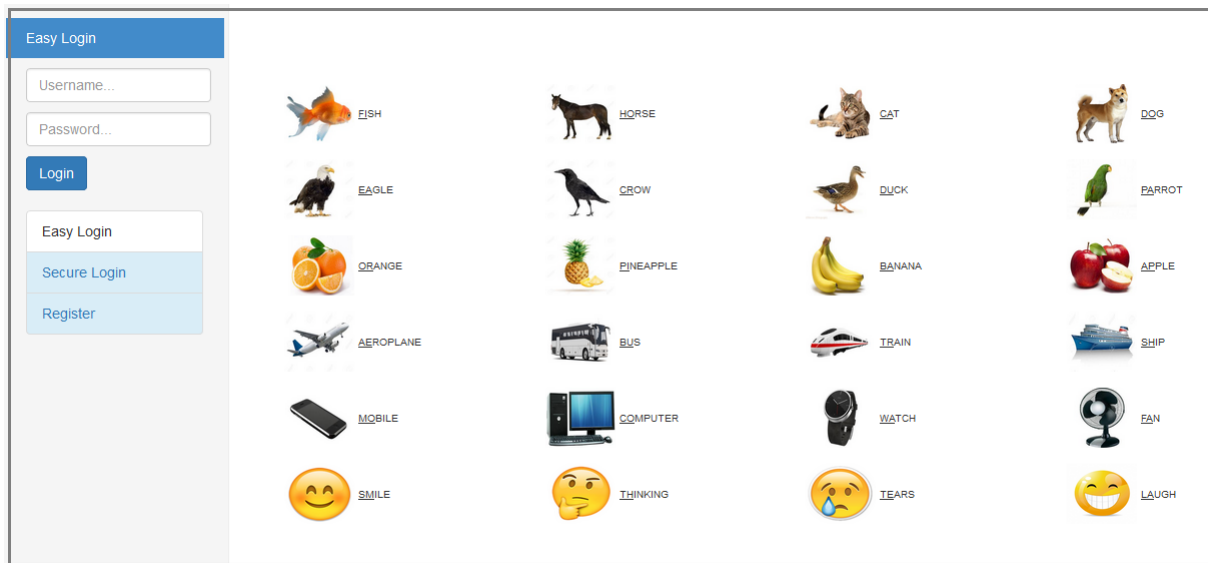


Fig. 3. Easy login screen.

Due to indirect insertion of passwords, the secure login process requires the passwords to be stored in two way encryption. Two way encryption is relatively weak than hashing or one way encryption. Therefore, passwords for the proposed scheme need to be secured with different techniques such as differential masking [21].

V. USABILITY AND MEMORABILITY ANALYSIS

In order to analyze usability and memorability aspects of the proposed scheme, a web based application was developed. A hidden process was created inside the application for calculating timings of registration and login activities. A log was also maintained for analyzing failed and successful login attempts.

For testing purpose, 50 participants were selected from different departments of Quaid-E-Awam university of Engineering Science and technology, Pakistan. Professionally the participants were students, teachers and administrative staff. The users were selected based upon their knowledge about

computer usage. All the users had basic knowledge about internet and its working. All the participants performed the registration and login activities inside the testing application.

A. Testing Procedure

Testing phase for the proposed scheme was consist of four sessions. In first session, users performed registration and login activities. While in remaining three sessions, users only performed login activities. Before conducting the tests, a demonstration was given for creating user account and sign-in through the testing application. When users fully understood the login and registration activities, then the users were asked to perform the activities inside the application.

Users were free to create any password (alphanumeric or image-based). Minimum length of passwords were set to eight elements and users were also asked to set passwords from at least two categories such as numbers and special characters. Second session was started after one day of registration, in this session the users only performed login activities with their

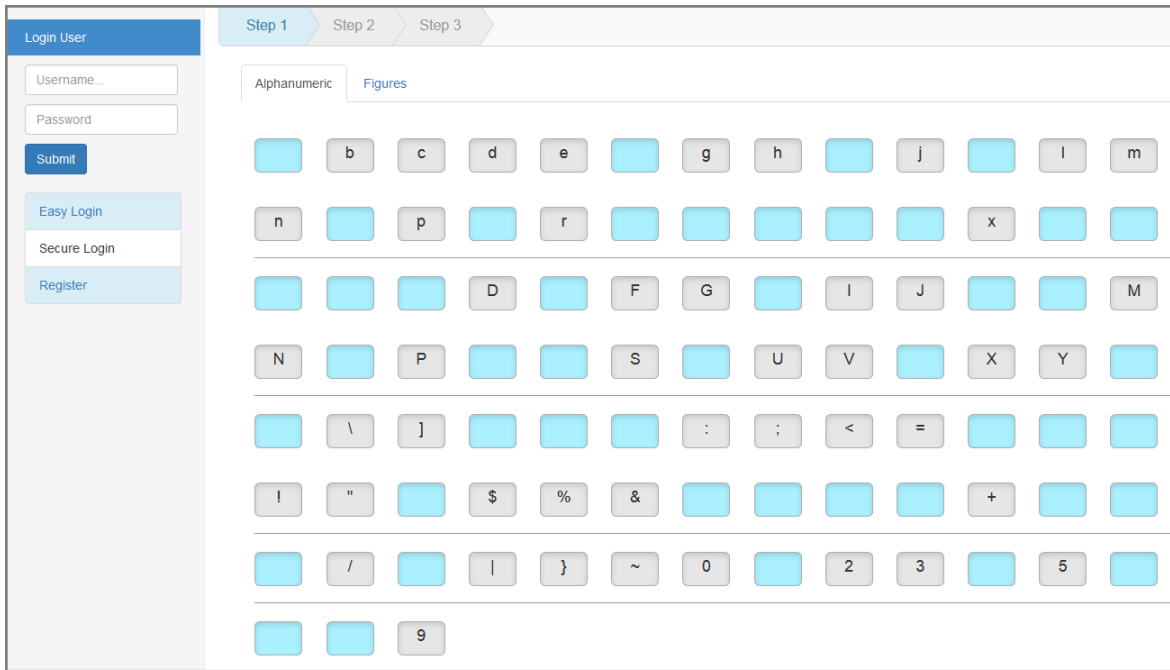


Fig. 4. Secure login screen showing alphanumeric characters.

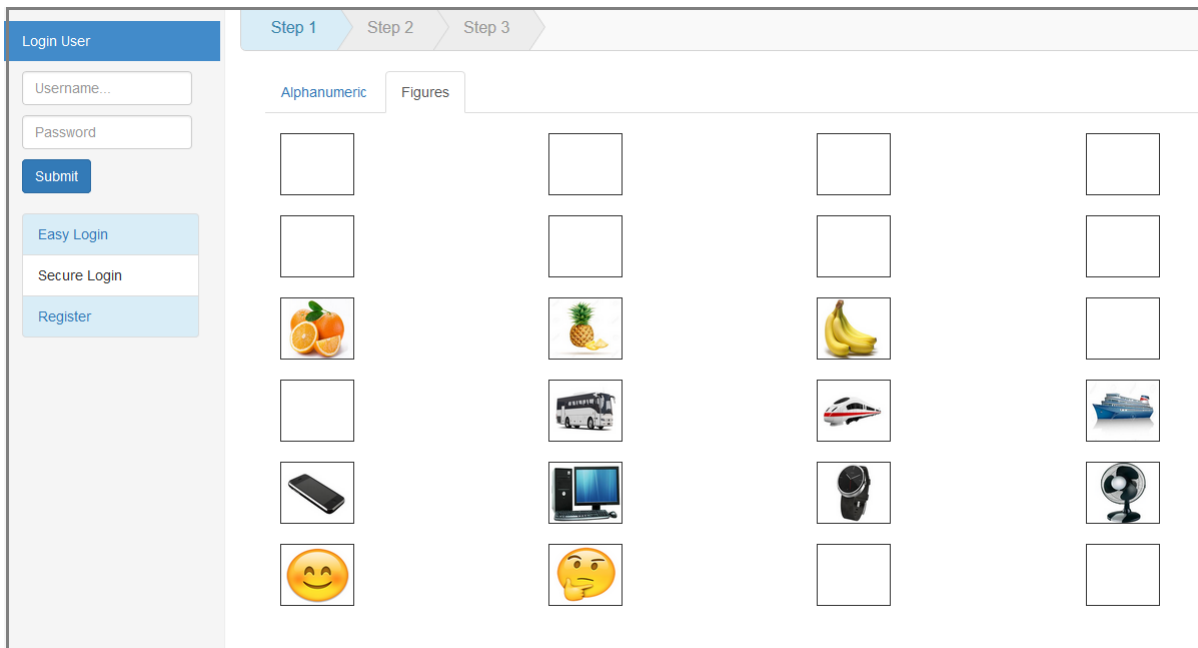


Fig. 5. Secure login screen showing images.

registered username and passwords. Third and fourth sessions were started after one and two weeks of registration, respectively. In both the sessions only those users were asked to login who have successfully authenticated in previous sessions.

B. Testing Results

The experiment data was analyzed to get the performance of the proposed scheme with respect to usability and memorability. Results show that passwords were mostly consist of alphanumeric characters. Out of 50 participants, thirteen users used combination of alphanumeric and image based passwords and three users selected passwords with only images. This behaviour was due to wide use of traditional textual passwords among the users.

Mostly users authenticated in first login attempt when they have remembered their passwords. The results of failed login attempts show that users had no difficulty in authentication inside the proposed scheme.

Authentication timing is the main factor for analyzing usability of an authentication scheme. Testing results showed that the proposed scheme requires 15.82 seconds for password registration or selecting the password elements. Average login time was 11.86 seconds in easy login screen and 32.73 seconds in secure login screen. Large time requires in secure login screen due to three steps of authentication.

The users who used only alphanumeric characters in their passwords took less time in easy login screen in comparison with the users who used images in their passwords. Image based passwords require more time for password insertion due to usage of shortcut keys. The users who set alphanumeric passwords took more time in secure login screen as compared to image based passwords. This behaviour is due to the more effort requires for searching the alphanumeric characters in secure login screen. Average length of the passwords were found 8.9 and average password entropy was found 55 bits.

Memorability tests were conducted immediately after registration, one day, one week and two weeks. The results are shown in Table I.

TABLE I. PASSWORD MEMORABILITY IN THE SCHEME

Duration	Password memorability
After registration	94%
After 1 Day	86%
After 1 Week	72%
After 2 Weeks	62%

VI. SECURITY ANALYSIS

Table II shows the status of different security attacks against textual password scheme, Android unlock scheme and proposed authentication scheme. In Table II value “Y” shows that the login screen is resilient to the particular attack, while the value N shows that the screen is not resilient to the attack. In the table value “Hard” shows that a very high level of effort is required to crack the password.

In the proposed scheme, 118 elements (alphanumeric characters and images) are used, while only 95 elements are used in traditional textual passwords scheme based upon American standard keyboard. Therefore, the proposed scheme provides

more password space than textual password scheme. Higher password space is better for security against brute force attack because an attacker needs to apply large number of combinations for password crack.

The proposed scheme also performs better with respect to dictionary attacks because users have option to select images along with alphanumeric characters. Due to the inclusion of images, password dictionaries are difficult to create for the proposed scheme. The attackers have to create the list of passwords with the combination of alphanumeric characters and images and they have to identify the Unicode symbols used for the images.

Shoulder surfing and spyware attacks can be applied in the easy login screen but these attacks are resisted in secure login screen. Passwords are indirectly inserted into the secure login screen, therefore the passwords can not be captured by applying shoulder surfing or spyware attacks in the secure login screen.

Man in the Middle attack depends upon the implementation of the proposed scheme, this attack can be resisted if secure communication channel is used for easy login screen such as SSL or TLS [22] [23]. While, secure login screen is resilient to this attack because users enter temporary numbers instead of actual passwords in the secure login screen.

In multiple recording attack, passwords are captured by recording information of multiple login sessions. In easy login screen, recording of single login session is enough for password heck because a user enters original password elements into the password field. In the secure login screen, different numbers are entered into the password field instead of exact password elements, therefore passwords can not be captured from multiple recording attack.

In blind guessing attack, an attacker randomly enters different passwords into the password field for authentication. In easy login screen this attack does not work because an attacker needs to apply very large number of passwords for authentication, which is not manually possible. In secure login screen, it is possible that an attacker enters three numbers which are the password for current login session. Chances of this threat can be reduced by deactivating the authentication process after three failed login attempts.

Table II shows that easy login screen is weak with respect to client side security attacks but it improves security against offline gessability attacks. While, secure login screen is resilient to most of the security attacks.

VII. CONCLUSION

The proposed scheme does not replaces the traditional textual password scheme but it enhances the security of textual password scheme in terms of password entry procedure and list of password elements. Therefore, proposed scheme can be easily deployed in existing applications which use the textual password scheme for authentication.

Blind guessing attack may work in the secure login screen. Chances of this attack can be reduced by increasing steps inside the secure login screen and adding policies such as allowing only three attempts for authentication.

TABLE II. STATUS OF SECURITY ATTACKS

Scheme	Brute Force	Dictionary	Shoulder Surfing	Spyware	Man in the Middle	Multiple Recording	Blind Guessing
Textual Password	Hard	N	Medium	N	N	N	Y
Android Unlock	N	N	N	N	N	N	Y
Easy login Screen	Y	Y	Medium	N	N	N	Y
Secure Login Screen	Y	Y	Y	Y	Y	N	Hard

In the proposed scheme, number of images can be increased to enhance the password space of the scheme. However, increment should be such that it should not affect the usability or memorability of the scheme. Large number of images can create messy look and feel of the authentication screens and users will face difficulty in finding their password images. Memorability may also be effected by adding or replacing the images with complex images, i.e. the images with less cues for password memorization.

REFERENCES

- [1] J. Yan, A. Blackwell, R. Anderson, and A. Grant, "Password memorability and security: Empirical results," *IEEE Security & privacy*, vol. 2, no. 5, pp. 25–31, 2004.
- [2] L. Tam, M. Glassman, and M. Vandenwauver, "The psychology of password management: a tradeoff between security and convenience," *Behaviour & Information Technology*, vol. 29, no. 3, pp. 233–244, 2010.
- [3] L. F. Cranor and S. Garfinkel, *Security and usability: designing secure systems that people can use*. O'Reilly Media, Inc., 2005.
- [4] R. English and R. Poet, "Towards a metric for recognition-based graphical password security," in *Network and System Security (NSS), 2011 5th International Conference on*. IEEE, 2011, pp. 239–243.
- [5] A. Das, J. Bonneau, M. Caesar, N. Borisov, and X. Wang, "The tangled web of password reuse," in *NDSS*, vol. 14, 2014, pp. 23–26.
- [6] A. S. Brown, E. Bracken, S. Zoccoli, and K. Douglas, "Generating and remembering passwords," *Applied Cognitive Psychology*, vol. 18, no. 6, pp. 641–651, 2004.
- [7] A. M. Eljetlawi and N. Ithnin, "Graphical password: Comprehensive study of the usability features of the recognition base graphical password methods," in *Convergence and Hybrid Information Technology, 2008. ICCIT'08. Third International Conference on*, vol. 2. IEEE, 2008, pp. 1137–1143.
- [8] A. Forget, S. Chiasson, P. C. van Oorschot, and R. Biddle, "Persuasion for stronger passwords: Motivation and pilot study," in *International Conference on Persuasive Technology*. Springer, 2008, pp. 140–150.
- [9] Y. Meng, "Designing click-draw based graphical password scheme for better authentication," in *Networking, Architecture and Storage (NAS), 2012 IEEE 7th International Conference on*. IEEE, 2012, pp. 39–48.
- [10] A. De Angeli, L. Coventry, G. Johnson, and K. Renaud, "Is a picture really worth a thousand words? exploring the feasibility of graphical authentication systems," *International Journal of Human-Computer Studies*, vol. 63, no. 1, pp. 128–152, 2005.
- [11] I. Jermyn, A. J. Mayer, F. Monrose, M. K. Reiter, A. D. Rubin *et al.*, "The design and analysis of graphical passwords," in *Usenix Security*, 1999.
- [12] P. Dunphy and J. Yan, "Do background images improve draw a secret graphical passwords?" in *Proceedings of the 14th ACM conference on Computer and communications security*. ACM, 2007, pp. 36–47.
- [13] R. Biddle, S. Chiasson, and P. C. Van Oorschot, "Graphical passwords: Learning from the first twelve years," *ACM Computing Surveys (CSUR)*, vol. 44, no. 4, p. 19, 2012.
- [14] S. Uellenbeck, M. Drmuth, C. Wolf, and T. Holz, "Quantifying the security of graphical passwords: the case of android unlock patterns," in *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*. ACM, 2013, pp. 161–172.
- [15] G. E. Blonder, "Graphical password," Sep. 24 1996, uS Patent 5,559,961.
- [16] S. Wiedenbeck, J. Waters, J.-C. Birget, A. Brodskiy, and N. Memon, "Passpoints: Design and longitudinal evaluation of a graphical password system," *International Journal of Human-Computer Studies*, vol. 63, no. 1, pp. 102–127, 2005.
- [17] S. Chiasson, A. Forget, R. Biddle, and P. C. van Oorschot, "User interface design affects security: Patterns in click-based graphical passwords," *International Journal of Information Security*, vol. 8, no. 6, pp. 387–398, 2009.
- [18] R. Dhamija and A. Perrig, "Deja vu-a user study: Using images for authentication," in *USENIX Security Symposium*, vol. 9, 2000, pp. 4–4.
- [19] S. Wiedenbeck, J. Waters, L. Sobrado, and J.-C. Birget, "Design and evaluation of a shoulder-surfing resistant graphical password scheme," in *Proceedings of the working conference on Advanced visual interfaces*. ACM, 2006, pp. 177–184.
- [20] D. Davis, F. Monrose, and M. K. Reiter, "On user choice in graphical password schemes," in *USENIX Security Symposium*, vol. 13, 2004, pp. 11–11.
- [21] S. Z. Nizamani, S. R. Hassan, and R. Naz, "A theoretical framework for password security against offline guessability attacks," *Indian Journal of Science and Technology*, vol. 10, no. 33, 2017.
- [22] A. Freier, P. Karlton, and P. Kocher, "The secure sockets layer (ssl) protocol version 3.0," 2011.
- [23] T. Dierks, "The transport layer security (tls) protocol version 1.2," 2008.

A High-Performing Similarity Measure for Categorical Dataset with SF-Tree Clustering Algorithm

Mahmoud A. Mahdi

Faculty of Computers and Information
Zagazig University, Egypt

Samir E. Abdelrahman

Faculty of Computers and Information
Cairo University, Egypt

Department of Biomedical Informatics
University of Utah, USA

Reem Bahgat

Faculty of Computers and Information
Cairo University, Egypt

Abstract—Tasks such as clustering and classification assume the existence of a similarity measure to assess the similarity (or dissimilarity) of a pair of observations or clusters. The key difference between most clustering methods is in their similarity measures. This article proposes a new similarity measure function called PWO “Probability of the Weights between Overlapped items” which could be used in clustering categorical dataset; proves that PWO is a metric; presents a framework implementation to detect the best similarity value for different datasets; and improves the F-tree clustering algorithm with Semi-supervised method to refine the results. The experimental evaluation on real categorical datasets, such as “Mushrooms, KrVskp, Congressional Voting, Soybean-Large, Soybean-Small, Hepatitis, Zoo, Lenses, and Adult-Stroke” shows that PWO is more effective in measuring the similarity between categorical data than state-of-the-art algorithms; clustering based on PWO with pre-defined number of clusters results a good separation of classes with a high purity of average 80% coverage of real classes; and the overlap estimator perfectly estimates the value of the overlap threshold using a small sample of dataset of around 5% of data size.

Keywords—Algorithm; clustering; similarity; measurement; categorical; F-Tree; SF-Tree

I. INTRODUCTION

General data mining applications have two types of data, categorical and numerical. Most clustering algorithms focus on numerical data whose inherent geometric properties can be exploited naturally to define distance functions between data points [1]. Categorical data refers to the data describing objects, which have only categorical (non-numerical) attributes [2]. Such data is often related to transactions involving a finite set of elements, or items, in a common item universe [3]. Transactional data is a kind of categorical data in which records can have different sizes. It is generated by many applications such as e-commerce, healthcare, and CRM [4]. It plays an important role in many fields like market basket data, web usage data, customer profiles, patient symptoms’ records, and image features. This paper focuses categorical and transactional data.

Clustering is a widely used technique in which data items are partitioned into groups (called clusters) based on their similarities or differences, such that data items in the same cluster are more similar among themselves than items in other

clusters [5]. It is usually difficult to deal with categorical attributes; therefore, clustering of categorical attributes has not received as much attention as its numerical counterpart [6]. Categorical attributes have unique features from the definition in [2]; therefore, the traditional approach to convert categorical data into numerical values does not necessarily produce meaningful results specially in the case where categorical domains are not sorted [2], [7], [8]. For example, hierarchical clustering algorithms may be unstable when used to cluster categorical data because the distance between the centroid of clusters of categorical data is not a good estimator of the similarity between the data [9]. Partition clustering algorithms may also be unsuitable because the sets of items that define clusters may not have the same sizes since the cluster may contain a small subset of the possible number of items. Thus, it is possible that a pair of transactions in a cluster have few items in common [7]. Moreover, clustering categorical data involve complexity that is not encountered in numerical data. In addition, different clustering algorithms hardly generate the same clustering result for the same dataset. For these reasons, there is an unmet need for algorithms that tackle these limitations during clustering categorical data [6].

One of the most important aspects of data mining problem is how similarity measure is defined [10] and calculated [11], since the similarity measures have the effect of clustering and classifying information with respect to data types. Clustering techniques for categorical data are very different from those for numerical data in terms of the definition of similarity measure [12]. It is also rare to find the boundaries of the clusters and avoid overlapping between them, which adds an additional constraint to researchers when choosing the optimal similarity measure that could be applied to a wide range of data types. Most of the clustering algorithms have two phases: allocation and refinement phases. The refinement phase has two drawbacks: 1) its results depend on the results of the allocation phase; and 2) its run time complexity is relatively high. It is known that the size of transactional data is usually large, so there is a great demand for fast and high quality algorithms to cluster large-scale transactional datasets.

This article extends our prior study of measuring the similarity between clusters of categorical (or transactional) data in [13]. The list of this article contributions are presented in (Tables I and II) and are summarized as follows.

A criterion function is described in details for similarity measure called PWO (Probability of the Weights of Overlapped items) for categorical data to overcome the problem of overlapping between clusters. A new algorithm which depends on PWO is provided for clustering categorical datasets with possibly different dimensions. A new framework is proposed to estimate the best similarity threshold parameters for different datasets that could be used as the proposed clustering algorithm. The similarity measure is tested on real-world datasets obtained from the UCI Machine Learning Repository [14] and applied to find similar groups in models constructed from different datasets. Inferences from the similar groups found to be logically meaningful. Finally, the algorithm is also compared versus different state-of-the-art algorithms in terms of the purity, the number of clusters, and the performance.

To sum up, this article extends significantly the earlier work [13] in the aspects described in Tables (I, and II). Here, the PWO is discussed in details, propose the PWO as a stand-alone algorithm, improve F-Tree algorithm, and implement the overlap estimator algorithm. The algorithms are evaluated in details versus different algorithms using addition datasets.

The reset of this paper is organized as follows. Section 2 summarizes the general notation and definition. Section 3 discusses the related work to this paper. Section 4 discusses the PWO similarity measure that is used to calculate the similarity between clusters. Section 5 discusses the approach to cluster categorical data based on PWO similarity measure. Section 6 presents the overlap estimator framework to determine the overlap threshold. Section 7 describes F-Tree clustering algorithm. Sections 8 and 9 describe the data mentioned in this research followed by a comprehensive set of experiments and related discussions. Sections 10 and 11 present the limitations and conclusion of this study.

II. NOTATION AND DEFINITION

In order to simplify the expressions throughout this paper, the following notations are used. Consider a categorical or transactional dataset D consisting of a set of transactions $\{t_1, t_2, \dots, t_n\}$ of size N . where, each transaction T contains a set of items or attributes $I = \{i_1, i_2, \dots, i_m\}$. Hence, Clustering $\{C_1, C_2, \dots, C_k\}$ is a partition of transactions $\{t_1, t_2, \dots, t_n\}$. Where, each C_k called a cluster and K is the total number of clusters. M_k , and N_k are used respectively to denote the number of distinct items, and the number of transactions in the cluster C_k . I_k represents the categorical items in a cluster C_k , where $I_k = \{i_{k1}, i_{k2}, \dots, i_{kM}\}$. S_k is the sum of occurrences of all items in cluster C_k . Θ is the minimum support or the minimum number of item's occurrence that should be present in each cluster.

III. RELATED WORK

The recent categorical clustering techniques are reviewed in this section. Each algorithm follows one concept of the three main concepts. First, clustering algorithms based on a predefined knowledge of the number of clusters such as COOLCAT [15], LIMBO [16], Fast clustering [17], Ensemble [18], and Hybrid [19]. Second, clustering algorithms without any knowledge about the clusters such as LargeItem [20], SLR [21], SEED [22], CACTUS [23], CLOPE [24], CLICKS [25]

TABLE I. EXTENDED EFFORTS

	Prior paper [13]	This article contribution
PWO Measure	Summary	More focus and metric proof
PWO Algorithm		Novel
Overlap Estimator	Proposed idea	implementation
F-Tree Algorithm	More description	Summary, and add predefined number of clusters
SF-Tree Algorithm		Novel

TABLE II. EXTENDED EXPERIMENTS

	Prior paper [13]	This article contribution
PWO Measure		Evaluation of metric function
PWO Algorithm	Minimum support vs number of clusters	Evaluation and analysis clustering algorithm with(out) fixed number of clusters
Overlap Estimator		Analysis precision and scalability
F-Tree Algorithm	Compare with 4 algorithms	Compare with more than 10 other algorithms
SF-Tree Algorithm		Analysis with predefined number of clusters, without predefined number of clusters, and analysis minimum fit of training dataset
Dataset	Mushroom and Votes	Nine Datasets (Table IV)

and DELTA [26]. The last type includes clustering algorithms that depend on the number of clusters at further step in order to refine and improve the clustering or to have an ability to work in the first place, such as ROCK [7], WCD [4], Squeezer [1], and SCCADDS [27]. In addition, it have been found some authors presented many techniques to find the best number of clusters.

Most algorithms generate clusters in the allocation phase then try to refine them in the refinement phase depending on the similarity measure function. The numbers of refinement steps are then state, as they affect the algorithm's performance. Measure function is applied on either local clusters or global clusters or both. Approaches based on local function compute the evaluation function between items inside the same cluster; the result shows the degree of how items inside a cluster are related to each other. On the other hand, global approaches compute the evaluation function between clusters; the result shows the degree of how clusters are dissimilar and more distinct. Finally, the measurement parameters and their numbers are stated; increasing the number of parameters will increase the complexity of the algorithm and the difficulty of the user's experience.

The LargeItem [20] uses the concept of large items to divide the transactions into clusters. An item marked as large in a cluster of transactions if its occurrence rate is larger than a minimum support parameter that is specified by the user. The LargeItem approach scans each transaction and either allocates it to an existing cluster or assigns it to a new cluster based on a cost function. The process of choosing a cluster for each transaction is based on the global goodness of clustering. This goodness is measured by minimizing the total cost function. Therefore, the LargeItem algorithm needs to set two parameters the minimum support Θ and the large item factor or weight w . In addition, the LargeItem algorithm is exhaustive in the decision procedure of moving a transaction t to the best cluster. The data structure used to handle clusters

is complex, and the approach taken to update the criterion function is not efficient; although the implementation uses the B-Tree structure to increase the performance of updating, but it consumes a lot of memory in case of handling the large dimensions of a dataset with large number of attributes. Moreover, the procedure of scanning transactions one at a time in each refinement phase and writing it back to the file is very I/O consuming.

The ROCK [7] is based on the number of links between two records of data items, instead of the distances between them. The links capture the number of other records that the two are both sufficiently similar to it. ROCK heuristically optimizes a cluster quality function with respect to the number of links in an agglomerative hierarchical fashion. ROCK has proved to be quite effective in categorical data clustering, but it is naturally inefficient in processing large databases [24]. The base algorithm is cubic in the dataset size, which makes it unsuitable for large problems. Therefore, ROCK may be suitable for small datasets. The ROCK data's format assumes that the similarities between data items are given. Hence, ROCK uses the similarity measure between two transactions as a number of common neighbors, but the computational cost is heavy and sampling has to be used when clustering large dataset [7]. The choice of $f(T)$ is critical in defining the fitness function, and the authors point out that the function depends on the dataset as well as on the kind of clusters that the user is interested in. Thus, the choice of the function is a weak and difficult task [15]. Besides, ROCK is difficult to fine-tune to find the right parameter T .

The COOLCAT [15] algorithm is based on the idea of entropy reduction within the generated clusters. Therefore, it does not rely on distance or arbitrary metrics. The algorithm groups points in the dataset trying to minimize the expected entropy of the clusters. This approach requires only parameter, which makes it stable and useful for larger datasets. However, the problem appears in the order in which the points are processed or grouped because of the point that appears to be a good fit for a cluster using a particular order of process may become a poor fit as more points are clustered using another order of process. To reduce this problem, the author added a re-processing step of a fraction of the points in the batch, so points are clustered in each batch and the worst fit points are re-clustered, while the number of occurrences for each of the attributes' values in a particular cluster is used to determine the goodness of the fit. However, this step increases the complexity of the algorithm specially in determining the number of fractions and worst fit points.

The CACTUS [23] is based on the concept of the common occurrences for the categories of different variables. The categories are considered strongly connected if the difference in the number of occurrences is greater than a user-defined threshold. The algorithm includes three phases: summarization, clustering and verification. In the summarization phase, the summary information is computed from the dataset. In the clustering phase, the summary information is used in discovering a set of candidate clusters. In the validation phase, the actual set of clusters is determined from the set of candidate clusters [6]. CACTUS could perform better, if the inter-attribute and intra-attribute summaries fit in the main memory. Like ROCK, this algorithm may be more suitable for small datasets. There

are main problems with CACTUS; first, it does not scale since it requires the calculation and storage of potentially large similarity matrices; second, it lacks stability when the data is re-shuffled in the similarity matrices, it includes an unnatural distinguishing set assumption; and third there is no extension step after the cluster projections is found.

Small-Large Ratio or SLR [21] uses the measure of the ratio between small to large items; the item is marked as large or small depending on the number of its occurrences in a cluster. The algorithm tries to minimize the ratio of the number of small items to that of large items in each cluster. The goal of this method focuses on designing an efficient algorithm for the refinement phase of the LargeItem algorithm [20]. The SLR algorithm compares the small to large items' ratios with the pre-specified SLR threshold α to decide the best cluster for each transaction. SLR needs to set the support Θ , the weight w , the maximal ceiling E , and the SLR threshold α . In general, the SLR algorithm must compute all the costs of clustering when transaction t is put into another cluster to use the small-large ratios, which adds additional computational steps, and the large number of parameters makes it difficult to adapt the algorithms. Thus, the algorithm did not reduce the memory and I/O consumption of the LargeItem method. The authors in [22] concluded that both the LargeItem and SLR method suffer a common drawback; that they may fail to give a good representation of the clusters.

The CLOPE [24] approach depends on the ratio between the height and the width. The height represents the number of transaction's item occurrence, while the width represents the number of clusters. The CLOPE tries to increase the height-to-width ratio of the cluster histogram. The larger height-to-width ratio of histogram the better intra-cluster similarity. CLOPE needs to set the repulsion r . There are two disadvantages of the r parameter. First, if there is no knowledge about the behavior of the dataset, it is difficult for users to expect the best value of repulsion r [4], so they must run the clustering phases more than once to get feedback on the best values for the clusters' numbers. Second, the CLOPE algorithm runs slower for non-integer repulsion r -values because of the computational overhead that comes with the floating point. The algorithm also requires two additional steps to handle adding and removing a transaction into and from a cluster in case of any refinement step.

The SEED [22] approach generates an initial seed of cluster centroid. The algorithm starts by finding the optimal number of clusters. SEED tries to maximize the fitness function value. This fitness measure calculates the average similarity between every transaction in a cluster to its centroid. The update of centroids will result in the need for clusters' re-organization. The process of centroid update and clusters' re-organization will be repeated until a suitable point of stability of the fitness function is reached.

The WCD [4] algorithm tries to preserve as many frequent items as possible within clusters and controls items' overlap between clusters. The WCD uses a partition-based clustering approach and tries to maximize the criterion function EWCD, "Expected Weighted Coverage Density". However, by default, when all transactions are considered in a single cluster, it will get the maximum EWCD, since this function cannot determine when the algorithm has to stop because merging

clusters maximizes the EWCD. Therefore, an additional phase prior to the clustering phases is required to determine the best number of clusters by taking a sample of data and running it on different values of K . This makes the algorithm's performance poor in a dynamic environment since the number of clusters can suddenly change.

The CLUC [28] algorithm depends on a similarity measure called cohesion that determines the degree with which items belong to clusters. The CLUC clusters data in two phases, initialization and refinement as most of the algorithms. But, this approach is most similar to LargeItem [20] except in the way of assigning the items in each cluster. The main drawback is found in multiple scanning of data to complete the clustering process.

The LIMBO [16] algorithm is based on the Information Bottleneck (IB) method, which uses the mutual information metric to define the measure for categorical clustering. Therefore, the algorithm works to minimize the information loss when grouping the items into clusters. The clustering approach works with three phases, Building Tree, Clustering, and associating tuples with clusters. The benefit of this approach is that it uses to cluster both tuples and attribute values; therefore, it can be classified as a hybrid algorithm.

The CBDT [29] approach is based on the distance between transactions. In this approach, the similarity between clusters is processed in three stages of calculating the distances: 1) between individual items; 2) between corresponding cells of different transactions; and 3) between transactions themselves. Therefore, the similarity measurements is considered time consuming. This is in addition to the large memory needed to store the pairwise results, which makes the algorithm's performance poor in large-scale data.

The Squeezer [1] approach works with one phase and is most similar to the allocation phase of the LargeItem algorithm with a similarity measure based on statistics. The algorithm reads transactions data in sequence and assigns it to the first maximum similar cluster or assigns it to a new cluster based on a minimum similarity. The output from this approach could change in case of change in the sequence of data input, as there is no refinement step. In addition, the similarity measure result will depend on the first clusters generation. The performance of this algorithm is better for uses with large-scale data; however, the approach use only local similarity computations to determine the maximum similar cluster.

Table III presents a comparison between the studied algorithms in the above. The clustering approach of each technique and the clustering phases are stated, as they affect the purity as well as the number of clusters.

IV. PROPOSED PWO SIMILARITY MEASURE

The goal of any clustering algorithm is to reach the final pure state of clusters, so an estimation function must be adapted to measure how many object in one cluster are different with each objects in other clusters and at the same time, the objects within a cluster are similar. It is noted that the key difference between most of the methods is in defining the criterion function for measuring similarity. However, the difficulty lies in proposing a good scenario to solve the

overlapping between clusters that depends on the underlying dataset. Therefore, to solve the overlapping problem, there are two requirements need to be considered in evaluation: the maximization of the frequent items within clusters and the minimization of the items that are overlapping between clusters.

A. The Overlapping Weight

It refers to the number of occurrences of an item in a cluster C as the item weight in C . For the purpose of the comparisons with the proposed measure function, the modified Jaccard theory [30] of similarity between clusters is described as in (1).

$$J_c(C_1, C_2) = \frac{|C_1 \cap C_2|}{|C_1 \cup C_2| - |C_1 \cap C_2| + 1} \quad (1)$$

In that sense, the similarity between two clusters increases along with the increase of the total intersection between them comparing with the total difference of the intra-join. It is noticed that Jaccard similarity neglects the weight (or support) of items in the clusters, which is significant in case of cluster's categorical dataset. Therefore, another measure that takes the weight of items is needed to take in consideration when measuring similarity, not only the number of items in the intersection.

B. The Probability of the Weights of Overlapped Items

The Probability of the Weights of Overlapped items (PWO) is introduced as a new measure function that estimates the goodness of clusters. Given a cluster C_k , suppose the number of distinct items is M_k , the items set of C_k is $I_K = I_{k1}, I_{k2}, \dots, I_{kM}$, and the sum of occurrences of all items in cluster C_k is S_k , as calculated by (2).

$$S_k = \sum_{j=1}^{M_k} |I_{kj}| \quad (2)$$

Now, the weight of an item, W_{I_j} , inside a cluster C_k is defined as the ratio of occurrences of an item I_j to the sum of occurrences of all items inside the cluster; in other words the probability of an item inside the cluster C_k , is as shown by (3).

$$W_{I_j} = P(I_j) = \frac{|I_j|}{S_k} = \frac{|I_j|}{\sum_{j=1}^{M_k} |I_{kj}|} \quad (3)$$

In this sense, the total probability of all items within a cluster is equal to one, (4).

$$\sum_{j=1}^{M_k} W_{I_j} = \sum_{j=1}^{M_k} P(I_j) = 1 \quad (4)$$

Let overlap O_{ij} be the list of mutual items between cluster C_i and cluster C_j , where $O_{ij} = C_i \cap C_j$, and $|O_{ij}|$ represents the number of mutual items. All items belonging to O_{ij} have two possible weights: its weight for each cluster separately, and, its weight depending on the group of transactions.

TABLE III. SUMMARY OF CATEGORICAL CLUSTERING ALGORITHMS

Algorithm	Clustering Approach	No. of Phases	Measurement Approach	Metric Parameters	No. of Parameters	No. of Classes
Squeezer	Assign data to the first max similar cluster.	1	Local Similarity	Minimum similarity s	1	N/Y
SLR	Minimize the small large Ratio between clusters	2	Local and Global Similarity	Minimum support Θ , weight w , maximal ceiling E , and the SLR threshold α	4	N
LargeItem	Increase items' frequency inside clusters			Minimum support Θ , and the weight w	2	N
CACTUS	Depend on shared items between clusters			Distinguishing number K , passes D		N
COOLCAT	Minimize the expected entropy for each clusters			Minimum Entropy	1	Y
CLOPE	Increasing the high-to-width ratio of the cluster histogram			Repulsion r		N
ROCK	Increase the number of links between items inside the cluster			Fitness function $f(T)$		N/Y
CLUC	Depend on the cohesion measuring similarity to assigned items to clusters			User-defined threshold α		N
CBDT	Pairwise less distance between transactions			Number of classes r		Y
WCD	Increase the coverage of large items inside the cluster	Number of classes K	Y			
LIMBO	Minimize the information loss.	Information loss threshold α	Y			
SEED	Generated of an initial seeding of cluster centroids		Global Similarity	Minimum support Θ	N	

Now, let $WO(C_i | C_j)$ represents the sum of weights of all items of a cluster C_i that overlap or intersect with cluster C_j , and expressed in (5).

$$WO(C_i | C_j) = \sum_{k=1}^{|O_{ij}|} W_{I_k \in (C_i \cap C_j)} \quad (5)$$

Similarly, $WO(C_i | C_j)$ is the sum of weights of all items of cluster C_j that overlap with cluster C_i . Now the definition of the probability of the weights of overlapped items between clusters C_i and C_j , $PWO(C_i, C_j)$, is presented as (6).

$$PWO(C_i, C_j) = WO(C_i | C_j) \cdot WO(C_j | C_i) \quad (6)$$

Hence, the similarity measure is defined by (7).

$$sim(C_i, C_j) = PWO(C_i, C_j) \quad (7)$$

Now, the similarity function $sim(C_i, C_j)$ captures the closeness between the pair of clusters C_i and C_j . Actually, the sim values range between zero and one, with larger values indicating that the clusters are more similar. The sim value is one for identically matching clusters and zero for very dissimilar clusters.

C. Proof of PWO Similarity Metric

In [31] defines the similarity measure S as a function with a non-negative real value that satisfies three properties (1-3). Moreover, if S satisfies properties (4) and (5) then S is called a metric similarity measure.

- 1) $\exists s_0 \in R : -\infty < S(x, y) \leq s_0 < +\infty, \forall x, y \in X$.
- 2) $s(x, x) = s_0 \forall x \in X$.
- 3) $s(x, y) = s(y, x) \forall x, y \in X$.
- 4) $s(x, y) = s_0 \leftrightarrow x = y \forall x, y \in X$.
- 5) $s(y, z) \leq [s(x, y) + s(y, z)]s(x, z) \forall x, y, z \in X$.

The following is a proof of PWO being a metric similarity measure. Proof. $S = S_{PWO}$:

- 1) Property 1 is satisfied by the properties of probabilities $0 \leq S_{PWO} \leq 1$, thus $s_0 = 1$.
- 2) Property 2 is trivially satisfied by the fact that S_{PWO} is a similarity measure, thus $S_{PWO}(x, x) = 1$.
- 3) Property 3 is also satisfied by the fact that S_{PWO} is a similarity measure, so $S_{PWO}(x, y) = S_{PWO}(y, x)$.
- 4) If $S_{PWO}(x, y) = 1$ then x and y are identical, which together with (2) satisfies Property 4.
- 5) Property 5 is also satisfied by the properties of probabilities.

V. PWO CLUSTERING

A. Clustering using PWO

In order to study the ability of PWO as a similarity measure in clustering, the cost function in the LargeItem algorithm [20] and the PWO similarity metric are compared. While the goal of LargeItem is to minimize the total cost of each cluster, the goal is to maximize the similarity between transactions in the same cluster. Therefore, the LargeItem algorithm it modified to adapt this goal. Another important point is that the cost function of LargeItem is relative to the $MinimumSupport(\theta)$ that is given by the user. Therefore, the similarity between a pair of clusters is only accepted if it is larger than or equal to the minimum support, as shown in (8). Thus, higher values of θ correspond to higher thresholds for the similarity between a pair of clusters before they are considered similar. Algorithm.1 illustrates an overview of the clustering algorithm.

$$sim(C_i, C_j) \geq Minimum - support(\theta) \quad (8)$$

B. Merging Clusters by Groups

To speed up the clustering process, a new strategy of merging clusters is applied. Instead of merging the most

Algorithm 1 PWO Clustering Algorithm

```
while not EndOfFile do
  Read the next transaction  $\langle t, - \rangle$ 
  Allocate  $t$  to an existing  $C_i$  with MAX similarity larger
  than MIN support or in a new cluster
  Write  $\langle t, C_i \rangle$ 
end while
{Refinement phase}
 $not\_moved \leftarrow \mathbf{true}$  { $not\_moved$  true if no transaction  $t$  is
moved between clusters.}
repeat
  while not EndOfFile do
    Read the next transaction  $\langle t, C_i \rangle$ 
    Move  $t$  to an existing non-singleton cluster  $C_j$  that has
    a MAX similarity with it
    if  $C_i \neq C_j$  then
      Write  $\langle t, C_j \rangle$ 
       $not\_moved \leftarrow \mathbf{false}$ 
      eliminate any empty cluster
    end if
  end while
until not_moved
```

Algorithm 2 PWO-M Clustering Algorithm

```
while not EndOfFile do
  Read the next transaction  $\langle t, - \rangle$ 
  Allocate  $t$  to an existing  $C_i$  with MAX similarity larger
  than MIN support or in a new cluster
end while
{Refinement phase}
 $no\_merge \leftarrow \mathbf{false}$ 
{ $not\_merge$ : false if no similar cluster found to merge.}
repeat
   $num\_cluster \leftarrow i$ ;
  group clusters with similarity  $\geq \alpha$ ;
  if  $num\_cluster == i$  then
     $no\_merge \leftarrow \mathbf{true}$ 
  end if
until  $not\_merge == \mathbf{true}$ 
```

similar pairs of clusters, the most similar groups of clusters are merged. First, each transaction is allocated in a single cluster and compute the similarity between clusters, and then each cluster is assigned to a group if it is similar to any cluster in the group. The similarity's rule is applied using (8). The process of merging group of clusters at once instead of merging pairs of clusters minimizes the number of merging steps without affecting the purity.

Algorithm 2 presents a fast way of clustering transactions with the use of PWO as the similarity measure as well as in group's merging. To differentiate it from the previous algorithm explained in Algorithm 1, this algorithm is called PWO-M, in which groups of similar clusters continue to be merged until no more merge is possible. This approach speeds up the clustering process while maintaining the same degree of purity and number of clusters.

Algorithm 3 PWO Clustering Algorithm with Predefined Number of clustering

```
while not EndOfFile do
  Read the next transaction  $\langle t, - \rangle$ 
  Allocate  $t$  to an existing  $C_i$  with MAX similarity larger
  than MIN support or in a new cluster
end while
{Refinement phase}
repeat
  group pair clusters with MAX  $sim(C_i, C_j)$ 
until  $num\_cluster == predefined$ 
```

C. Clustering using a pre-defined number of clusters

In case there is a pre-defined number of clusters before clustering, it can modify the algorithm presented in Algorithm 2 can modified to be as the one presented in Algorithm 3 by adding one additional step in the refinement phase. Merging pairs is continued with maximum similarity until the number of clusters reaches the required number. Results of this modification presented in the Experiments.

VI. OVERLAP ESTIMATOR

The similarity of clusters changes according to the dataset, so the overlapping between clusters varies depending on the behavior of the dataset. An automated framework is proposed to specify the best threshold value for determining the cluster's neighbor, i.e. similar clusters, according to the training dataset.

A. Cluster's Neighbours

Initially, A cluster's neighbors are those clusters that are considerably similar to it, and therefore they can be merged with it forming a large cluster.

Now, Let $sim(C_i, C_j)$ be a similarity function that normalizes and captures the degree of similarity between the pair of clusters C_i and C_j . The sim values are between zero and one, with larger values indicating that the clusters are more similar.

Given a threshold α between 0 and 1, a pair of clusters C_i, C_j are defined to be neighbors if (9) holds

$$sim(C_i, C_j) \geq \alpha \quad (9)$$

In (9), α is a parameter that can be used to control how close a pair of clusters must be in order to be considered neighbors; it is called the "neighborhood threshold".

Accordingly, higher values of α correspond to higher thresholds for the similarity between the pair of clusters before they are considered neighbors. Assuming that sim is one for matching clusters and zero for very dissimilar clusters, a value of one for α constrains clusters to be neighbors to only identical clusters. On the other hand, a value of zero for α permits any arbitrary pair of clusters to be neighbors. Equation (8) is the same as (9) if $\alpha = \theta$; θ is a user parameter while α is the clusters' neighbor threshold.

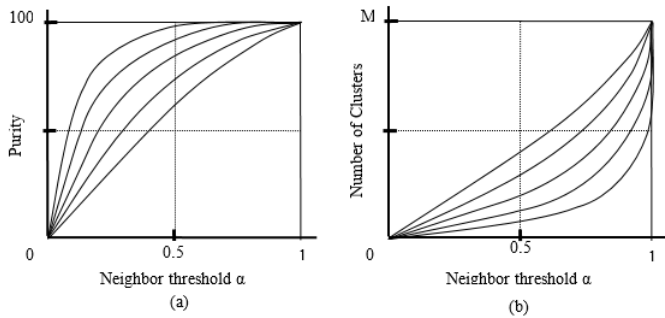


Fig. 1. The Neighbour's threshold relationship with purity and number of clusters.

B. Characteristics of the Neighbor's Threshold

The neighbor's threshold α has the following characteristics:

- Increasing (decreasing) the neighbour's threshold α increases (decreases) purity and increases (decreases) the number of clusters.
- Its result depends on the data type. A value of α that produces a purity of 100% on data D_1 does not necessarily produce a purity of 100% on data D_2 .

Fig. 1 summarizes the relationship between the neighbor's threshold from one side, and purity and number of clusters from another side. In addition, it represents different dataset type curves. The growth rate with purity is appearing logarithmically in Fig. 1(a), while it is exponential with the number of clusters as shown in Fig. 1(b).

C. Estimating the overlap parameter

The best value of the neighbor's threshold α is the value that would get high purity with minimum number of clusters. The best value is called the overlap threshold α . The overlap estimator is used to estimate the best value of the neighbor's threshold α . The framework tries to find the minimum value of closeness of cluster's neighbor that will produce 100% purity of clusters, and uses this value in clustering the transactions of dataset.

To estimate the best overlap threshold α value, the following approach is applied:

- 1) Starting by clustering the training dataset at a minimum support of 100%.
- 2) Merging the clusters using different values of the neighbour's threshold α .
- 3) Test the purity of the output clusters based on the training set classes values.
- 4) Repeating steps (2-3) until reaching the minimum value of the neighbour threshold α getting 100% purity value.

In addition, it is noticed that starting from $\alpha = 0$ ascending to one is faster than starting from $\alpha = 1$ descending to zero in computing the similarities. Algorithm 4 shows the overlap estimator algorithm.

Algorithm 4 The Overlap Estimator Algorithm

Ensure: Input: The training dataset D_n , Classes C_n
 $\alpha \leftarrow 0$ { α is the threshold parameter}
while Purity $\neq 1$ **do**
 $\alpha \leftarrow \alpha + 0.1$
 PWO-M (D_n, α , Output)
 Check(C_n , Output, Purity)
end while
return $\alpha - 0.1$

VII. F-TREE CLUSTERING

Unlike traditional data, categorical clustering requires transactions to be partitioned across clusters in such a manner that instances within a cluster share a common set of large items, where the concept of the large follows the same meaning attributed to frequent items in association rule mining [32]. Thus, it is clear that categorical clustering requires a fundamentally different approach from the traditional clustering technique. F-Tree [13] is a summarization clustering algorithm that clusters categorical data based on a new tree structure.

A. F-Tree Clustering Algorithm

The basic F-Tree approach consists of the four main steps as follows:

- 1) Calculating items' frequencies: it scans the input dataset to rank all items.
- 2) Building a F-Tree: it inserts all transactional items of the dataset into F-Tree structure; it uses the frequencies of items to reorder the transactional items before inserting it into the F-Tree as discussed in the previous example.
- 3) Extracting initial clusters: initial clusters are generated using F-Tree by pruning the F-Tree at some level based on the minimum support.
- 4) Refining clusters: it applies the merging algorithm operated with PWO measurement to merge similar clusters that are extracted from the previous step.

All the above steps are divided in two phases. The allocation phase is concerned with the first three steps; while the refinement phase is concerned with only the last step.

1) *F-Tree Data Structure:* The F-Tree data structure is designed to compress the categorical dataset. As the categorical dataset contains a set of records, the F-Tree groups the records using their shared items or attributes in the tree. In the beginning, the global item frequencies are computed. Then all transactions' items are inserted in the F-tree using the item as the node key. The insertion of item is based on the global frequency order. As a result, the groups of items that get in the path starting from the tree root to any leaf node composes a single record, and so all paths from root to leaves nodes compose all transactions in the dataset. Thus, each sub-path can represent a set of transactions, which share the same path prefix. The following example illustrates the F-Tree process.

Suppose the following is the transaction records {ACF, ACH, BDE, AE, BF}. Computing the global items frequency gives {A(3), C(2), B(2), F(2), E(2), D(1), H(1)}. Now sorting the records' items based on global frequency gives {ACF,

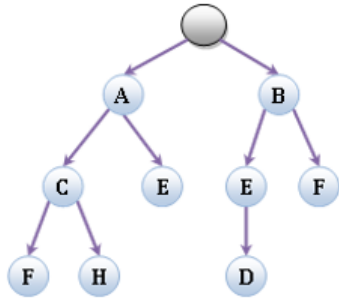


Fig. 2. An example diagram of an F-Tree structure.

ACH, BED, AE, BF}. Fig. 2 shows the F-Tree structure of these records; then, it can absolutely be seen that node A is a shared item between three records {ACF, ACH, AE}. In addition, Node B is a shared item between two records {BED, BF}.

2) *Generate Clusters from F-Tree*: Extract the clusters from the F-Tree depending on the F-Tree levels since each inner node shares all items in the upper level. Hence, the minimum support parameter is replaced by F-Tree depth, and for generalization, (10) is used.

$$ClusterLevel = MinimumSupport * tree_{depth} \quad (10)$$

To illustrate how the clusters are extracted from the previous graph in Fig. 2, for instance if the $minimum_support\theta = 75\%$, and $tree_depth = 3$ then the cluster level = 2. At the second level (depth) there are four nodes. Each of those nodes contains one or more transaction in it or in its children. Therefore, there are four clusters as {{ACF, ACH}, {AE}, {BED}, {BF}}.

While if extracting clusters at the first level (depth), there are only two clusters as well as there are two nodes. These clusters are {{ACF, ACH, AE}, {BED, BF}}. At this point, it easy to notice that the number of clusters decreases as the level goes up to the root of the tree.

3) *Merging Clustering Algorithm*: The major steps of the merging algorithm are defined as the following:

- Computing the similarity list between clusters.
- Creating the group of neighbour clusters.
- Merging all clusters in the same group.
- Repeating these steps until there is no further merging.

The similarity list between clusters are computed using the merging overlap threshold α . Then, any similar cluster will belong to the same group in which the group of clusters contains only neighbour pairs of clusters. Lastly, it merges all clusters' neighbours inside the same group. A new generation of clusters could be more likely similar or dissimilar based on the result of merging. Therefore, the refinement procedure will repeat the merging algorithm until there is no more merge done or no new cluster's neighbour found.

Algorithm 5 The SP-TREE Algorithm

Ensure: The dataset D_n , Training set S_m
 F-Tree ($D_n + S_m, \alpha, C_i$)
for all s Cluster in Seed S_m **do**
 Merge all clusters c contains any items $\in s$
end for
return Clusters C_k

B. Semi-Supervisor F-Tree Clustering

The F-Tree generates a large set of pure clusters, and although the refinement step breaks down the number of generated clusters, the refinement step needs a learning process to minimize the number of pure clusters without losing its precision specially when there is a predefined number of clusters. So, a training data base could be used as a seed in the merging phase. These seeds will be used to guide the merge algorithm with correct similar sets to speed up and correct the fitting of the merging algorithm. Algorithm 5 shows the ST-Tree algorithm.

VIII. EXPERIMENTS

In this section, the accuracy is analyzed, precision and execution time of the proposed measure metric and clustering algorithms with real-life datasets. Several experiments are conducted for clustering to evaluate the general purity of the clustering algorithm using PWO. The version of LargeItem [20] algorithm is implemented for comparing performance and precision of clustering purpose, as well as LargeItem algorithm is a tree based techniques. Others algorithms results are collected from authors references.

A. Empirical Datasets

labeled datasets in Table (IV) are obtained from the UCI Machine Learning Repository [14] are used. The number of clustering presented in this table are used in clustering evaluation.

B. Pre-processing The Dataset

The input dataset is converted into a format that can be processed by algorithms handling transactional datasets. First, the class label is removed from all datasets. Second, each record of dataset is converted to a list of distinct items by mapping each property character for any attribute to a distinct numerical number for all, since each record of transactional data always contains a list of distinct items, and the record of a dataset has a list of properties' characters that may be duplicate throughout different attributes.

For instance, if there are the following dataset record {{F, G, F, F, F}, {F, G, T, N, F}, {Y, N, T, N, F}} with five attributes, then after mapping the transaction record would be {{1, 3, 5, 7, 9}, {1, 3, 6, 8, 9}, {2, 4, 6, 8, 9}}. The same mapping process is applied if the dataset contained Boolean values. For instance, if the dataset record is {{1, 0, 0, 1, 0}, {1, 0, 0, 1, 1}, {0, 1, 1, 0, 0}}, then the equivalent transaction record would be {{1, 3, 5, 7, 9}, {1, 3, 5, 7, 10}, {2, 4, 6, 8, 9}}.

TABLE IV. DATASET PROPERTIES

Dataset	Size	No. of classes		No. of features	Total Attribute	Missing Values
		Num.	Cat.			
Mushroom	8124	2	1	22	23	2480
Chess	3196	2	0	36	36	0
Car	1728	4	0	6	6	0
Pima diabetes	768	3	8	0	8	0
Breast cancer	699	2	0	9	9	0
Vote	435	2	0	16	16	288
Wine	178	3	13	0	13	0
Iris	150	3	4	0	4	0
Zoo	101	7	1	15	16	0

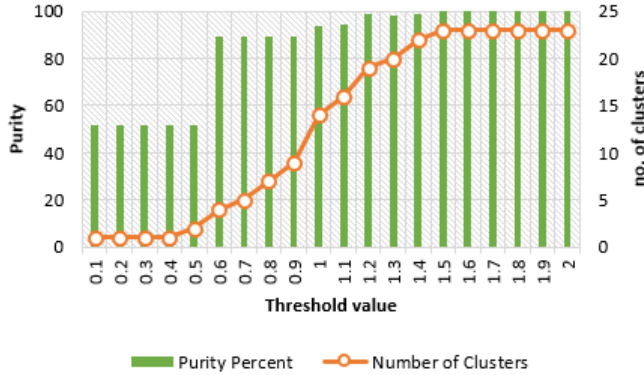


Fig. 3. The effect of threshold value on the data purity after merging clusters with Jaccard metric on the mushroom dataset.

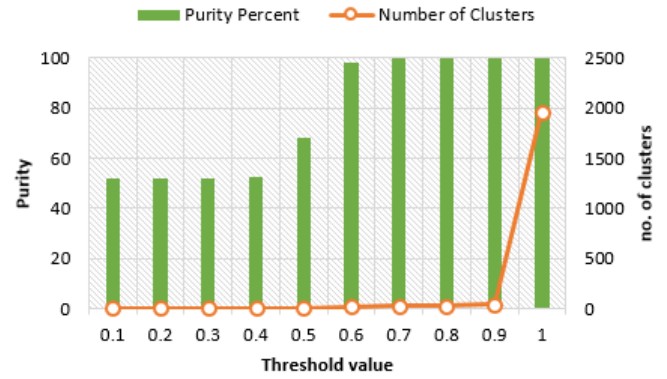


Fig. 4. The effect of threshold value on the data purity after merging clusters with PWO metric on the mushroom dataset.

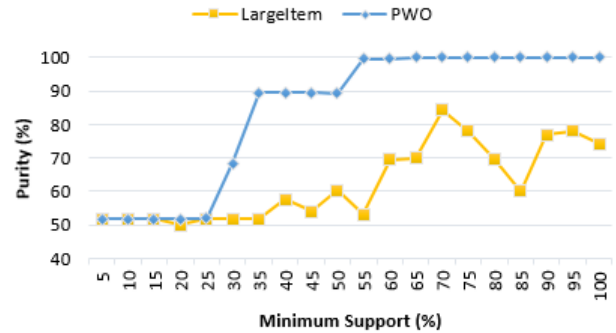


Fig. 5. Clusters Purity Results LargeItem vs. PWO on the mushroom dataset.

IX. EVALUATION STUDY

A. PWO Evaluations

1) *Experiment 1: Evaluation of PWO Metric Function:* To illustrate the advantage of PWO metric function, an analytical test is performed, in which the PWO and Jaccard capability are compared in grouping similar clusters. First, the base clusters are generated, then the similar clusters are computed using the Jaccard coefficient and then the PWO metric to merge the similar clusters in groups. Similar clusters are determined if the metric result is above the threshold value. Fig. 3 and 4 shows the result of purity with number of clusters using the Jaccard coefficient and PWO metric respectively applied on the mushroom dataset. In Fig. 4, the “100%” purity is reached when threshold value of Jaccard is equal to 1.5 and the number of clusters is “23”. While in Fig. 4, the “100%” purity is reached at threshold value of PWO equal to 0.7 and the number of clusters is 23. Comparing between the two figures indicate the normalization strength of PWO to measure the similarity.

2) *Experiment 2: Evaluation of PWO Based Clustering Algorithms:* In this experiment, the LargeItem algorithm is modified by changing its cost function with PWO.

First, the purity of resulting clusters are compared from both algorithms on the mushroom dataset. Fig. 5 shows the result of this test, and it is observed that PWO reaches purity of 90% with minimum support of 35%. As a conclusion it is found that PWO is a powerful clustering similarity measure.

Second, the resulting number of clusters is compared. It is a fact that the purity of clusters is proportional to the number

of clusters generated. Fig. 6 shows the numbers of clusters generated by LargeItem and PWO. it is observed that when the minimum support equals to 100% the PWO algorithm handles each transaction in an individual cluster because PWO equals to one.

To minimize the number of clusters, the merge strategy is applied with with PWO-M algorithm, and then compare the number of clusters with PWO algorithm. Fig. 7 illustrates the effect of group merging on minimizing the number of clusters depending on the minimum overlap value. Next, LargeItem, PWO, and PWO-M are put in comparison of their final number of clusters, as illustrated in Fig. 8. There are two tests for LargeItem when (Intra=1) and when (Intra=10). It is obvious that LargeItem produces two very different results. However, PWO-M returns a reasonable number of clusters while illustrating high accuracy and best stable result.

3) *Experiment 3: Comparison of Different Clustering Algorithms vs. PWO:* In this experiment, the following algorithms: PWO, PWO-M, LargeItem, and CLOPE are compared using the datasets and parameters in Table V.

Fig. 9 illustrates their comparison in terms of purity of resulting clusters. It was seen that PWO and PWO-M are more accurate and more stable. Second, the output number of clusters between the four algorithms are compared, illustrated in Fig. 10. It is noticed that LargeItem returns the minimum number of clusters for all datasets, while PWO and PWO-M return a very large number of clusters for Voting and

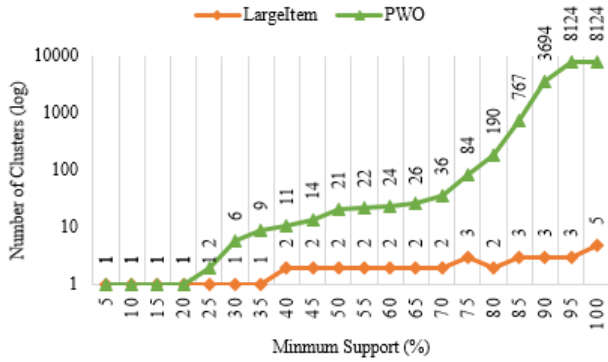


Fig. 6. Final No. of clusters using LargetItem vs. PWO on the mushroom dataset.

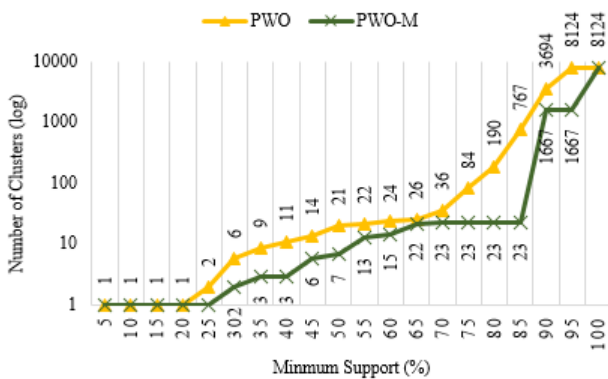


Fig. 7. Final No. of clusters using PWO vs. PWO-M on the mushroom dataset.

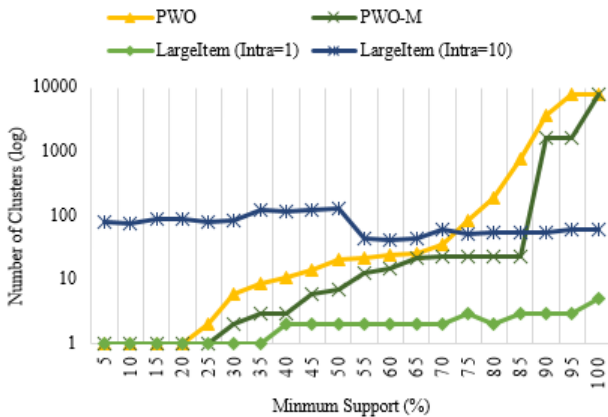


Fig. 8. Number of clusters vs. minimum support on the mushroom dataset.

Hepatitis datasets, but for the remaining datasets, they return a reasonably larger number of clusters if it is taken into consideration the gain of clusters' purity. From this test, it is concluded that PWO is measurable and its strength appears in the purity of resulting clusters.

4) Experiment 4: Evaluation of PWO Algorithm using Fixed Number Clustering : As previously explained, PWO-M is adapted to work with the case when there is a pre-defined

TABLE V. DATASETS AND SELECTED PARAMETERS

Dataset	Size	Overlap/Minimum Support
Mushroom	8124	0.7
Voting	435	0.8
Hepatitis	155	0.5
Zoo	101	0.75
Lenses	24	0.6
Adult-Stretch	20	0.6

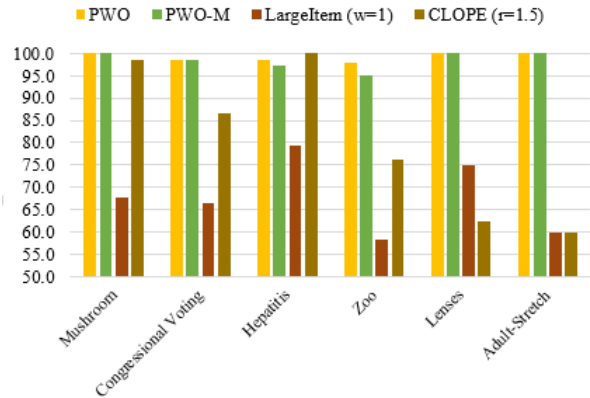


Fig. 9. Data purity in all algorithms using different datasets.

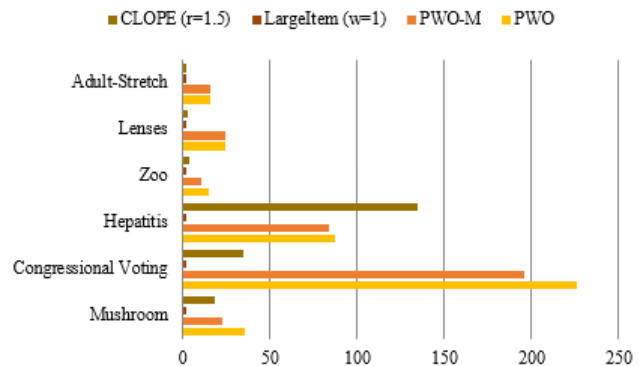


Fig. 10. Number of clusters resulting from algorithms using different datasets.

TABLE VI. SUMMARY OF FIXED-N CLUSTERING RESULTS

Dataset	No. of Classes	Purity	Avg. Classes Coverage
Mushroom	2	83.26	88.20
Zoo	7	88.12	74.28
Hepatitis	2	83.87	63.24
Voting	2	89.20	90.65

number of clusters. The result of fixed-N clustering of the Zoo dataset is 88.12% purity. The Mushroom dataset output classes having an average of 83.26% purity, while the result of the Congressional vote is 89.19% of clusters' purity. Table VI is a summary of the fixed-N clustering results. One can notice that the total average of clusters' purity is above 85%, while the average of classes' coverage is above 75%.

5) Experiment 5: Evaluation of the Overlap Estimator's Precision: To evaluate the precision of the proposed overlap estimator, at first a small training sample of the mushroom dataset is used as an input to the algorithm, 800 out of 8124

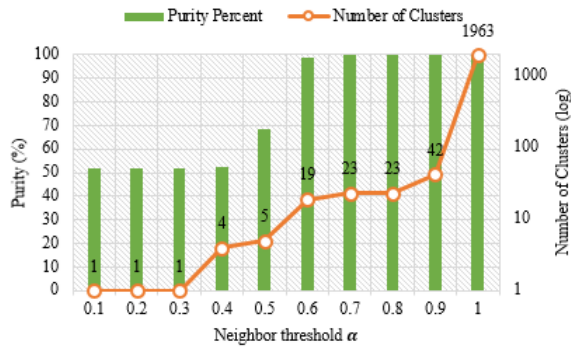


Fig. 11. Number of Clusters vs. purity in case of different neighbour's threshold values.

TABLE VII. OVERLAP THRESHOLD VALUES FOR DIFFERENT DATASET

Dataset	Total Size	Sample Size	Overlap threshold
Mushroom	8124	100	0.70
Voting	435	20	0.80
Hepatitis	155	16	0.50
Zoo	101	10	0.75
Lenses	24	12	0.60
Adult-Stretch	20	10	0.60

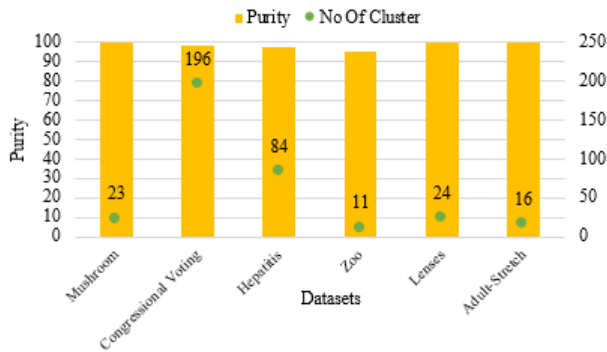


Fig. 12. Purity of clustering different datasets using the estimated overlap threshold values.

transactions selected randomly from the mushroom dataset. The overlap estimator selects a value of $\alpha = 0.7$ as an overlap threshold. Second, all the mushroom dataset is clustered and then the clusters are merged using different values of the neighbor's threshold, and the resulting number of clusters and the purity of clusters in each case are compared. The result is shown in Fig. 11, which illustrates that the value of $\alpha = 0.7$ returns the minimum number of clusters with high purity.

The same experiment is repeated using different datasets. The overlap threshold values for the datasets is listed in Table VII, while the purity and number of clusters for each dataset is displayed in Fig. 12. It is noticed that the clusters' purity for all datasets is above 95%; this indicates that the overlap estimator is very effective in detecting the behavior of data.

6) *Experiment 6: Evaluation of Overlap Estimator Scalability:* In this experiment, the scalability degree or the effect of the sample size is evaluated to estimate accurate neighbor's threshold. In Fig. 13, the clusters' purity is measured versus

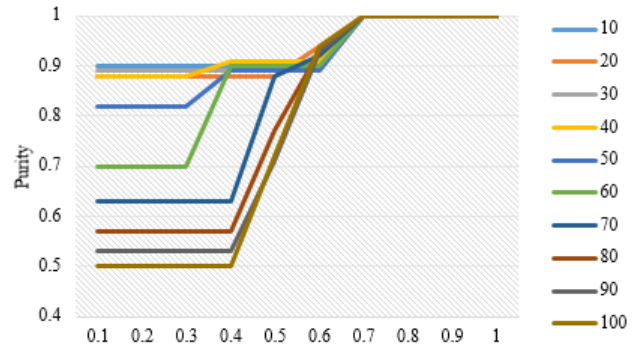


Fig. 13. Purity vs. Neighbour's threshold value using different sizes of the mushroom dataset.



Fig. 14. Overlap Estimation Based on Sample Sizes of different datasets.

neighbor's threshold values on different percentages of the mushroom dataset size. It is noticed that all dataset sizes reach 100% purity when the neighbor's threshold value of 0.7 is used. So, the perfect estimate of overlap threshold for the mushroom dataset should be 0.7 regardless of the dataset size.

Now, the overlap estimator is applied using different sizes of dataset sample to compare the accuracy of the overlap threshold. The overlap estimator has perfectly estimated the value of the overlap threshold starting from 10% sample of the dataset, although with a 5% sample dataset (around 406 transactions) the overlap threshold is $\alpha = 0.6$ and the purity of the clusters is 93.6%, which is also acceptable.

Fig. 14 illustrates the effect of three sample sizes (10%, 50%, and 100%) on estimating the overlap threshold for the different datasets. The Adult-Stretch dataset is very small and therefore the different sample sizes produced very different estimated thresholds, i.e. it is difficult to estimate the clusters' behavior. From this experiment, it is observed that the overlap estimator could expect the best overlap threshold using a small sample if the sample is normally distributed on all classes.

Fig. 15 illustrates the clusters' purity when 10%, 50% and 100% of dataset size is used as a sample for estimation for different datasets using the PWO-M algorithm. It is seen that a nearly 100% purity is achieved from using only 10% of the dataset as a sample in the case of the Mushroom, Hepatitis, Zoo and Lenses datasets. However, this is not the case for other datasets due to the following reasons: 1) there are missing

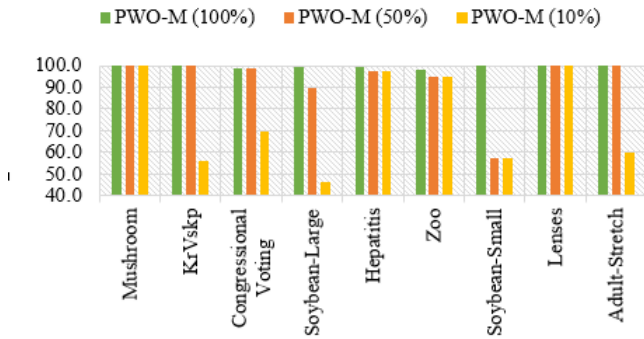


Fig. 15. Clusters' purity vs. Sample Sizes for different datasets using PWO-M.

TABLE VIII. MUSHROOM CLUSTERING PURITY VS NO. CLUSTERS

Clustering Algorithm	Purity	No. of Clusters
F-Tree ($\theta = 0.8$)	100	23
Hybrid	100	23
CLUC	100	24
CLOPE	100	30
CLICKS	100	553
Fast Clustering	99.90	23
Squeezer	99.90	24
ROCK	99.60	21
SCCADDS	99.00	19
CLICKS	97.00	19
LargeItem	95.62	14
F-Tree ($\theta = 0.3$)	95.42	16
CLICKS	87.10	14

attributes' values as in the Voting dataset. 2) The dataset could be very small as in the Adult-Stretch dataset that has 20 transactions and Soybean-Small that has 47 transactions. 3) There is a large number of clusters such as in Soybean-Large that has 19 classes. 4) There is a large number of attributes per transaction, such as in the KrVskp dataset that has 36 attributes and Soybean-Large that has 35 attributes.

Finally, three results are concluded. First, the overlap threshold value can be change for different data types. Second, the estimation of the overlap threshold using a small size of data would depend on the data type and number of clusters in the dataset. Third, the overlap estimator could improve the clusters' purity.

B. F-Tree Clustering Evaluations

1) *Experiment 1: Comparing clustering algorithms without pre-defined number of clusters:* In this experiment, almost all of the algorithms are compared for best purity with closer number of clusters. The algorithms are run without pre-defined number of clusters. Table VIII lists different clustering algorithms in order to evaluate the clustering purity versus the number of clustering. This table is sorted by purity in descending order and by number of clusters in ascending order. The algorithms that are given the closer number of clusters with high purity are both F-Tree [13] and Hybrid [19], which gives 100% of purity with only 23 clusters. It is also figured that LargeItem returned the minimum number of clusters but with higher purity than CLICKS [25].

TABLE IX. DIFFERENT DATASET CLUSTERING PURITIES

Algorithm	Mushroom	Car	Zoo	Hepatitis	Voting	Cancer
ROCK	77.00	77.08	-	99.35	79.00	97.20
COOLCAT	76.00	-	-	-	87.00	-
Squeezer	53.60	-	89.00	-	61.80	86.20
LIMBO	89.00	44.50	-	84.52	87.12	69.93
DELTA	89.02	30.15	-	69.67	89.43	70.97
SCCADDS	89.00	-	-	-	88.00	-
Ensemble	89.00	-	93.00	-	87.00	96.70
SF-Tree (10% seed)	99.00	89.90	93.00	79.00	87.00	96.20
CBDA	89.02	-	-	-	91.95	-
DILCA M	89.02	70.08	-	83.22	91.95	74.47
DILCA RR	89.02	70.08	-	69.67	89.43	74.47

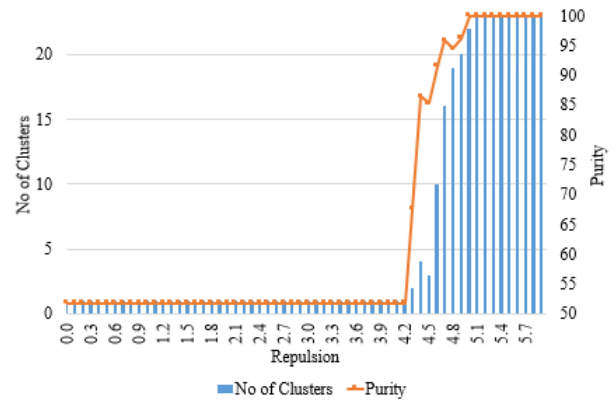


Fig. 16. Clusters' purity vs. No of clusters for different values of Repulsion of CLOPE algorithm.

2) *Experiment 2: Comparing clustering algorithms with pre-defined number of clusters:* In this experiment, almost all of the algorithms are compared for best purity but with pre-defined number of clusters given prior to the clustering process. In Table IX, the results are presented across different datasets. it could notice that SF-Tree is the best algorithm across different full categorical datasets such as Mushroom, Car, and Zoo. However Voting dataset is categorical but contains a lot of missing data which affects the purity of algorithm. But, in numeric dataset such as Hepatitis and Cancer, It fail to get a competitive result as this version of F-Tree is based on exact matching of feature's value or records.

3) *Experiment 3: Evaluation of CLOPE cost function:* It is implemented the CLOPE clustering algorithm to analyze the effect of CLOPE [24] cost function on merging pure clusters. In this experiment, there is a try to minimize the number of clusters generated from F-Tree using CLOPE cost function. On mushroom dataset, F-Tree generates 23 pure clusters. These clusters are input to CLOPE algorithm as existing clusters, then run algorithm to minimize the number of clusters. there is a try of a range of r-parameter (from 0 to 6) of CLOPE algorithm and measure the number of clusters and purity each time. In Fig. 16, it is noticed that at value of (4.9) the number of clusters is decreased, but with an effect on purity. The CLOPE cost function is failed to minimize the number of clusters in addition to the main problem of determining the best value of r-parameter or "Repulsion".

TABLE X. ANALYSIS OF SF-TREE ON DIFFERENT DATASETS

Dataset	Mushroom	Cancer	Voting	Zoo
Number of samples	100	50	25	20
DB Size	8124	699	435	101
Seed Percent(%)	1%	7%	6%	20%
Class	2	2	2	7
Entropy	1.003	0.942	1.012	2.201
CAIR	0.972	0.868	0.762	0.946
Info-Loss	0.014	0.066	0.118	0.032
E-min	0.010	0.044	0.090	0.059
Precision	0.990	0.956	0.910	0.941
Recall	0.990	0.956	0.916	0.976
F-Measure	0.990	0.956	0.911	0.940
Purity	99.00	95.60	91.00	94.10

4) *Experiment 4: Analysis of Clustering with SF-Tree:* The goal of this experiment is to test the dependence between the dataset and seed percentage, it is noticed that the seed can be used randomly but it is effective if the dataset size is smaller. Table X shows that in Zoo dataset it is used 20% of dataset as a seed in order to gain the purity because the Zoo is very small dataset around 101 records, while in mushroom the use of 1% of dataset is sufficient to gain a good purity value because the dataset is large enough.

X. LIMITATIONS

This clustering approach using PWO was unable to sufficiently minimize the number of resulting clusters and further research is needed to overcome this drawback. The overlap estimator framework is designed for categorical datasets; extending it to domains with continuous values will be a challenging task.

XI. CONCLUSION AND FUTURE WORK

In this paper, a new similarity measure, “PWO” is proposed to overcome the overlapping between clusters. From the experiments, it is concluded that PWO is applicable to different categorical datasets and generates acceptable degree of clusters’ purity. New clustering algorithms using PWO is presented and experiments showed that this approach is effective. PWO also can be applied in many applications, so it can be used in 1) measuring the similarity between clusters of categorical or transactional dataset, 2) measuring the best pair of clusters, and 3) classifying the dataset based on the similarity between categorical items. Since it is important to determine the best similarity threshold for different datasets, the overlap estimator framework based on the training dataset is proposed. Experiments show that only 10% of the total datasets is sufficient to detect the best similarity threshold value.

REFERENCES

[1] Z. He, X. Xu, and S. Deng, “Squeezer: an efficient algorithm for clustering categorical data,” *Journal of Computer Science and Technology*, vol. 17, no. 5, pp. 611–624, 2002.
[2] Z. Huang, “A fast clustering algorithm to cluster very large categorical data sets in data mining,” in *DMKD*. Citeseer, 1997.
[3] D. Gibson, J. Kleinberg, and P. Raghavan, “Clustering categorical data: An approach based on dynamical systems,” *Databases*, vol. 1, 1998.

[4] H. Yan, K. Chen, and L. Liu, “Efficiently clustering transactional data with weighted coverage density,” in *Proceedings of the 15th ACM international conference on Information and knowledge management*. ACM, 2006, pp. 367–376.
[5] N. Ye, *The handbook of data mining*. Lawrence Erlbaum Associates Mahwah, NJ, 2003, vol. 24.
[6] M. W. Berry and M. Browne, *Lecture notes in data mining*. World Scientific, 2006.
[7] S. Guha, R. Rastogi, and K. Shim, “Rock: A robust clustering algorithm for categorical attributes,” in *Data Engineering, 1999. Proceedings., 15th International Conference on*. IEEE, 1999, pp. 512–521.
[8] H. D. Margaret, “Data mining introductory and advanced topics,” *LPE: Pearson Education Publishing*, 2003.
[9] S. Guha, R. Rastogi, and K. Shim, “Cure: an efficient clustering algorithm for large databases,” in *ACM SIGMOD Record*, vol. 27, no. 2. ACM, 1998, pp. 73–84.
[10] D. J. Hand, H. Mannila, and P. Smyth, *Principles of data mining*. MIT press, 2001.
[11] S. Sharma, *Applied multivariate techniques*. John Wiley & Sons, Inc., 1995.
[12] K. Chen and L. Liu, ““best k”: critical clustering structures in categorical datasets,” *Knowledge and information systems*, vol. 20, no. 1, pp. 1–33, 2009.
[13] M. A. Mahdi, S. E. Abdel-Rahman, R. Bahgat, and I. A. Ismail, “F-tree: An algorithm for clustering transactional data using frequency tree,” in *Al-Azhar University Engineering Journal, Eleventh International Conference*, vol. 5, no. 8. Al-Azhar University, Cairo, Egypt. <https://arxiv.org/abs/1705.00761>: JAUE, December 21 - 23 2010, pp. 101–123.
[14] C. Blake and C. J. Merz, “UCI repository of machine learning databases,” 1998.
[15] D. Barbara, Y. Li, and J. Couto, “Coolcat: an entropy-based algorithm for categorical clustering,” in *Proceedings of the eleventh international conference on Information and knowledge management*. ACM, 2002, pp. 582–589.
[16] P. Andritsos, P. Tsaparas, R. J. Miller, and K. C. Sevcik, “Limbo: Scalable clustering of categorical data,” in *Advances in Database Technology-EDBT 2004*. Springer, 2004, pp. 123–146.
[17] L. Xia, J. Sheng-Yi, and S. Xiao-Ke, “A novel fast clustering algorithm,” in *Artificial Intelligence and Computational Intelligence, 2009. AICI’09. International Conference on*, vol. 4. IEEE, 2009, pp. 284–288.
[18] J. Al-Shaqsi and W. Wang, “A clustering ensemble method for clustering mixed data,” in *Neural Networks (IJCNN), The 2010 International Joint Conference on*. IEEE, 2010, pp. 1–8.
[19] S.-Y. Jiang and X. Li, “A hybrid clustering algorithm,” in *Fuzzy Systems and Knowledge Discovery, 2009. FSKD’09. Sixth International Conference on*, vol. 1. IEEE, 2009, pp. 366–370.
[20] K. Wang, C. Xu, and B. Liu, “Clustering transactions using large items,” in *Proceedings of the eighth international conference on Information and knowledge management*. ACM, 1999, pp. 483–490.
[21] C.-H. Yun, K.-T. Chuang, and M.-S. Chen, “An efficient clustering algorithm for market basket data based on small large ratios,” in *Computer Software and Applications Conference, 2001. COMPSAC 2001. 25th Annual International*. IEEE, 2001, pp. 505–510.
[22] Y. S. Koh and R. Pears, “Transaction clustering using a seeds based approach,” in *Advances in Knowledge Discovery and Data Mining*. Springer, 2008, pp. 916–922.
[23] V. Ganti, J. Gehrke, and R. Ramakrishnan, “Cactus clustering categorical data using summaries,” in *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 1999, pp. 73–83.
[24] Y. Yang, X. Guan, and J. You, “Clope: a fast and effective clustering algorithm for transactional data,” in *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2002, pp. 682–687.
[25] M. J. Zaki and M. Peters, “Clicks: Mining subspace clusters in categorical data via k-partite maximal cliques,” in *Data Engineering, 2005. ICDE 2005. Proceedings. 21st International Conference on*. IEEE, 2005, pp. 355–356.

- [26] A. Ahmad and L. Dey, "A method to compute distance between two categorical values of same attribute in unsupervised learning for categorical data set," *Pattern Recognition Letters*, vol. 28, no. 1, pp. 110–118, 2007.
- [27] E. Abdu and D. Salane, "A spectral-based clustering algorithm for categorical data using data summaries," in *Proceedings of the 2nd Workshop on Data Mining using Matrices and Tensors*. ACM, 2009, p. 2.
- [28] A. Nemalhabib and N. Shiri, "Cluc: a natural clustering algorithm for categorical datasets based on cohesion," in *Proceedings of the 2006 ACM symposium on Applied computing*. ACM, 2006, pp. 637–638.
- [29] K. Lu, "Clustering transactions using context-based distance," no. 143, pp. 123–138, 2012.
- [30] G. Ivchenko and S. Honov, "On the jaccard similarity test," *Journal of Mathematical Sciences*, vol. 88, no. 6, pp. 789–794, 1998.
- [31] F. Alqadah and R. Bhatnagar, "Similarity measures in formal concept analysis," *Annals of Mathematics and Artificial Intelligence*, vol. 61, no. 3, pp. 245–256, 2011.
- [32] Y. S. Koh and R. Pears, "Rare association rule mining via transaction clustering," in *Proceedings of the 7th Australasian Data Mining Conference-Volume 87*. Australian Computer Society, Inc., 2008, pp. 87–94.

Efficient Community Detection Algorithm with Label Propagation using Node Importance and Link Weight

Mohsen Arab, Mahdiah Hasheminezhad*
Department of Computer Science
Yazd University, Yazd, Iran

Abstract—Community detection is a principle tool for analysing and studying of a network structure. Label Propagation Algorithm (LPA) is a simple and fast community detection algorithm which is not accurate enough because of its randomness. However, some advanced versions of LPA have been presented in recent years, but their accuracy need to be improved. In this paper, an improved version of label propagation algorithm for community detection called WILPAS is presented. The proposed algorithm for community detection considers both nodes and links important. WILPAS is a parameter-free algorithm and so requires no prior knowledge. Experiments and benchmarks demonstrate that WILPAS is a pretty fast algorithm and outperforms other representative methods in community detection on both synthetic and real-world networks. More specifically, experiments show that the proposed method can detect the true community structure of real-world networks with higher accuracy than other representative label propagation-based algorithms. Finally, experimental results on the networks with millions of links reveal that the proposed algorithm preserve nearly linear time complexity of traditional LPA. Therefore, the proposed algorithm can efficiently detect communities of large-scale social networks.

Keywords—Label propagation; node importance; link weight

I. INTRODUCTION

Many complex systems can be modelled as networks with nodes for entities and edges for the connections between them. Many real-world networks have community structure. Communities can be found in many complex systems such as social and biological networks, the internet, food webs and so on. Nodes of a community have often several characteristics in common.

By now, many different methods have been proposed for community detection. In 2002, Newman and Girvan devised a divisive algorithms using centrality indices to find community boundaries [1]. This index called edge betweenness and it refers to the number of shortest paths between all pairs of nodes that run along the edge. The edge with highest edge betweenness is removed in iterative steps until no edges remain. This process takes $O(m^2n)$ which makes it impractical to be run on the networks with more than 30,000 nodes. In 2004, a measure called modularity was introduced to evaluate a given partition of a network into communities [2]. So many methods were presented for modularity optimization [3], [4], [5]. Aside from modularity optimization, a variety of different algorithms such as graph partition-based methods [6], [7], [8] and density-based methods[9], [10] and label propagation algorithm (LPA)

[11] have been presented for community detection.

Among all the community detection methods, LPA is one of the fastest algorithms. LPA algorithm is simple and its time complexity is nearly linear time. However because of randomness, the detected communities have poor stability. That is, LPA may find different communities in different runs. In some runs, small communities are merged with big ones forming “monster” communities which is a drawback of LPA [12].

The LPA can be described as follows. Initially, each node is assigned a unique numeric label. At each iterative step, each node updates its label to the most frequent label from its neighbours in a random order. When there are multiple most frequent labels, the node will randomly pick one of them. Relabeling continues until the label of each node is its most frequent label among its neighbours. Finally, the nodes with the same label are considered in the same community. In fact, there are two sources of randomness in LPA which make it unstable and inaccurate. First source is random update order of nodes and the second one is randomly selecting one label when there are multiple most frequent labels to choose.

In this paper, a novel label propagation method for community detection called WILPAS is introduced. WILPAS algorithm has two stages. Let $l(v)$ be the label of node v . In the first stage, two sources of randomness of LPA are eliminated to increase accuracy. That is, firstly, random node sequence for label updating of LPA is replace by one specific update order. Secondly, WILPAS presents a novel label updating mechanism based on both node importance and link strength which makes the second source of randomness very unlikely to happen. The first stage of WILPAS is called weighted importance label propagation algorithm (WILPA).

Resulted communities from the first stage (WILPA) might be sub-communities of real ones. Therefore, in stage two of WILPAS, detected labels of nodes during the first stage are injected as a seed into a method called LPA_d . In fact, LPA_d is the same as traditional LPA in using random update order and the traditional label updating formula, but with one difference. When half of the neighbours of a node v have label $l(v)$, LPA_d does not update its current label $l(v)$. As it will be shown later, this change can avoid possible label oscillations in stage two of WILPAS.

Extensive experimental studies demonstrate that WILPAS is a pretty fast algorithm and it can get better community

detection results comparing with several label propagation based algorithms on both synthetic and real-world networks.

This paper is structured as follows. In Section II, related works in the field are listed. Some notions are defined in Section III. In Section IV the proposed method (WILPAS) is presented. The time complexity of proposed method is stated in Section V. Experimental results of comparing the proposed method with some famous methods in this area are discussed in Section VI. Finally, conclusion is given in Section VII.

II. RELATED WORKS

In 2007, Raghaval et al.[11] proposed Label Propagation Algorithm (LPA) for community detection. LPA can be summarized as four following steps:

- 1) Initialize every node with a unique label.
- 2) Arrange the nodes in a random order.
- 3) For every node in that random order, set its label with the one which is the most frequent label among its neighbours.
- 4) If every node has a label that the maximum number of their neighbours have, then stop the algorithm; else go to step 2.

The formula of label updating for LPA is as follows:

$$l(v) = \operatorname{argmax}_l \sum_{u \in N^l(v)} 1, \quad (1)$$

where $N^l(v)$ indicates the set of neighbours of node v with label l . This is LPA's asynchronous version. Since synchronous version has potential label oscillations as discussed in [11], this version is not considered. As discussed earlier LPA has two types of randomness. Unfortunately, randomness of LPA may result in missing small communities and even getting trivial solution in which all nodes are assigned the same label [12]. Moreover, it makes the algorithm unstable such that different communities may be detected in different runs of the algorithm.

Zhang et al. generalized LPA to weighted networks by calculating the probability value of every label [13]. The label updating formula in this case is changed as follows:

$$l(v) = \operatorname{argmax}_l \sum_{u \in N^l(v)} w_{vu}, \quad (2)$$

where w_{vu} indicates the weight of the edge between nodes v and u .

Barber and Clark proposed modularity-specialized algorithm (LPAm) to constrain the label propagation process [14]. Their algorithm is near-linear time, but it may get stuck in poor local maxima in the modularity space. To scape local maxima, Liu et al. introduced an advanced modularity-specialized label propagation algorithm called LPAm+ [15]. LPAm+ combines LPAm with multistep greedy agglomerative algorithm to get higher modularity values. Thus, LPAm+ does not guarantee near-linear time complexity [16]. Xing et al. presented a node influence based label propagation algorithm called NIBLPA [17]. NIBLPA defines two concepts node

influence and label influence for specifying node orders and label choosing mechanism respectively. Zhang et al. proposed a label propagation algorithm with prediction of percolation transition named LPAP [16]. They transformed the process of label propagation into network construction process. Using this prediction process of percolation transition, they tried to delay the occurrence of trivial solutions. Sun et al. proposed a centrality-based label propagation called CenLP [18]. They presented a new measure for computing the centrality of nodes. Based on these centrality values, one specific update order in addition to node preference values are specified in order to improve traditional LPA.

III. TERMINOLOGY

Let $G = (V, E)$ be an undirected network. The number of nodes and links of G is denoted by n and m , respectively. Let d_v be the degree of node v in the network. Degrees of node v within and outside of its community are denoted by d_v^{in} and d_v^{out} , respectively. Mixing parameter μ for each node v is defined as $\frac{d_v^{out}}{d_v}$. The set of all neighbours of node v is denoted by $N(v)$. Internal and external links respectively refers to the links within and between communities.

IV. PROPOSED METHOD (WILPAS)

The proposed algorithm has two stages. At first stage, in order to increase the quality of detected communities of LPA, one specific node order for label updating and a novel formula for selecting labels for nodes is introduced. The novel formula for label updating is based on the weights of links and importance of nodes. Therefore, the first stage of the proposed algorithm using these two modifications in traditional LPA is called weighted importance label propagation algorithm (WILPA). The detected communities resulted from stage one might be sub-communities of real ones. Therefore, in stage two, found labels of nodes resulted from stage one (WILPA) will be injected as a seed into a method similar to traditional LPA. The second stage which has a slight difference with traditional LPA is called LPA_d . By presenting these two stages, the proposed method is completed. Since detected labels of WILPA algorithm are injected as a seed into LPA_d algorithm, the proposed method is called WILPAS.

A. Weighting Measure for Links

There are several normalized similarity measure to assign weights to an edge (u, v) such as cosine [19]. Cosine similarity measure between two nodes u and v is defined as follows:

$$\operatorname{cosine}(u, v) = \frac{|N(u) \cap N(v)|}{\sqrt{|N(u)| |N(v)|}}, \quad (3)$$

Where $||$ indicates the cardinality of a set. Using cosine may result in assigning zero values to some links. Thus, instead an extended version of cosine [18] is chosen to assign non-zero weights to links. This measure is called structural similarity and is defined as follows:

$$\sigma(u, v) = \frac{|\Gamma(u) \cap \Gamma(v)|}{\sqrt{|\Gamma(u)| |\Gamma(v)|}}, \quad (4)$$

where $\Gamma(u) = N(u) \cup \{u\}$.

B. Stage One of WILPAS (WILPA)

The stage one of WILPAS (WILPA) improves traditional LPA with two modifications: The first modification is presenting one specific node order for updating labels instead of random order. The second one is presenting a novel formula for selecting new labels of nodes. This novel label updating formula considers both importance values of neighbour nodes and weights of neighbour links of a node to select its new label.

1) *Specific update order:* In the proposed method, nodes are rearranged such that important nodes update their labels first. The degree of each node v (i.e d_v) is chosen as its importance value. Among several nodes with equal degrees, those whose neighbours have higher degrees are more important. Thus, an extended version of importance value of each node v ($EI(v)$) is defined as follows:

$$EI(v) = d_v + \sum_{u \in N(v)} d_u \quad (5)$$

Therefore, order of nodes for label updating in the proposed method is specified in descending order of extended importance values of nodes.

2) *Novel label updating formula:* In WILPA, instead of selecting the most frequent label among neighbours of a node as its new label, a novel label choosing mechanism is adopted. This mechanism considers both node importance and link importance for selecting the new label.

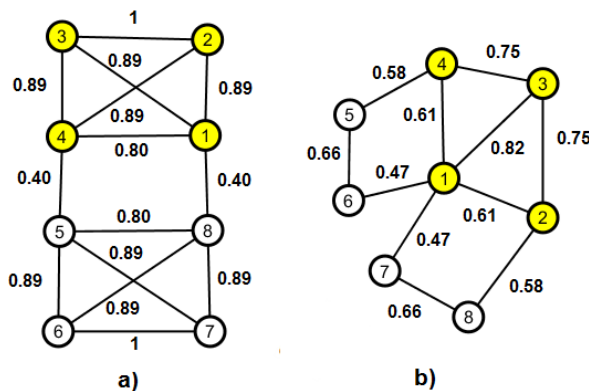


Fig. 1. Two sample networks.

Consider Fig. 1 a with two real communities $\{1, 2, 3, 4\}$ and $\{5, 6, 7, 8\}$. Suppose that during traditional label propagation process (using label updating formula 1) nodes 1-4 have label 1 and nodes 5-8 have a label equal to their own numbers. That is, denoting label of node v by $l(v)$: $l(1) = l(2) = l(3) = l(4) = 1$ and $l(5) = 5, l(6) = 6, l(7) = 7, l(8) = 8$. Moreover, suppose that the update order of nodes is 8, 5, 6, 7. That is, node 8 should update its label first, then node 5 and so on. Node 8 has four neighbours with labels 1, 5, 6 and 7. If node 8 selects label 1 randomly, then labels of nodes 5, 6 and 7 will be 1 as well, since label 1 will be the most frequent neighbour label for them. Thus, trivial solution (forming one big community) will be obtained.

To avoid trivial solution, the chance of propagation of labels between different communities should be decreased. One way to do this is using weights of links in the label propagation process. The idea behind using weights is that two endpoints of an internal link share more common friends to each other than two endpoints of an external link. Thus, if weight of a link is defined based on the ratio of common friends between its two endpoint nodes, then internal links are more likely to get higher weight than external ones. Thus, by considering the weights of links in label propagation process, one can expect that propagation of labels between two different communities will be less likely.

In Fig. 1(a), let this time take into account the structural similarity weight 4 of links and choose formula 2 as label updating mechanism. In this situation, the weights of the links connecting labels 1, 5, 6 and 7 to node 8 are 0.40, 0.80, 0.89 and 0.89. Therefore, node 8 will select one of two labels 6 or 7, because their corresponding weight 0.89 is maximum. Either of two labels 6 or 7 is chosen by node 8, the other three nodes 5, 6 and 7 will choose that label as well. Therefore, two real communities will be detected correctly.

However, using weighted label propagation can decrease the chance of propagation of labels between different communities, but in sparse real-world networks, this strategy may cause real communities to break apart into several sub-communities. For example, consider the network in Fig. 1(b) with one single community. Like previous example, let consider nodes 1-4 have label 1 and other nodes have label equal to their own numbers. Using weighted label propagation strategy will result in finding three communities $\{1, 2, 3, 4\}, \{5, 6\}, \{7, 8\}$. This is because the weights of two links (5, 6) and (7, 8) are greater than their neighbour links. Therefore, nodes 5 and 6 will choose the labels of each other. Similarly, both nodes 7 and 8 will adopt the same label 7 or 8 as their final label. Therefore, weighted label propagation strategy may divide some communities of a real-world network into several sub-communities.

To resolve the mentioned problem, one idea is to consider degrees of nodes as their importance values in label updating formula. This solution is based on this intuitive idea that in each network, there are some important nodes with high degree which play crucial role in spreading information, viral marketing, etc. Therefore, nodes with higher degrees are more likely to be centers of communities [18]. It is obvious that in social networks, a famous person or a celebrity with more friends and connections has more impact on each of his friends than a person with just a few friends.

Therefore, on the one hand, with weighted label propagation external links would have low effect in spreading labels between different communities. Thus, this idea can reduce the formation of monster communities. On the other hand, most important nodes (such as nodes with high degrees) play very crucial role in formation of communities. Therefore, the degrees of nodes should be considered in label updating formula as well. By taking into consideration both weights of links and degrees of nodes, the label updating formula of the proposed method is defined as follows:

$$l(v) = \operatorname{argmax}_l \sum_{u \in N^l(v)} w_{vu} * d_u \quad (6)$$

Therefore, each neighbour u of node v has an impact in defining $l(v)$ based on its importance value (d_u) and the weight of corresponding link (w_{vu}). As discussed above, the degree of each node is considered as its importance value. Therefore, the first stage of the proposed method is completed. The pseudo-code of WILPA is presented in Algorithm 1.

C. Stage Two of WILPAS (LPA_d)

Resulted communities from the first stage (WILPA) might be sub-communities of real ones. Therefore, in stage two of WILPAS, detected labels of nodes during the first stage are injected as a seed into a method similar to original label propagation algorithm. To be more accurate, LPA_d is the same as original LPA in using random update order and label updating formula, but with one difference. When half of the neighbours of a node v is the same as $l(v)$, LPA_d keeps its current label.

Consider the network in Fig. 2. WILPA algorithm as the stage one of WILPAS method detects two communities on this network which are shaded with colors green and yellow. If original LPA is applied on these found labels, final labels of two nodes 6 and 11 will be either green or yellow. This is because of the fact that two nodes 6 and 11 are connected to two different communities with equal number of links. In this situation, if LPA_d algorithm is used instead of traditional LPA, then labels of two nodes 6 and 11 will be fixed as green. Hence, LPA_d algorithm by avoiding possible label oscillations and unnecessary iterations can increase stability of detected communities and reduce the number of iterations of the proposed method.

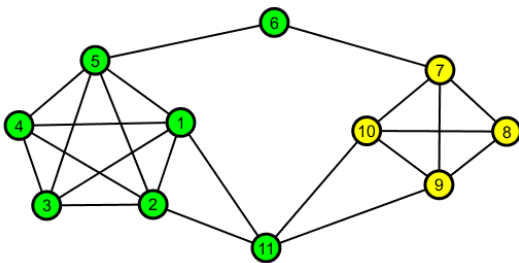


Fig. 2. A network with two detected communities by stage one of WILPAS.

D. Pseudo-code for WILPAS

Pseudo-code of WILPA, LPA_d and WILPAS are presented in Algorithms 1, 2 and 3, respectively.

V. TIME COMPLEXITY

In this section, the time complexity of the proposed method is discussed.

```

1 For each node  $v$  set  $l(v) = v$  /* Initialization of labels. */
2 Arrange nodes based on descending order of their EI values
  (formula 5) and put them into array X.
3 repeat
4   foreach vertex  $v$  in X do
5     Update  $l(v)$  using novel proposed label updating
     formula 6
6   end
7 until labels of nodes do not change any more;
```

Algorithm 1: WILPA

```

1 Arrange nodes randomly and put them into array Y.
2 repeat
3   foreach vertex  $v$  in Y do
4     if  $|N^{l(v)}(v)| < \frac{d_v}{2}$  then
5       Update  $l(v)$  using traditional label updating
       formula 1
6     end
7   end
8 until labels of nodes do not change any more;
```

Algorithm 2: LPA_d algorithm

```

1 Compute structural weights of all links using formula 4
  (e.g. Algorithm 4 can be used for computing weights).
2 WILPA() (Algorithm 1) /* Stage 1*/
3  $LPA_d$ () (Algorithm 2) /* Stage 2*/
```

Algorithm 3: WILPAS Algorithm

A. Time complexity of weighting all links

For computing the weight of a link, the number of common friends between its two endpoints should be counted. Algorithm 4 is a simple algorithm to do that. As it can be seen from the algorithm, for computing the weight of link (v, u) it is enough to explore the set C of u 's neighbours and then count the number of nodes in C which are neighbours to v as well. This is done in $O(d_u)$. Since there are d_u neighbour links for u , computing the weights of all of its neighbour links can be done in $O(d_u^2)$. Therefore, total time complexity of this simple weighting algorithm is

$$\sum_{u=1}^n d_u^2 \quad (7)$$

Space complexity of this algorithm using adjacency list is

$$O(m) \quad (8)$$

Checking whether two nodes z and v are neighbours can be done in $O(d_z)$ using adjacency list. But, in order to do that in $O(1)$, an extra array named 'mark' is used as follows. For each node v , at first, in line 3 of Algorithm 4, each of its neighbour u is marked as v . Then, in line 8 the adjacency of two nodes z and v is checked in $O(1)$ by comparing the content of 'mark' array of index z with v .

B. Time Complexity of Two Stages of WILPAS

In the first stage (WILPA), at first, all nodes are arranged based on their EI values. This can be done with time complexity $O(n \log n)$. Time complexity of each iteration of label

```

1 foreach vertex v do
2   foreach neighbor u of v do
3     mark[u]=v
4   end
5   foreach neighbor u of v do
6     Cfriends=0;
7     foreach neighbor z of u do
8       if mark[z]==v then
9         ++Cfriends;
10      end
11    end
12    /* compute weight of edge(u,v) using equation 4 */
13  end
14 end

```

Algorithm 4: A Simple Weighting Algorithm

updating in WILPA is the same as traditional LPA which is $O(m)$ [11]. This is because time complexity of computing new label $l(v)$ in formula 6 and formula 1 is the same. Therefore, time complexity of WILPA is

$$O(n \log n) + O(R_1 m), \quad (9)$$

where R_1 is the number of iterations of WILPA. Similarly, each iteration of LPA_d requires $O(m)$ time. Thus, time complexity of LPA_d is

$$O(R_2 m), \quad (10)$$

where R_2 is the number of iterations of LPA_d .

C. Total Time Complexity of WILPAS

Total time complexity of WILPAS is the summation of time complexities of computing weights, WILPA and LPA_d which is as follows:

$$O(n \log n) + O(R_1 m) + O(R_2 m) + O\left(\sum_{u=1}^n d_u^2\right) \quad (11)$$

It is important to note that in practice in most cases both R_1 and R_2 are less than 10. Moreover, real networks are often sparse, i.e. $m = O(n)$. In addition, as it will be shown in experiments section, the weighting Algorithm 4 consumes less than 25 seconds for finding the weights of all links of a network with 500,000 nodes and around 10 million links. Therefore, as it will be demonstrated later, WILPAS is pretty fast in practice, even with existing term $\sum_{u=1}^n d_u^2$.

VI. EXPERIMENTS

This section evaluates the effectiveness and the efficiency of the proposed algorithm. Several experiments on both synthetic networks and well-known real-world networks are conducted. Moreover, the performance of WILPAS with LPA, CenLP, LPAp, LPAm and NIBLPA are compared. All the simulations are carried out in a desktop pc with Pentium Core2, 1.8 GHZ processor and 4GB of RAM under Windows 8.1 OS.

In this paper, normalized mutual information (NMI) [20] is used as the evaluation measure which is currently widely used in measuring the quality of detected communities. NMI allows us to measure the amount of information common to

two different network partitions. Accordingly, if a network has a known community structure, one can explore the efficacy of the algorithm by comparing known real partition with the partition found by that algorithm. When the found partition matches the real one, then $NMI=1$, and when two partitions are independent of each other, then $NMI=0$.

A. Test on Synthetic Networks

In this section, LFR benchmark networks [21] are chosen which are currently the most commonly used synthetic networks in community detection. The parameters of LFR benchmark networks are as follows: number of nodes n , the average degree k , maximum degree $maxk$, mixing parameter μ . Moreover, $minc$ and $maxc$ refer to the minimum and maximum values for community sizes, respectively.

Three ranges for different community sizes are used which are indicated by the letters S (stays for small), B (stays for big) and VB (stays for very big). The ranges of community sizes for three letters S, B and VB are $[min, max] = [10, 50]$, $[min, max] = [20, 100]$ and $[min, max] = [200, 1000]$, respectively. For each type of networks, 10 samples are generated and on each sample, each tested label propagation-based algorithm is run 10 times. Then, the average of these 100 NMI values are reported as output. In this paper for all the networks with $n \geq 100,000$, the average degree $k = 40$ and the letter VB are used, i.e. community sizes of these networks range between 200 and 1000 where average degree of nodes is 40.

Fig. 3 and 4 show the accuracy of the mentioned methods on the networks with size of 1000. One can observe that for $n = 1000$, when $\mu \leq 0.50$ three methods WILPAS, LPAm and CenLP find communities pretty well. However, when communities are big, for $\mu > 0.50$, LPAm gets better results (see Fig. 4).

Fig. 5 and 6 show the accuracy of methods when $n = 10,000$. From these two figures it can be observed that when $n = 10,000$, WILPAS outperforms other methods. CenLP is the second most accurate method for community detection on this network. On this network, NIBLPA shows poor performance in community detection.

Fig. 7 demonstrates the NMI results for the three most accurate tested label propagation methods i.e. WILPAS, CenLP and LPAm for a network with $n = 100,000$, $k = 40$, $[min, max] = [200, 1000]$. As it can be observed from this figure, WILPAS achieves higher NMI values than CenLP and LPAm. CenLP shows more accuracy than LPAm except for $\mu = 0.70$. The detailed information about the results is displayed in Table I.

TABLE I. NMI RESULTS OF THREE METHODS WILPAS, CenLP AND LPAm ON THE NETWORK WITH $n = 100,000$

μ	LPAm	CenLP	WILPAS
0.40	0.9963	1	1
0.45	0.9955	1	1
0.50	0.9946	1	1
0.55	0.9927	1	1
0.60	0.9818	0.9999	1
0.65	0.9527	0.9970	1
0.70	0.8277	0	0.9997

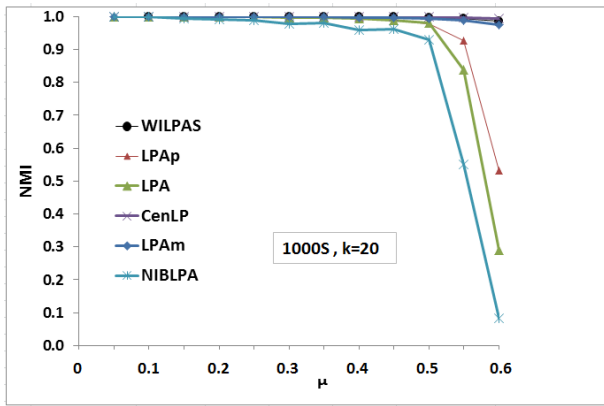


Fig. 3. Comparing different label propagation-based algorithms on the network with $n = 1,000$ where communities are small.

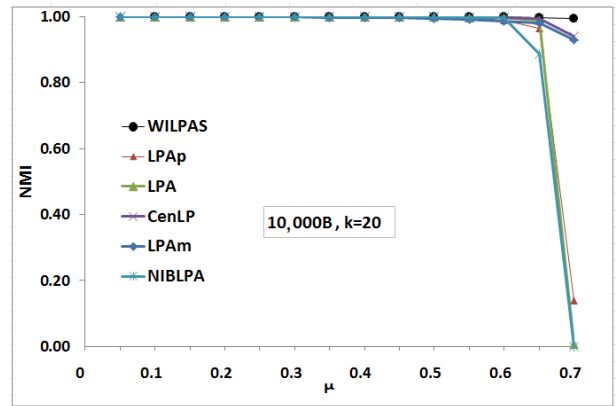


Fig. 6. Comparing different label propagation-based algorithms on the network with $n = 10,000$ where communities are big.

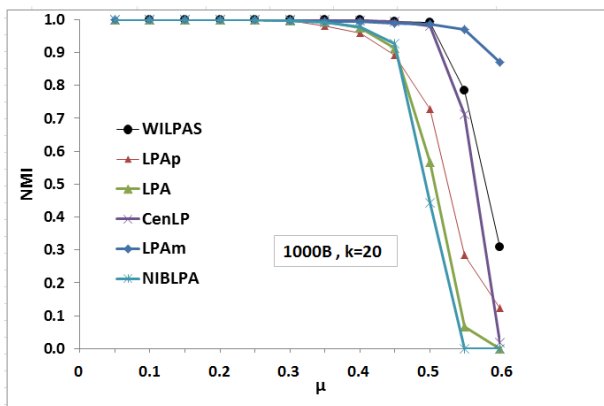


Fig. 4. Comparing different label propagation-based algorithms on the network with $n = 1,000$ where communities are big.

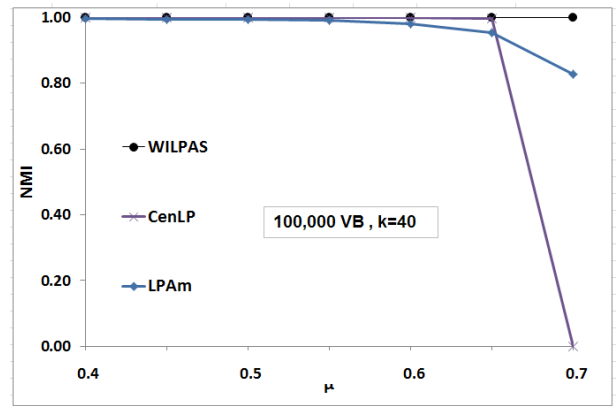


Fig. 7. Comparing different label propagation-based algorithms on the network with $n = 100,000$ where communities are very big.

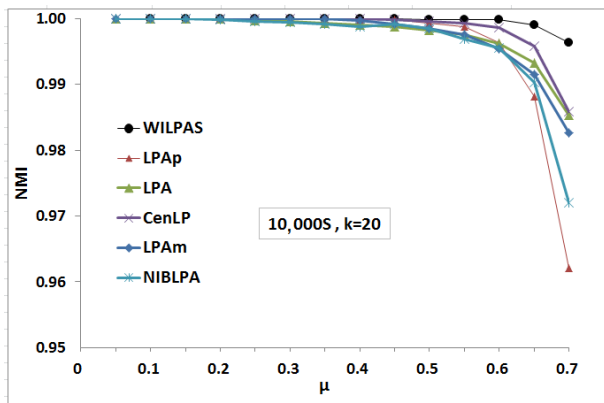


Fig. 5. Comparing different label propagation-based algorithms on the network with $n = 10,000$ where communities are small.

are shown in Table II. The NMI results of all tested label propagation-based methods are displayed in Table III.

Each method is run 10 times on each real network, then the average NMI results are reported. The number in the {} for CenLP, NIBLPA and WILPAS in Table III shows the number of found communities by these three deterministic methods. Since LPA, LPAP and LPAM detect different partitions on the same network for each run, they are ignored. The maximum resulted NMI values on each network has been bold in Table III.

TABLE II. REAL-WORLD NETWORKS WITH KNOWN COMMUNITY STRUCTURES

Network	Nodes	Links	Communities
Karate [22]	34	78	2
Dolphin [23]	62	159	2
Football [1]	115	615	12
Polblog [24]	1490	16715	2

B. Experiment on Real-world Networks

In this section, the evaluation of the above methods on real-world networks which their communities are already known is discussed. Zachary Karate club [22], American college football [1], Dolphin social network [23] and Polblog [24] are four famous networks in the field. The details of these networks

1) *Zachary Karate club*: The well-known Karate club network of Zachary [22] is a standard benchmark for community detection. Zachary observed 34 members of a karate club in the United States over two years. Because of a disagreement between administrator and instructor of the club, a new club was formed by the instructor by taking about the half of the

original club members. The edge between nodes (members) of this network represent the social interactions between the members outside the club. These two original communities are specified with the shapes 'square' and 'circle' in Fig. 8.

As it can be observed from Table III, WILPAS is the only method that finds exactly the two real communities of Karate club network with NMI=1. CenLP is the second best method with NMI=0.60 with finding four communities. NIBLPA has poor performance on Karate network with NMI=0.21. The sets of sizes of detected communities by WILPAS, CenLP and NIBLPA are {16,18}, {12,5,4,13} and {2,3,29}, respectively. Fig. 8 shows the two detected communities by WILPAS on Karate club network with different colors.

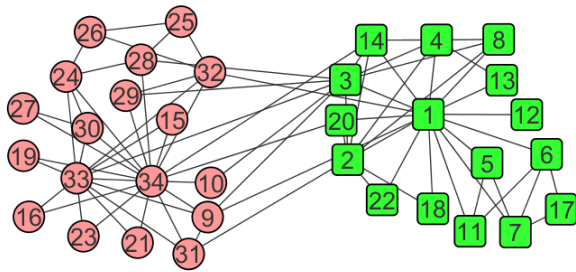


Fig. 8. Result of community detection by WILPAS on Karate network. Two real communities are specified with shapes 'circle' and 'square'. Two found communities are shaded with different colours. WILPAS method detects the two real communities exactly as it is.

2) *American college football*: Another well known benchmark for community detection is American college football network compiled by Girvan and Newman [1]. This network represents Division I games for the 2000 season. Nodes represent teams and the edges represent the games between teams. The teams belong to the conferences with 8 to 12 teams each. Since, games between the teams of the same conference are usually more frequent than the games between the teams of different communities, this network has community structure. As one can see from Table III both WILPAS and CenLP finds 13 communities on this network. In fact, both WILPAS and CenLP gain the maximum NMI value 0.90 on this network. After these two methods, LPAm is the third accurate method with NMI=0.89.

3) *Dolphin social network*: Dolphin network [23] shows the frequent associations between 62 dolphins living in Doubtful Sound, New Zealand. Nodes are dolphins and the edges between nodes shows that the two corresponding dolphin were seen together more than expected by chance. After leaving one of dolphins, they separated in two communities. Two original communities are specified with shapes 'circle' and 'square' in Fig. 9. Three communities which are detected by WILPAS are specified with different colors.

From Table III one can observe that WILPAS achieves higher NMI value than other methods on the Dolphin network. Moreover, the number of detected communities by WILPAS is more close to two real communities of Dolphin network. LPAm fails to detect true communities with getting lowest NMI value 0.45.

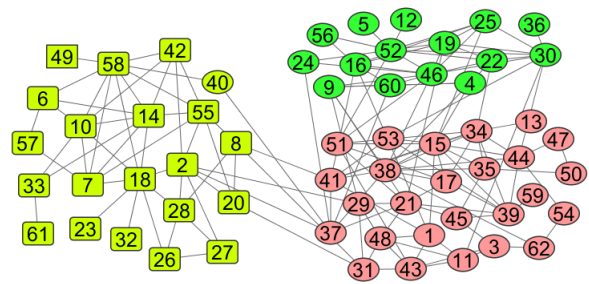


Fig. 9. Result of community detection by WILPAS on Dolphin network. Two real communities are specified with shapes 'circle' and 'square'. Three found communities are shaded with different colours.

4) *Polblogs network*: This network represents the links between weblogs about US politics preceding the US Presidential Election of 2004 [24]. The links were automatically extracted from a crawl of the front page of the weblogs. Each blog is labelled with '0' or '1' to indicate whether they are "liberal" or "conservative". This network can be considered both directed or undirected. In this paper, the undirected version of this network which has 1490 nodes and 16715 links is considered. Since nodes with degree zero makes this network disconnected, when comparing the performance of methods, these nodes are ignored. Thus, by removing 266 nodes with degree zero, the resulted network with 1224 nodes is considered for testing and comparing community detection methods. By doing this, the sizes of two real communities of Polblog are 588 and 636.

CenLP and WILPAS achieve NMI values 0.71 and 0.70 on Polblog network respectively. However CenLP gets a little more NMI value than WILPAS, but the number of detected communities of WILPAS is more close to two real communities of Polblog network. The sets of sizes of detected communities by WILPAS and CenLP are {552,2,670} and {559,2,4,659}, respectively.

In summary, when dealing with community detection on real networks, WILPAS outperforms other methods on Karate and Dolphin network, while CenLP has a little better accuracy than WILPAS on Polblog network. Both of these two methods has the same accuracy on Football network with finding 13 communities. The superiority of WILPAS on Karate and Dolphin networks is remarkable while superiority of CenLP on Polblog network is negligible. Moreover, while both of these two methods find 13 communities on Football network, the numbers of found communities of WILPAS on Karate, Dolphin and Polblog networks are more close to the numbers of real communities of these networks. These show that the proposed method WILPAS is the most accurate label propagation method in comparison to other tested methods for community detection on the real networks.

TABLE III. NMI RESULTS OF THE METHODS ON FOUR REAL NETWORKS WITH KNOWN COMMUNITY STRUCTURES

Networks/ methods	LPAm	LPap	WILPAS	LPA	NIBLPA	CenLP
Karate	0.55	0.56	1 ,{2}	0.70	0.21 {3}	0.60, {4}
Dolphin	0.45	0.55	0.66 ,{3}	0.52	0.50 {5}	0.61,{4}
Polblog	0.45	0.61	0.70,{3}	0.70	0.20 {9}	0.71 ,{4}
Football	0.89	0.88	0.90 ,{13}	0.87	0.78 {9}	0.90 {13}

C. Efficiency Analysis

To illustrate the running time of the proposed algorithm WILPAS and compare it with other algorithms, 10 networks using LFR software are produced, where the number of nodes $n = 100,000$ and the average degree $k = 40$ and $[minc, maxc] = [200, 1000]$. Fig. 10 plots the average running time of the proposed method WILPAS on these 10 synthetic networks compared with other five label propagation algorithms: LPA, LPAm, LPAp, CenLP and NIBLPA. As one can see from Fig. 10, method WILPAS is faster than LPAm but slower than LPA, LPAp and NIBLPA. In addition, it has a comparative execution time with CenLP.

Fig. 11 illustrates the running time of the weighting Algorithm 4 where n ranges from 100,000 to 500,000. As it can be seen from this figure, the weighting Algorithm 4 consumes less than 3.1 seconds for finding weights of this network with 100,000 nodes and around 2 million links. With increasing the number of node n to 500,000, the consumed time increase near linearly. Therefore, finding weights of all links of a network with 500,000 nodes and around 10 million links requires less than 25 seconds.

Similarly, for evaluating the scalability of WILPAS, the average running time of WILPAS on 10 LFR networks is reported where n ranges from 100,000 to 500,000. From Fig. 12 one can observe that the execution time of WILPAS scales approximately linearly with n , while it is less than double of execution time of LPA. As one can see from Fig. 12, WILPAS consumes less than 104 seconds for community detection on the network with 500,000 nodes and around 10 million links. This shows the efficiency and scalability of WILPAS in community detection.

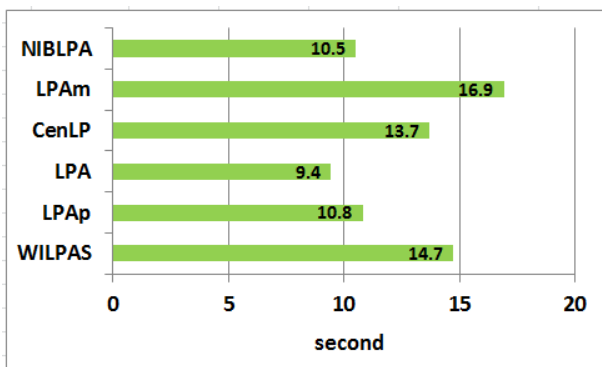


Fig. 10. The execution times of different methods on a network with $n=100,000$, $k=40$, $[minc, maxc] = [200, 1000]$.

VII. CONCLUSION

In this paper, a new label propagation algorithm called WILPAS is proposed. WILPAS presents specific update order and a novel label choosing formula in order to increase the accuracy of community detection. WILPAS is parameter-free that requires no prior knowledge. Experimental results on both synthetic and real-world tested networks demonstrate that WILPAS is the most accurate label propagation algorithm, while it is pretty fast. Moreover, finding communities of networks with around 10 million links in less than two minutes

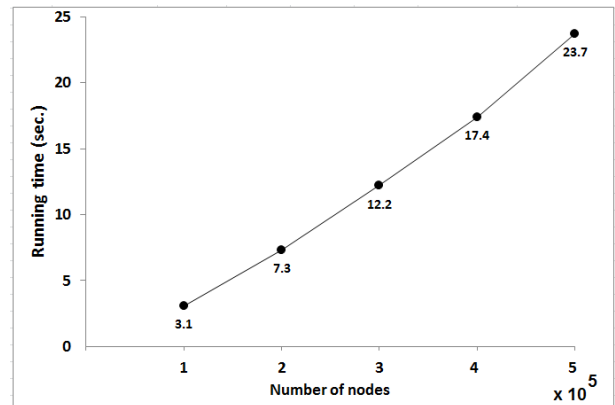


Fig. 11. The execution time of simple weighting Algorithm 4 with increasing n . The average degree $k=40$, $[minc, maxc] = [200, 1000]$. The number of nodes n ranges from 100,000 to 500,000.

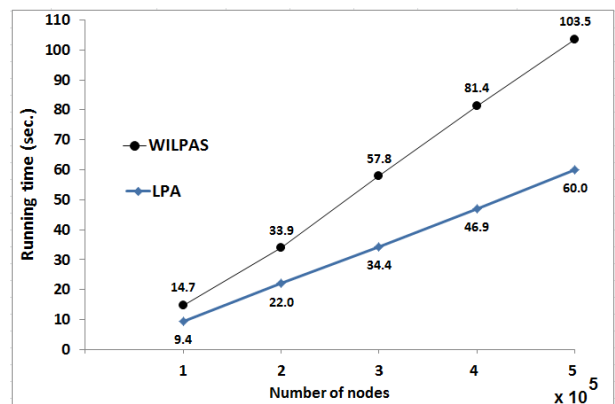


Fig. 12. The execution time in second for LPA and WILPAS on LFR benchmark with $k=40$, $[minc, maxc] = [200, 1000]$. The number of nodes n ranges from 100,000 to 500,000.

shows its scalability. Finally, experiments on several well-known real-world networks demonstrate that WILPAS outperforms other tested label propagation algorithms in finding true community structures of networks. In this paper, the communities should be distinct from each other. As future work, this algorithm can be extended to be used for overlapping (or fuzzy) community detection where each node may belong to several different communities.

REFERENCES

- [1] M. Girvan, M. E. Newman, Community structure in social and biological networks, Proceedings of the national academy of sciences 99 (12) (2002) 7821–7826.
- [2] M. E. Newman, M. Girvan, Finding and evaluating community structure in networks, Physical review E 69 (2) (2004) 026113.
- [3] G. Agarwal, D. Kempe, Modularity-maximizing graph communities via mathematical programming, The European Physical Journal B 66 (3) (2008) 409–418.
- [4] L. Bennett, S. Liu, L. Papageorgiou, S. Tsoka, A mathematical programming approach to community structure detection in complex networks, in: Symposium on Computer, no. June, 2012, pp. 17–20.
- [5] M. Arab, M. Afsharchi, Community detection in social networks using hybrid merging of sub-communities, Journal of Network and Computer Applications 40 (2014) 73–84.

- [6] B. W. Kernighan, S. Lin, An efficient heuristic procedure for partitioning graphs, *Bell system technical journal* 49 (2) (1970) 291–307.
- [7] G. W. Flake, S. Lawrence, C. L. Giles, Efficient identification of web communities, in: *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2000, pp. 150–160.
- [8] S. White, P. Smyth, A spectral clustering approach to finding communities in graph., in: *SDM*, Vol. 5, SIAM, 2005, pp. 76–84.
- [9] X. Xu, N. Yuruk, Z. Feng, T. A. Schweiger, Scan: a structural clustering algorithm for networks, in: *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, 2007, pp. 824–833.
- [10] H. Sun, J. Huang, J. Han, H. Deng, P. Zhao, B. Feng, gskeletonclu: Density-based network clustering via structure-connected tree division or agglomeration, in: *2010 IEEE International Conference on Data Mining*, IEEE, 2010, pp. 481–490.
- [11] U. N. Raghavan, R. Albert, S. Kumara, Near linear time algorithm to detect community structures in large-scale networks, *Physical review E* 76 (3) (2007) 036106.
- [12] I. X. Leung, P. Hui, P. Lio, J. Crowcroft, Towards real-time community detection in large networks, *Physical Review E* 79 (6) (2009) 066107.
- [13] ping Zhang, A., Ren, G., Cao, H., zhu Jia, B. and bin Zhang, S., 2013, May. Generalization of label propagation algorithm in complex networks. In *Control and Decision Conference (CCDC)*, 2013 25th Chinese (pp. 1306-1309). IEEE.
- [14] Barber, M.J. and Clark, J.W., 2009. Detecting network communities by propagating labels under constraints. *Physical Review E*, 80(2), p.026129.
- [15] Liu, X. and Murata, T., 2010. Advanced modularity-specialized label propagation algorithm for detecting communities in networks. *Physica A: Statistical Mechanics and its Applications*, 389(7), pp.1493-1500.
- [16] Zhang, A., Ren, G., Lin, Y., Jia, B., Cao, H., Zhang, J. and Zhang, S., 2014. Detecting community structures in networks by label propagation with prediction of percolation transition. *The Scientific World Journal*, 2014.
- [17] Xing, Y., Meng, F., Zhou, Y., Zhu, M., Shi, M. and Sun, G., 2014. A node influence based label propagation algorithm for community detection in networks. *The Scientific World Journal*, 2014.
- [18] Sun, H., Liu, J., Huang, J., Wang, G., Yang, Z., Song, Q. and Jia, X., 2015. CenLP: A centrality-based label propagation algorithm for community detection in networks. *Physica A: Statistical Mechanics and its Applications*, 436, pp.767-780.
- [19] Z. Liu, P. Li, Y. Zheng, M. Sun, Community detection by affinity propagation, *Tech. rep.*, Technical Report (2008).
- [20] L. Danon, A. Diaz-Guilera, J. Duch, A. Arenas, Comparing community structure identification, *Journal of Statistical Mechanics: Theory and Experiment* 2005 (09) (2005) P09008.
- [21] A. Lancichinetti, S. Fortunato, F. Radicchi, Benchmark graphs for testing community detection algorithms, *Physical review E* 78 (4) (2008) 046110.
- [22] W. W. Zachary, An information flow model for conflict and fission in small groups, *Journal of anthropological research* (1977) 452–473.
- [23] D. Lusseau, M. E. Newman, Identifying the role that animals play in their social networks, *Proceedings of the Royal Society of London B: Biological Sciences* 271 (Suppl 6) (2004) S477–S481.
- [24] Adamic, L.A. and Glance, N., 2005, August. The political blogosphere and the 2004 US election: divided they blog. In *Proceedings of the 3rd international workshop on Link discovery* (pp. 36-43). ACM.

An Indefinite Cycle Traffic Light Timing Strategy

Ping Guo, Daiwen Lei, Lian Ye
College of Computer Science
Chongqing University
Chongqing, 400044
China

Abstract—Intelligent transportation signal control plays an important role in reducing traffic congestion and improving road capacity. The key of signal control is to adjust the traffic lights appropriately according to the traffic flow, which is an adaptive control. In this paper, we propose a new timing strategy. This strategy includes green time optimization and lane combination calculation. According to the real-time traffic flow, we optimize green time and calculate lane combination to adjust the cycle and then we can get the timing plan. The simulation results of random data and actual traffic data show that the strategy we proposed can increase traffic efficiency by more than 15% at intersections, reduce vehicle detention, and relieve traffic congestion.

Keywords—Traffic light; timing strategy; adaptive control; signal control; intelligent transportation

I. INTRODUCTION

With the development of the economy and the expansion of the city scale, the rapid growth of urban car has caused many problems such as traffic congestion and environmental pollution. The intelligent transportation system provides effective support for solving traffic congestion. It integrates multiple disciplines in the field of information and aims to establish a full-coverage, real-time, accurate, and efficient integrated transportation and management system [1], [2].

In recent years, the research on intelligent transportation systems has achieved many results. Collotta M. proposed a novel approach to dynamically manage the traffic lights cycles and phases in an isolated intersection [3]. This method is a traffic lights dynamic control system that combines Wireless Sensor Network (WSN) for real-time traffic monitoring with multiple fuzzy logic controllers. This system outperforms other solutions in the literature, since it significantly reduces the vehicles waiting times. Li C. focused on an optimum route search in the in-vehicle routing guidance system. For the dynamic route guidance system (DRGS), it should provide dynamic routing advice based on real-time traffic information and traffic conditions, such as congestion and roadwork [4]. Darwish T. S. J. applied big data to intelligent transportation and proposed a real-time ITS (intelligent transportation system) big data analysis method for vehicle internet environment [5]. The Intelligent Transportation Project RoadEye was proposed to solve various traffic problems and make traffic safer. The project was able to detect weather conditions, maintain a safe distance and fixed speeds between vehicles [6]. Hassouneh Y. [7] proposed an ITS that can adapt its behavior in response to environmental changes. Cao Z. proposed a pheromone-based traffic management framework to reduce traffic congestion and unify dynamic vehicle routing and traffic control strategies [8].

The intersection is the node that transforms the traffic flow and plays an important role in the road network [9]. The control of traffic lights at intersections is the basis and guarantee for efficient and orderly operation of urban traffic. The merits of signal control methods directly determine the smoothness of traffic flow at intersections [10]. Traffic light control [11] at the intersection mainly includes timing control, induction control and adaptive control [12]. The timing control method is based on the use of a fixed green light time in each phase, which is simple and easy to maintain. Induction control is to set up a vehicle detector on the entrance of the intersection, and the traffic light timing plan changes with traffic flow in real time. This control is more suitable for situations where traffic flow at intersections is not obvious, saturation is not high, or the traffic flow in each phase is quite different. Adaptive control is to collect the traffic flow information in all directions of the intersection in real time, and calculate the green light time according to the prediction model to adapt to the change of intersection traffic. Its advantage is that it can reduce the delay time of vehicles and improve the communication efficiency of vehicles at intersections.

There are about 20 kinds of common traffic light control systems. They include OPAC (Optimized Policies for Adaptive Control) proposed by Nathan Gartner [13], [14], SCOOT (Split Cycle Offset Optimization Technique) proposed by the Transport Research Institute of the United Kingdom [15], SCATS (Sydney Coordinated Adaptive Traffic System) proposed by Road Traffic Bureau of New South Wales, Australia [16] and RHODESReal-time, Hierarchical, Optimized, Distributed and Effective System developed by the University of Arizona [17].

The key of intersection traffic light control is the traffic light timing calculation. Gtlich S. proposed the use of the traffic flow conservation law to calculate the optimal traffic light timing plan at the intersection [18]. Younes M. B. designed an efficient dynamic traffic light timing algorithm that adjusts the optimal green time for each traffic flow based on the real-time traffic situation around the intersection. The algorithm also takes into account the presence of emergency vehicles so that they can quickly pass through the intersections [19].

In this paper, based on adaptive control, a single intersection traffic light timing model is studied and a set of timing plan calculation algorithm is proposed. Simulation experiments show that the proposed timing plan can improve the traffic efficiency at the intersection by about 15% to 20%. The rest of the paper is organized as follows: Section II describes the intersection model and some basic assumptions of this paper; Section III presents a set of timing strategies with indefinite

cycle and indefinite phases; Simulation experiments and results analysis are given in Section IV. Finally, the paper is concluded in Section V.

II. ASSUMPTIONS AND DEFINITIONS

The intersection is the point of convergence between people and vehicles. People and vehicles from different directions arrive at their destinations through intersections, which determines the importance and complexity of traffic control at intersections in traffic systems. A typical intersection model is shown in Fig. 1. To simplify the discussion, we assume that

- 1) The right turn of the vehicle is not restricted by the traffic light and it is immediately released.
- 2) Do not consider pedestrians' influence on traffic.
- 3) The lane $L_i (1 \leq i \leq 8)$ can have different travel sequences and green light time during an intersection signal cycle.
- 4) At most two lanes of vehicles are allowed to travel at the same time.

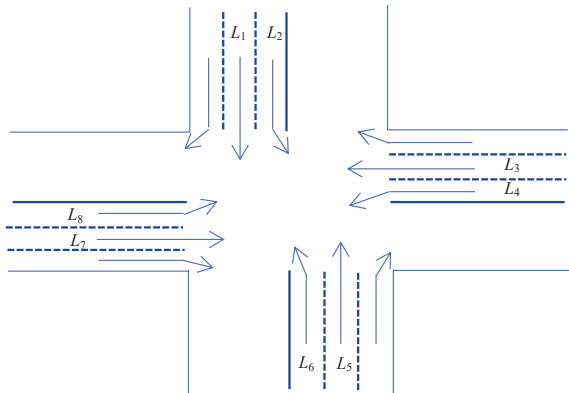


Fig. 1. Intersection lane layout.

The combination of any two lanes that do not have collisions (that is, no cross traffic) is called a phase, and the 8 lanes in Fig. 1 can make up to 12 phases. For example, L_1 and L_2 may form one phase, and L_1 and L_5 may also form one phase.

For the convenience of description in this paper, we give the following definitions:

Definition 1. As for any lane L_i , when lane $L_j (j \neq i)$ does not affect the traffic of lane L_i , L_j is called L_i 's compatible lane. Each lane has multiple compatible lanes. For example, lane L_6 have compatible lanes: L_2, L_3, L_5 .

Definition 2. In two lanes of a phase, the lane first determined the right to travel is called the main lane.

Definition 3. The traffic rate refers to the number of vehicles passing through the unit time. Obviously, the more vehicles that pass during the unit time, the higher the traffic rate is. The traffic rate for each lane is denoted by $q_i (i = 1, 2, \dots, 8)$, and the traffic rate of the intersection is denoted by q_0 .

Definition 4. In all lanes, the maximum number of vehicles that pass through the intersection without stagnating during the unit time is called intersection saturation flow q [20].

Obviously, saturation flow is greater than or equal to the traffic rate. When the traffic rate reaches saturation flow, the traffic rate is the highest. In this paper, we assume that the saturation flow of the lane is the same as the saturation flow of the intersection.

In practical applications, the time spent by different types of vehicles passing through the intersection is not the same. For ease of calculation and comparison, we use Table I to convert all types of vehicles into standard vehicle. The vehicles described later in this article refers to the standard vehicle.

TABLE I. VEHICLE CLASSIFICATION AND VEHICLE CONVERSION FACTOR

Vehicle Type	vehicle conversion factor	Load and power	Explanation
Small trucks	1.0	Load quality ≤ 2 tons	-
Medium truck	1.5	2 tons < load quality ≤ 7 tons	Including cranes
Large trucks	2.0	7 tons < load quality ≤ 14 tons	-
Heavy-duty truck	3.0	Load quality > 14 tons	-
Trailer	3.0	-	Including semi-trailer, flatbed trailer
Container car	3.0	-	-
small-sized buses	1.0	Rated seats ≤ 19	-
large-sized buses	1.5	Rated seats > 19	-
motorcycle	0.4-0.6	-	-
tractor	4.0	-	-

III. TIMING STRATEGY

A. Cycle and Traffic Light Control

Intersection traffic light cycle (hereinafter referred as cycle) is an important parameter of the timing strategy. Short cycles can easily cause the traffic lanes to alternate too frequently, affecting traffic flow through the intersections and excessively long cycles lead to increased waiting time for vehicles, causing a backlog of vehicles. Therefore, choosing an appropriate cycle is important for improving traffic efficiency.

Due to the uncertainty of the vehicle arrival time at the intersection, we use the traffic rate as the goal to calculate the green light time of each lane and cycle. The lower limit of the green light time is T_g , and the upper limit of the cycle is T_{max} . Thus, the signal control strategy is:

- 1) In each cycle, each lane has a green light right to travel;
- 2) The green light time of each lane is not less than T_g , and the traffic sequence is not fixed;
- 3) Each cycle time is not fixed, but not greater than T_{max} ;
- 4) Only vehicles of two compatible lanes are allowed to pass at the same time.

In order to implement the above strategy and facilitate control, the vehicle traffic data collected in the i -th cycle will be calculated in the $(i+1)$ th cycle to form a timing plan, which is implemented in the $(i+2)$ th cycle. This means that

the vehicle arrival data of the $(i+2)$ th cycle is predicted by the traffic data of the i -th cycle, and the efficiency of the $(i+2)$ th cycle is improved by adjusting the timing plan.

B. Timing Calculation

The timing calculation is to give the green light time and traffic sequence of each lane in a cycle.

1) *Green time optimization:* Assume that the green light time of lane L_j ($j = 1, 2, \dots, 8$) in the i -th cycle is T_j^i and the number of vehicles passing through the intersection is S_j^i . The traffic rate of lane L_j is $q_j^i = \frac{S_j^i}{T_j^i}$.

If the number of vehicle arrivals in the $(i+2)$ th cycle is the same as the i -th cycle, the traffic rate q_j^{i+2} can be increased by calculating the green light time T_j^{i+2} of the $(i+2)$ th cycle in (1).

$$T_j^{i+2} = \begin{cases} \frac{S_j^i}{\mu q}, & q_j^i \geq \mu q \\ \frac{S_j^i}{\lambda q}, & q_j^i \leq \lambda q \\ T_j^i, & \lambda q < q_j^i < \mu q \end{cases} \quad (1)$$

In (1), to make the green light time robust, the parameters μ ($0 < \mu \leq 1$) and λ ($0 < \lambda < 1, \lambda < \mu$) are set to limit the adjustment of the green light time:

- when $\lambda q < q_j^i < \mu q$, set $T_j^{i+2} = T_j^i$. That is, do not adjust the green light time.
- when $q_j^i \geq \mu q$, it means that there are too many vehicles in the lane L_j during i th cycle and the green light time is relatively short, and it is necessary to increase the green light time to allow more vehicles to pass.
- when $q_j^i \leq \lambda q$, it means that there are few vehicles in the lane L_j during i th cycle and the green time is longer. And the green light time needs to be shortened to increase the traffic rate of the lane. In order to ensure that each lane has the right to travel during the cycle, the minimum green light time is T_g .

To facilitate the operation of the control system, the green light time after the optimization of (1) is adjusted to a multiple of 5. This gives the predicted green light time T_j^{i+2} ($1 \leq j \leq 8$) for the $(i+2)$ th cycle.

2) *Lane combination calculation:* The green light time obtained in the previous section is considered from the perspective of the traffic of each lane. Vehicle traffic at intersections also needs to consider the cycle and the traffic sequence in order to obtain a timing plan that can be implemented.

In the timing plan, the traffic sequence at the intersection is called the lane combination. A lane combination consists of two lane sequences, as in (2).

$$\Pi = \begin{cases} L_{i1}L_{i2} \cdots L_{ik} \\ L_{j1}L_{j2} \cdots L_{jk} \end{cases} \quad (2)$$

In Π , $L_{ip}, L_{jq} \in \{L_1, L_2, \dots, L_8\}$, L_{ip} and L_{jq} are mutually compatible lanes, which belong to lane sequence one and lane sequence two, respectively. And they all travel at the

same time. Remark $R(L_i)$ is the remaining green time of the lane L_i . The calculation algorithm of the lane combination is as follows:

- If every lane has obtained the right to travel, then stop.
- Randomly select a lane L_m that has not obtained the right of travel as the main lane and give it the right to travel.
- If there are L_m compatible lanes that do not obtain the right of travel, then randomly select a lane L_i to give it the right to travel with the main lane L_m at the same time. Therefore L_i and L_m travel simultaneously and form a lane combination segment (L_m, L_i) , and then go to step(d); otherwise, L_m travel alone and form a lane combination (L_m) , then go to step(a).
- If $R(L_m) = R(L_i)$, then go to step(a).
- If $R(L_m) > R(L_i)$, let $R(L_m) \leftarrow R(L_m) - R(L_i)$ and L_m as main lane; otherwise, let $R(L_i) \leftarrow R(L_i) - R(L_m)$ and L_i as main lane. Then go to step(c).

Obviously, in the above algorithm, different choices of L_m and L_i will result in different lane combinations. The time spent traveling in the traffic sequence is called the cycle of lane combination, and the smallest cycle of the lane combination recorded as T_{min} . When $T_{min} \leq T_{max}$, all lane combinations with a T_{min} cycle are called candidate lane combinations; otherwise, the green time of the lane is compressed as described below until $T_{min} \leq T_{max}$.

- Randomly select a lane combination with a cycle T_{min} ;
- Calculate the difference between T_{min} and T_{max} , then we can get the current time needed to compress:

$$\Delta T = T_{min} - T_{max} \quad (3)$$

- In Π , the required compression time for each lane is according to the ratio of the initial traffic rate of the lane. That is:

$$T_i = \Delta T * \frac{q_i}{q_1 + q_2 + q_3 + \cdots + q_8} \quad (4)$$

3) *Lane traffic plan:* In the candidate lane combinations calculated in 2) of the part B above, there may be multiple kinds of lane combinations. There may be 4 pairs, 3 pairs, 2 pairs, and 1 pair that the green time of main lanes and compatible lane ends at the same time. In the actual traffic light control at intersection, the green light time of the two lanes in the same phase, ends at the same time. Therefore, in the candidate lane combination, the lane combination having the most pairs that the green time of main lanes and compatible lane ends at the same times is selected as the implemented lane traffic plan.

For example, if the optimized green light time of lanes L_1 to L_8 is 35s, 15s, 35s, 15s, 30s, 15s, 20s and 10s respectively, we can get two lane combinations:

$$\Pi_1 = \begin{cases} L_5L_2L_2L_8L_6L_4 \\ L_1L_1L_7L_7L_3L_3 \end{cases}$$

$$\Pi_2 = \begin{cases} L_2L_5L_5L_3L_3L_3L_8 \\ L_1L_1L_6L_6L_4L_7L_7 \end{cases}$$

In Π_1 there are two pairs of compatible lanes (L_8 and L_7 , L_4 and L_3) that green time ends at the same time. And in Π_2 ,

there is only one pair of compatible lanes (L_8 and L_7) that green time ends at the same time. So we choose Π_1 as the implemented lane traffic plan.

IV. EXPERIMENTS AND ANALYSIS

In this section, simulations of the timing strategies proposed in this paper are carried out by using the vehicle's simulation data and measured data at intersections.

A. Experimental Environment

The environment for the simulation experiments in this paper is as follows:

- (a) hardware configuration:
 - (i) processor: Inter(R) Core(TM) i5-2400 CPU @3.10GHZ
 - (ii) RAM: 8.00G
- (b) related Software
 - (i) operating system: Windows10
 - (ii) software tools: Microsoft matlab2014
 - (iii) programming language: c++

B. Experimental Data

The simulation data of this paper is divided into two categories: the first group of vehicle arrival data is randomly generated, and the second group of vehicle arrival data is real data we measured.

- (a) Random data
In each cycle, suppose that the number of arriving vehicles S_i in lane L_i is satisfied with the inequality (5).

$$0 \leq S_i \leq q * T_i \quad (5)$$

where T_i is the green time of lane L_i . This assumption guarantees that there is no stranded vehicle in each lane at the end of the cycle.

- (b) Measured data
We use the real data as vehicle arrival data. Assume that the vehicle arrives randomly at an intersection and obeys a Poisson distribution (6). The arrival rate is the average arrival rate in the measured data sample.

$$P(k) = \frac{(\rho k)_t e^{-(\rho t)}}{k!} = \frac{m^k e^{-m}}{k!} \quad (6)$$

Where ρ denotes the average vehicle arrival rate (vehicles/second) and $m = \rho t$ denotes the average number of vehicles arriving within time t .

In order to facilitate calculation and comparison, the actual measured data is converted into standard vehicles according to Table I, as shown in Table II. The vehicle arrival rate is shown in Table III.

TABLE II. THE TOTAL NUMBER OF STANDARD VEHICLES AT INTERSECTIONS

Time	L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8
7:30-7:45	87	56.5	306.5	144	85.5	63.5	379.5	73
7:45-8:00	114	44.5	332	129.5	100.5	83	313	57
8:00-8:15	86	41.5	343	101.5	77	48.5	284.5	50
8:15-8:30	79.5	36.5	304	102.5	54	46.5	276.5	60

TABLE III. MEASURED VEHICLE ARRIVAL RATE

Time	L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8
7:30-7:45	0.097	0.063	0.341	0.160	0.095	0.069	0.422	0.081
7:45-8:00	0.127	0.049	0.369	0.214	0.112	0.092	0.348	0.063
8:00-8:15	0.096	0.046	0.392	0.113	0.086	0.054	0.316	0.056
8:15-8:30	0.088	0.041	0.338	0.114	0.060	0.052	0.307	0.067

C. Experimental Results

We compared the timing plan proposed in this paper (abbreviated as ITP, Improved timing plan) and the fixed time plan (FTP) through the simulation experiments. Experimental results are evaluated using the following indicators:

- (a) Number of traffic vehicles: The sum of the number of traffic vehicles in each lane during a cycle.
- (b) Traffic efficiency: The ratio of number of traffic vehicles to cycle.
- (c) Number of detention vehicles: The sum of the number of vehicles staying in each lane after the end of the cycle.

1) *Random data simulation results analysis:* Simulation experiments conducted 50 cycles. In the first cycle, it is assumed that the green time of lanes L_1 to L_8 are: 30s, 15s, 30s, 15s, 30s, 15s, 30s, 15s. The number of arrival vehicles is calculated according to (5), and the lane combination is:

$$\begin{cases} L_1 L_2 L_3 L_4 \\ L_5 L_6 L_7 L_8 \end{cases}$$

From second to the 50th cycle, the ITP is the timing plan described in Section III. The number of arrival vehicles in each lane is still calculated according to (5). The FTP still adopts the first cycle timing plan, and the number of arrival vehicles in each lane is the same as the ITP timing plan. Other assumptions for the ITP: $T_{max}=90$ seconds, $q=0.5$ (vehicles/second), $\lambda=0.4$, $\mu=0.7$, $T_g=10$ seconds. The simulation results are shown in Fig. 2 and 3.

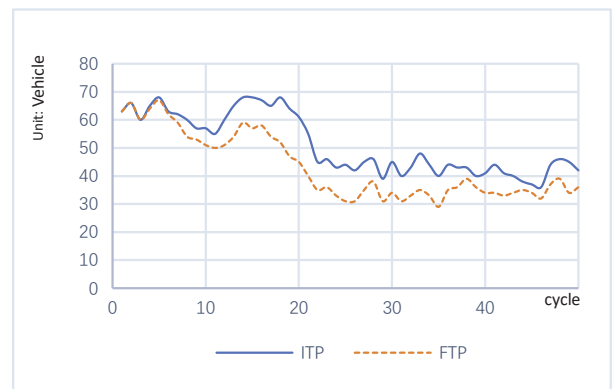


Fig. 2. Number of traffic vehicles comparison.

From Fig. 2, we can see that after several cycle adjustments, the number of traffic vehicles in the timing plan proposed in this paper is significantly increased compared to the number of vehicles in fixed timing plan. And the average number of the traffic vehicles increased by 20.3%. In the adjustment process, there are some cycles time in ITP are

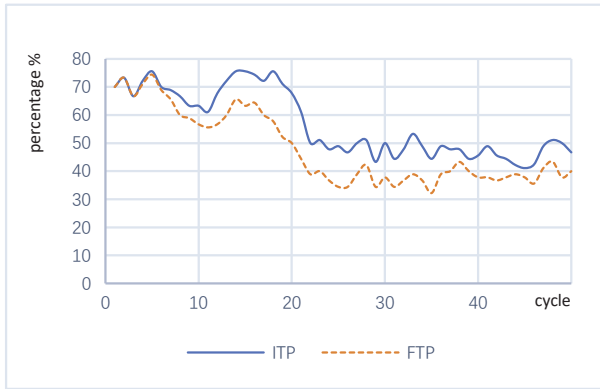


Fig. 3. Traffic efficiency comparison.



Fig. 5. Traffic efficiency comparison.

smaller than T_{max} . This will make the cycle change faster and make vehicles wait less time, then traffic efficiency is improved. From Fig. 3, traffic efficiency increased by an average of 19.9%. This shows that the timing strategy proposed in this paper significantly improves the traffic efficiency.

2) *Measured data simulation results analysis:* The simulation experiment is conducted for 1 hour. In the first cycle, it is assumed that the green light time of lanes L_1 to L_8 are: 30s, 15s, 30s, 15s, 30s, 15s, 30s, 15s. The number of arrival vehicles is the measured data in Table III and the lane combination is:

$$\begin{cases} L_1 L_2 L_3 L_4 \\ L_5 L_6 L_7 L_8 \end{cases}$$

From the second to the 50th cycle, the number of arrival vehicles in each lane is obtained from Table III. The FTP still adopts the first cycle timing plan, and the number of arrival vehicles in each lane is the same as the ITP. Parameter settings for the ITP: $T_{max} = 90$ seconds, $q = 0.5$ (vehicles/second), $\lambda = 0.4$, $\mu = 0.7$, $T_g = 10$ seconds. The simulation results are shown in Fig. 4, 5 and 6.

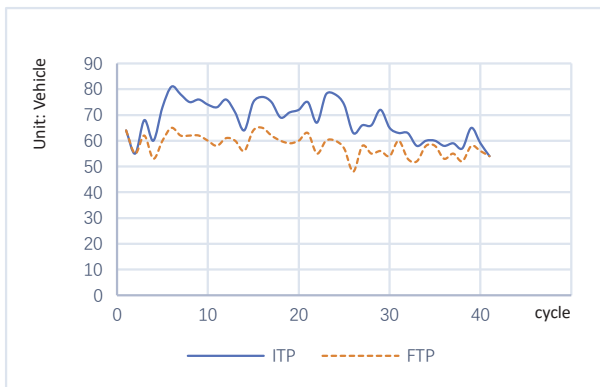


Fig. 4. Number of traffic vehicles comparison.

From Fig. 4, after several cycle adjustments, the number of traffic vehicles in ITP is higher than the number of traffic vehicles in FTP. And the average number of traffic vehicles

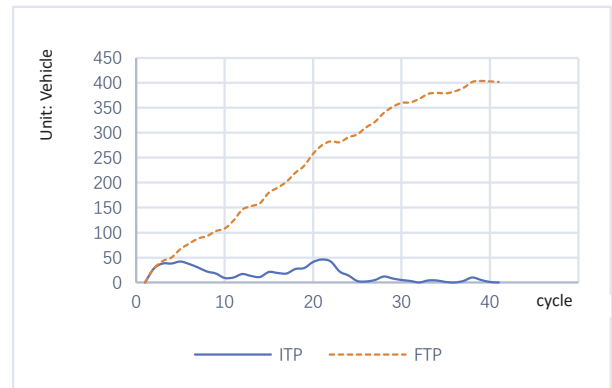


Fig. 6. Comparison of the number of retained vehicles.

increased by 16.6%. At the same time, from Fig. 5, the traffic efficiency in one hour increased by an average of 17.2%. It can also be seen from the Table III that the traffic rate of L_3 and L_7 are relatively large. When ITP is adopted, the green light time of large traffic rate lane will be appropriately extended, and the green light time of small traffic rate lane will be reduced to a sufficient extent. As shown in Fig. 6, the gap between the number of retained vehicles in ITP and the number of retained vehicles in FTP is increasing.

In summary, after the adjustment of ITP for about 5 cycles, its superiority gradually emerged. The ITP is significantly better than the FTP in terms of the number of passing vehicles, the efficiency of traffic, and the number of retained vehicles.

V. CONCLUSION

We proposes an adaptive timing strategy IPT with indefinite cycle and indefinite phase. On the one hand, the strategy can automatically adjust the green light time of each lane according to the traffic flow collected in real time. So that the green light time of lanes with less vehicle or no car is shortened, and the green light time of lanes with more cars is longer. On the other hand, the strategy calculates the appropriate lane combination through the concept of main lanes, compatible lanes, etc. to facilitate traffic light control at intersections. The

ITP improves the traffic efficiency of vehicles at intersections by increasing the utilization of the green light. The simulation experiment results show that the ITP increases the traffic efficiency compared to the traditional fixed timing plan by an average of 15% to 20%. The use of the ITP timing strategy can effectively shorten the waiting time of the vehicle, reduce the vehicle stagnation and improve the traffic condition of the intersection.

For further work, firstly, when the optimized cycle exceeds the maximum cycle, we can study the compression method so that the cycle satisfies both the cycle requirement and the maximum traffic rate. Secondly, in the lane combination calculation, we can consider the green time as the information of selecting the main lane, so that the implemented lane combination can be calculated more quickly.

REFERENCES

- [1] Gang Xiong, Fenghua Zhu, Xiwei Liu, Xisong Dong, Wuling Huang, Songhang Chen, Kai Zhao. "Cyber-physical-social System in Intelligent Transportation". IEEE/CAA Journal of Automatica Sinica, 2015, 2(03):320-333.
- [2] Li D, Deng L, Cai Z, et al. "Intelligent Transportation System in Macao based on Deep Self Coding Learning[J]. IEEE Transactions on Industrial Informatics, 2018", doi: 10.1109/TII.2018.2810291.
- [3] Collotta M, Bello L L, Pau G. "A novel approach for dynamic traffic lights management based on Wireless Sensor Networks and multiple fuzzy logic controllers". Expert Systems with Applications, 2015, 42(13):5403-5415.
- [4] Li C, Anavatti S G, Ray T. "Analytical Hierarchy Process Using Fuzzy Inference Technique for Real-Time Route Guidance System". IEEE Transactions on Intelligent Transportation Systems, 2014, 15(1):84-93.
- [5] Darwish T S J, Bakar K A. "Fog Based Intelligent Transportation Big Data Analytics in The Internet of Vehicles Environment: Motivations, Architecture, Challenges and Critical Issues". IEEE Access, 2018, doi: 10.1109/ACCESS.2018.281598
- [6] Ibrahim M, Riad M, El-Abd M. "RoadEye The Intelligent Transportation System". 2017 International Conference on Computer Systems and Applications (AICCSA), Hammamet, Tunisia, IEEE, 2017:21-22.
- [7] Hassouneh Y, Tumar I and Abu-Issa A. "Case study about capturing uncertainty in adaptive intelligent transportation systems". 2017 International Conference on Engineering and Technology (ICET), Antalya, Turkey, 2017:1-2.
- [8] Cao Z, Jiang S, Zhang J, et al. "A Unified Framework for Vehicle Rerouting and Traffic Light Control to Reduce Traffic Congestion". IEEE Transactions on Intelligent Transportation Systems, 2017, 18(7):1958-1973.
- [9] Khelafa I, Ballouk A, Baghdad A, et al. "Development of control algorithm for urban traffic". 2017 International Conference on Electrical and Information Technologies (ICEIT) Zhuhai China, 2017:1-5.
- [10] Collotta M, Pau G, Scat G, et al. "A dynamic traffic light management system based on wireless sensor networks for the reduction of the red-light running phenomenon". Transport & Telecommunication Journal, 2014, 15(1):1-11.
- [11] F. Faheem, Z. Zuraidah, A. Kayani and A. K. Aminuddin. "Optimization of vehicle actuation and multiplan algorithms for urban traffic control systems". 2017 IEEE Conference on Systems, Process and Control (IC-SPC), Malacca, Malaysia, 2017: 59-64.
- [12] Denisova, L.A. Meshcheryakov, V.A. "Automatic parametric synthesis of a control system using the genetic algorithm", Automation and Remote Control, 2015, 76(1):149156.
- [13] Gartner N, Pooran F, Andrews C. "Optimized Policies for Adaptive Control Strategy in Real-Time Traffic Adaptive Control Systems: Implementation and Field Testing". Transportation Research Record Journal of the Transportation Research Board, 2002, 1811(1):148-156.
- [14] Gartner N, Pooran F and Andrews C. "Implementation of the OPAC adaptive control strategy in a traffic signal network". 2001 IEEE Intelligent Transportation Systems, Oakland, CA, 2001: 195-200.
- [15] Bretherton D, Bodger M, Baber N. "SCOOT - the future [urban traffic control]". 12th IEE International Conference on Road Transport Information and Control, 2004. RTIC 2004. London, UK, 2004:301-306.
- [16] Luk, J.Y.K.. "Two Traffic-responsive Area Traffic Control Methods: SCAT and SCOOT ". Traffic Engineering & Control, 1984, 251:14-22.
- [17] Mirchandani P, Wang F Y. "RHODES to intelligent transportation systems". IEEE Intelligent Systems, 2005, 20(1):10-15.
- [18] Gttlich S, Potschka A, Ziegler U. "Partial Outer Convexification for Traffic Light Optimization in Road Networks". Siam Journal on Scientific Computing, 2017, 39(1):B53-B75.
- [19] Younes M B, Boukerche A. "An efficient dynamic traffic light scheduling algorithm considering emergency vehicles for intelligent transportation systems". Wireless Networks, 2017 DOI: 10.1007/s11276-017-1482-5.
- [20] Yunfeng Gao, Xiaoguang Yang, Hua Hu. "The theoretical model of basic saturation flow rate based on the car-following model ". 2005 Cross-Strait Intelligent Transportation System Seminar and Tongzhou Traffic Forum Shanghai China 2005:3-8.

Search Manager: A Framework for Hybridizing Different Search Strategies

Yousef Abdi

Department of Computer Engineering
Payame Noor University, Tabriz, Iran
<http://orcid.org/0000-0002-8517-8769>

Yousef Seyfari

Department of Computer Science
University of Tabriz, Tabriz, Iran
<https://orcid.org/0000-0002-2393-8814>

Abstract—In the last decade, many of the metaheuristic search methods have been proposed for solving tough optimization problems. Each of these algorithms uses its own learn-by-example mechanism in terms of “movement strategy” to evolve the candidate solutions. In this paper, a framework, called Search Manager, is proposed for hybridizing different learn-by-example methods in one algorithm, which is inspired by the organizational management system in which managers change their management method by viewing performance reduction in their managerial organization. The proposed framework is verified using standard benchmark functions and real-world optimization problems. Further, it is compared with some well-known heuristic search methods. The obtained results indicate not only the optimization capability of the proposed framework, but also its ability to obtain accurate solutions and to achieve higher convergence precision.

Keywords—Global optimization; metaheuristic; organization management; hybridizing search methods

I. INTRODUCTION

In different areas of science such as industry, engineering, and management there are many complex problems, also known as optimization problems, such that there is no exact algorithm to solve them in polynomial time. Due to this fact, to find a relatively optimal solution for these kinds of problems, in the past decades, a significant number of different optimization algorithms have been introduced by researchers.

Optimization algorithms can be divided into two broad groups, deterministic and stochastic. Deterministic algorithms such as the Nelder–Mead search method [1], the tunnelling method [2], and renormalization group methods [3] perform based on gradients and second-order derivatives. A remarkable advantage of deterministic optimization methods is the fast convergence; however, for high dimensional and multimodal functions they may fall into a local optimum. In stochastic algorithms, although the quality of the obtained solution cannot be guaranteed, they are more efficient and flexible than deterministic approaches. Other advantages are the capability of escaping from a local optimum, good performance, and ease of implementation.

Stochastic optimization algorithms generally are population-based algorithms that begin with a set of randomly generated candidate solutions and by applying some specified rules iteratively, gradually evolve initial solutions. The rules are usually inspired by the behaviors of biological and physical

systems in nature, culture, society or politics. Based on the source of inspiration, they can be classified into three main groups: (1) Evolution-based, (2) Swarm-based, and (3) Human-based algorithms.

Evolution-based algorithms mimic natural biological evolution and selection. The Genetic Algorithm (GA) has been the most popular in this class of optimization algorithms. It uses the Darwinian theory of natural selection, crossover, and mutation [4]. Other algorithms that can be classified in this group are Evolution Strategy (ES) [5], Genetic Programming (GP) [6], Differential Evolution (DE) [7], Evolutionary Programming (EP) [8], and Biogeography-Based Optimizer (BBO) [9].

Swarm-based algorithms get their inspiration from the collaborative conduct of a group of animals, such as ant colonies, honey bees, and bird flocks. They are typically made up of a population of agents (swarm individuals) interacting together to fulfill the main goal of the system. To date, several swarm based optimization algorithms have been proposed in the literature. Particle Swarm Optimization (PSO) is the most popular and significant algorithm that mimics the behavior of a flock of migrating birds heading for an unknown destination [10]. So many variants of PSO algorithm such as MLPSO [11], FST-PSO [12], etc. with the aim of improving its performance in different types of problems have been presented so far. Other examples of this class of algorithms are: Ant Colony Optimization [13], Artificial Bee Colony [14], Bacterial Foraging [15], Cat Swarm Optimization [16], Elephant Herding Optimization (EHO) [17], Bat Algorithm (BA) [18], etc.

Human-based algorithms imitate human social behaviors. For example, Fireworks Algorithm (FA) inspired by observing fireworks explosion [19], Harmony Search algorithm (HS) inspired by harmony improvisation process of musicians [20]. Teaching Learning Based Algorithm (TLBO) imitates the interaction between a teacher and her students [21]. Some of the other popular algorithms in this group are Cultural Algorithm (CA) [22], Imperialist Competitive Algorithm (ICA) [23], Exchange Market Algorithm (EMA) [24], Soccer League Competition (SLC) [25], [26], and Brain Storm Optimization (BSO) [27], World Competitive Contests (WCC) [28].

There are some extra optimization algorithms that cannot be assigned to the above three groups. They imitate some additional rules that exist in the world and universe. For

instance, Optics Inspired Optimization (OIO) uses the law of reflection [29] or Ray Optimization (RO) was developed based on Snell's light refraction law [30]. Some of the other algorithms are: Water Wave Optimization (WWO) [31], Big-Bang Big-Crunch (BBBC) [32], Charged System Search (CSS) [33], Artificial Chemical Reaction Optimization Algorithm (ACROA) [34], Curved Space Optimization (CSO) [35], Central Force Optimization (CFO) [36]. However, we think this kind of optimization algorithms is more complex for understanding than the specified three groups of algorithms.

Regardless of the inspiration's source, still no universal optimization algorithm has emerged; nor does it appear likely that one ever will [37] and the No Free Lunch Theorem of Optimization (NFLT) supports this view [38]. Besides, although many optimization methods have been proposed so far, one or another method presents a better solution on a specific problem. Therefore, finding more efficient algorithms is still in progress and we will see further optimization algorithms with the development of human identity from nature. However, a powerful optimization algorithm can solve many problems.

Considering the aforementioned facts, in this paper we intend to propose a framework, Search Manager, for combining different search abilities of different optimization algorithms in one algorithm to effectively solve the optimization tasks. The proposed framework accomplishes this by imitating the changing management style that is applied by managers in real-world organizations. In any organization when managers cannot obtain better results, they have to change their management style to improve the performance of the organization.

The simplicity and capability of using variable movement strategies are powerful aspects of the proposed framework. Although in this study we use four simple movement strategies, the others can be proposed or obtained from the other optimization algorithms to get better performance for different optimization problems.

The remainder of this paper is organized as follows: Section II gives a general view of the main behavior of the proposed algorithm. Section III reviews some of the main principles about organizations and their management process. In Section IV, the proposed framework is introduced. Section V presents experimental investigation of the proposed method and its comparative study with other metaheuristic algorithms. Finally, some concluding remarks and future works are presented in Section VI.

II. MOTIVATION

Almost all of the previous population-based optimization methods iteratively evolve and optimize a population of individuals (candidate solutions) according to some criteria to reach a final population which has a near-optimal solution to a problem. In fact, they use learn-by-example mechanisms in terms of 'movements' to evolve the candidate solutions, but the main problem of these methods is the lack of alternative movement strategies for different situations.

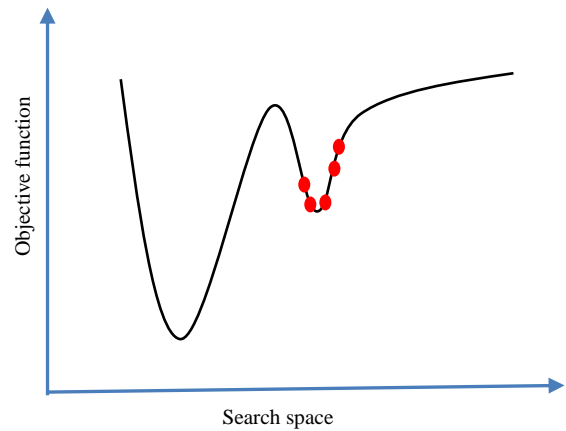


Fig. 1. Search space of an objective function.

The proposed framework overcomes the aforementioned limitation by combining different movement strategies. If the current strategy does not improve solutions within the population, another one will be selected. For example, consider the situation in Fig. 1 that represents the search space of an objective function, in which there exist a local and global minimum at inflection points. In this situation, the population of six individuals has got stuck in the local minimum. If there is not enough diversity in the population like this situation, an optimization algorithm would get trapped in the local optimums due to the lack of an alternative strategy. In the proposed framework, it is more likely that the population with not enough diversity could escape from local minimum by using different movement strategies.

III. ORGANIZATION MANAGEMENT BACKGROUND

Firm infrastructure, considered as organizations' structure and its management process, is a branch of social science, which includes a lot of relative concepts. However, in this section, we only review a few important concepts that are used in the proposed framework.

In the society around us, there are various organizations such as hospitals, schools, social institutions, etc. to achieve certain objectives in social life. In the words of Leavitt, an organization is "a particular pattern of structure, people, tasks and techniques" [39]. In the view of Katz and Kahn, an organization is "a system which is composed of a set of subsystems" [40], put the other way round, an organization is a unified system, whose functionality is divided into different key subsystems. These subsystems are all parts of an organization working together for a common purpose.

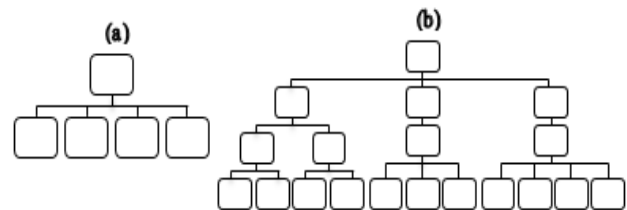


Fig. 2. (a) Flat management structure (b) Hierarchical management structure.

Each organization in the society needs a system of management. Management is defined in different ways by different people. It may be viewed as a process of leading and controlling the activities of organization employees [41], or as a decision-making process in which the manager is a decision-maker that leads the organization to achieve the objectives [42]. In all of the organizations, employees look up to managers for guidance, thus for making appropriate management decisions and guiding organization members, managers adopt a management style or a form of leadership, primarily by their abilities, personalities, and values [43]. However, half of the decisions made by managers within organizations fail [44]. Failure in decisions signals the need for new decisions or changing the management style that decisions were taken under it. Failed decisions can be a source of learning and they can help the manager to take effective and operational decisions in the future [45]. Management style is a form of leadership that specifies how employees should follow organization's policy. There are many types of management styles. Managers move in and out of these various styles as the need arises.

Management structure is another aspect of the management process. It is the manner in which the management of a company or organization is organized. It determines the scope and nature of how leadership is disseminated throughout the organization [42], [43]. Organizations commonly adapt either a flat or hierarchical structure. In the flat structure, there are a few or no levels of middle management between top managers and employees. But in the hierarchical structure, there are a number of hierarchical levels between top managers and employees. Fig. 2 shows these structures. The proposed framework uses the aforementioned concepts for combining different learn-by-example methods in one algorithm for solving optimization problems.

IV. PROPOSED METHOD

In this section, according to the basic principles of organization management outlined in the previous section, a framework is proposed for combining different search methods. The core of the Search Manager is based on a basic principle in the management of today's organizations: "managers need to adopt a new management style when organization's performance goal gets worse". The framework consists of three steps as follows:

- Step 1: Initialization
- Step 2: Movement
- Step 3: Repeat Step 2 until the stop criterion is satisfied.

Fig. 4 shows the flowchart of the Search Manager and details of its steps are represented in the following subsections.

A. Step 1: Initialization

The goal of the optimization task is to select n decision variables x_1, x_2, \dots, x_n (known as candidate solution) from a feasible region in such a way as to optimize a given objective function. The values of the decision variables are represented as floating point numbers and the cost of a candidate solution is obtained by evaluating the objective function f at these variables like the following equation.

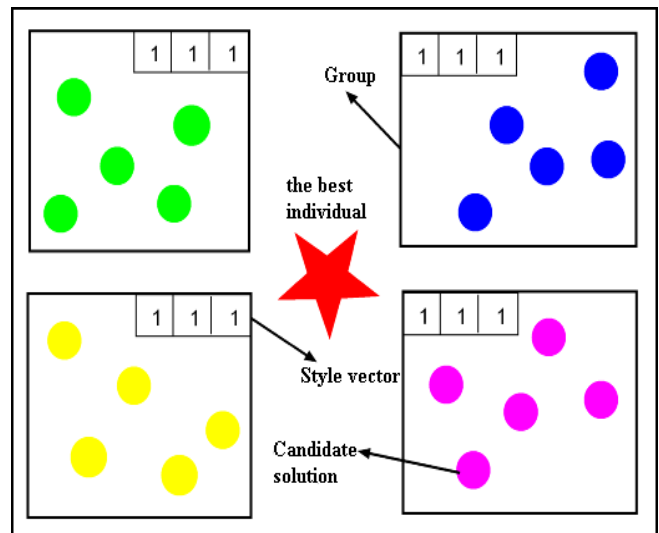


Fig. 3. A sample of population structure in Search Manager.

$$Cost = f(\text{Candidate solution}) = f(x_1, x_2, \dots, x_n)$$

To start the optimization algorithm, a population of size N_{pop} is generated, and then the best individual is selected. The remainder of solutions is divided among m distinct groups. For example, Fig. 3 shows four distinct groups of population, in which the best individual from the population is specified.

Each group will have current and previous average cost. The average cost of a group is defined as

$$GC_n = \text{mean}\{c_1, c_2, \dots, c_i\}$$

Where GC_n is the group cost of the n th group of candidate solutions and c_i is the cost of the i th solution. Furthermore, each group has an t -dimensional style vector, where t is the number of movement styles. The initial values of these vector components are set to 1 and each component holds the score value for the specified movement style. More explanations on this vector are presented in the following subsection.

B. Step 2: Movement

In the movement step, candidate solutions are evolved for the next generation using movement styles. The framework uses four simple movement styles inspired from other optimization algorithms and their mathematical formulas are represented in the following subsections.

As mentioned in subsection 4-A, each group of the candidate solutions has a vector that holds score values for different movement styles. For each group, a style (movement method) is selected by using the values of style vector and the roulette wheel selection method. The selection probability is calculated as (1).

$$P_i = \frac{\text{ScoreValue}_i}{\sum_{j=1}^t \text{ScoreValue}_j} \quad (1)$$

Where, P_i is the probability of being selected i th style from the t styles. It is obvious that the style with a higher score is more probable for selection than the others.

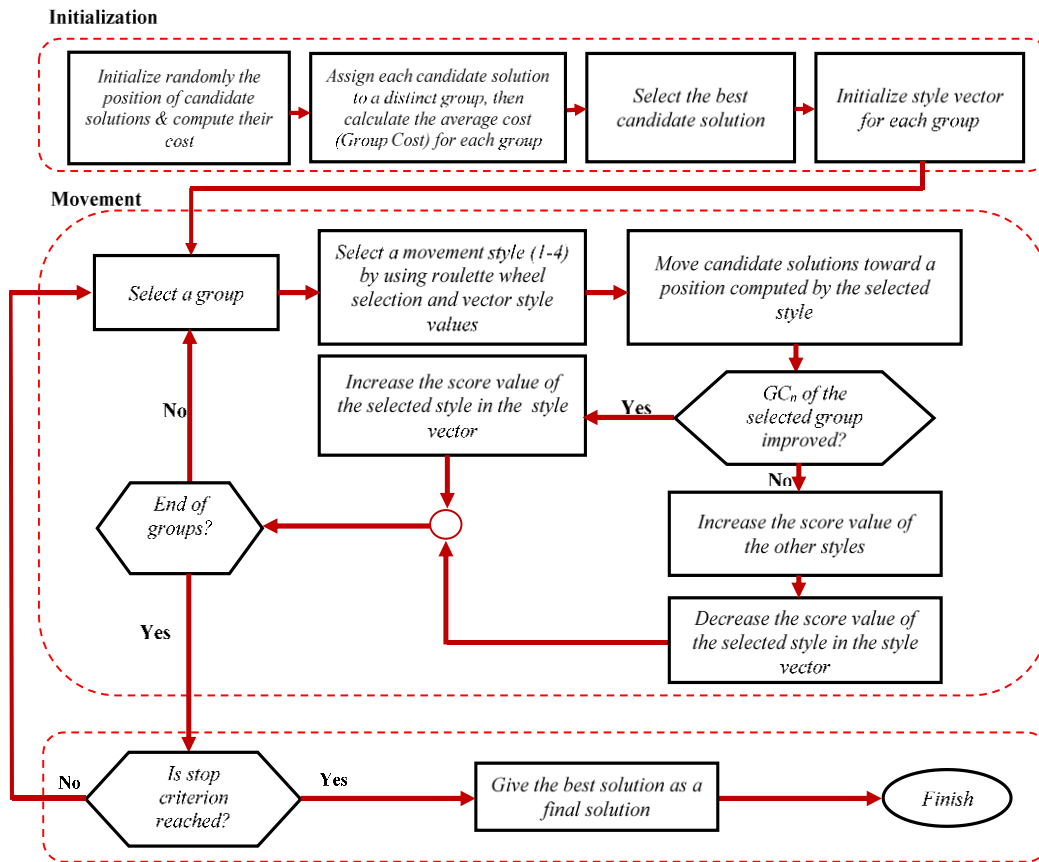


Fig. 4. The flowchart of the proposed framework.

At the beginning of the movement step, for each group one of the styles is selected randomly, then all candidate solutions in the group move toward a position calculated by using the selected movement style, and then finally, the group cost (GC_n) is computed and compared with the previous group cost, if an improvement is seen, the selected style is rewarded by increasing its score value in the style vector. Otherwise, the selected style is penalized by dividing its score value by a random number generated between 2 and the number of population. At the same time, the score values of the other styles are increased by a random number between 0 and 1. For example, Fig. 5 shows a sample of style vector for a group. *Style3* is more probable to be selected for the next iteration. If *Style3* is selected and the group cost is improved, the values of the vector will be changed like Fig. 6. It is observable that the score value of *Style3* has been increased.

Style 1	Style 2	Style 3
8	5	25

Fig. 5. A sample of style vector for holding scores of movement styles.

Style 1	Style 2	Style 3
8	5	25.3

Fig. 6. A sample of style vector after punishing and rewarding styles.

Style 1	Style 2	Style 3

8.5	5.5	2.5
-----	-----	-----

Fig. 7. A sample of style vector after punishing and rewarding styles.

Now, based on the values in Fig. 5, if *Style3* does not improve the group cost or make it worse, the values of the vector style will change like Fig. 7. It can be seen that *Style3* has been decreased and the others have been increased.

The used movement styles are inspired from other optimization algorithms and their equations are formulated as follows:

1) *The First Movement Style*

By the first movement style, solutions accept some of the variable values from the best individual. This style is applied by using (2).

$$\Phi = \text{randomly selected decision variables} \quad (2)$$

$$X_i(\Phi) = X_{\text{the best individual}}(\Phi)$$

Where, X_i is a selected solution from a group, and Φ is a set of randomly chosen decision variables.

2) *The Second Movement Style*

In the second movement style, solutions in the group are updated as

$$\begin{aligned} T &= X_{\text{the best individual}} \\ T &= T + \{\text{rand}\} \times \{T - X_r\} \\ X_i^{\text{new}} &= X_i^{\text{old}} + \{\text{rand}\} \times \{T - X_i^{\text{old}}\} \end{aligned} \quad (3)$$

Where, {rand} is an n-dimensional random vector and its values are in [0, 1]. X_r corresponds to the randomly selected solution from the group and X_i^{old} corresponds to the i th solution in the group.

3) The Third Movement Style

In the third movement style, solutions in the group are updated as

$$\begin{aligned} T &= X_{\text{the best individual}} \\ T &= T + \{\text{rand}\} \times (T - V) \\ X_i^{\text{new}} &= X_i^{\text{old}} + \{\text{rand}\} \times (T - X_i^{\text{old}}) \end{aligned} \quad (4)$$

Where, V is the variance of solutions in the group.

4) The Fourth Movement Style

The fourth movement style uses (5) for updating solutions in the group. Some of the selected variables of the best individual are changed in a temporary vector and then the other solutions in the group move towards this vector.

$$\begin{aligned} \Phi &= \text{randomly selected decision variables} \\ T &= X_{\text{the best individual}} \\ T(\Phi) &= \{\text{rand}\} \\ X_i^{\text{new}} &= X_i^{\text{old}} + \{\text{rand}\} \times (T - X_i^{\text{old}}) \end{aligned} \quad (5)$$

Where, {rand} is a generated random vector from the search space.

After applying movement operation, the best individual is selected from the population.

V. COMPARATIVE STUDY

In this section, the ability of the Search Manager is assessed from two perspectives: one is by applying it to optimization of 44 benchmark functions; another is through 15 real-world problems.

In the first perspective, Search Manager is applied on a wide range of nonlinear benchmark functions, 14 functions from the CEC 2005¹ [46] and all of the 30 functions from the CEC 2014 [47]. These benchmark suites include a diverse set of problem features such as unimodality, multimodality, separability, non-separability, rotation, scalability, etc. for single objective optimization. They are based on classical benchmark functions, such as Rosenbrock's, Rastrigin's, Swefel's, Griewank's, and Ackley's function. It should be noted although these test functions are organized by the community of heuristic algorithms in the framework of a workshop, in this study they are used as standard test functions for comparison. The total functions in CEC 2005 and CEC 2014 can be divided into some groups according to their

characteristics. Table I shows these groups. The other properties of these functions have been defined properly in their corresponding references, and hence they are not repeated here. The experimental results on these benchmark functions may reveal how Search Manager performs on various functions as well as can be compared to those obtained by other algorithms.

In the second perspective, Search Manager is applied to 15 real-world problems, which are derived from CEC 2011 [48].

To evaluate the performance of the proposed method, obtained results using Search Manager is compared with the result of some well-known and recently proposed social and nature-inspired optimization algorithms from the computational viewpoint. The selected algorithms are listed below.

- **CA:** Cultural Algorithm models social evolution and learning in agent-based societies [22].
- **ABC:** Artificial Bee Colony inspired by foraging behavior of honey bee swarm [49].
- **GSA:** Gravitational Search Algorithm uses the law of gravity and the notion of mass interactions based on the laws of gravity and motion [50].
- **FOA:** Forest Optimization Algorithm inspired by few trees in the forests which can survive for several decades, while other trees could live for a limited period [51].
- **WOA:** Whale Optimization Algorithm mimics the social behavior of humpback whales. The algorithm is inspired by the bubble-net hunting [52].
- **DE:** Differential Evolution is a type of standard genetic algorithm [7].

Computation code of all above-mentioned algorithms is taken from web pages dedicated to these algorithms. The MATLAB code for all algorithms is available at: <https://data.mendeley.com/datasets/f5jvxbw8xb/1> (DOI: **10.17632/f5jvxbw8xb.1**)

A. Experimental Settings

The experimental environment is a computer of Intel Core i3, 4GB DDR2 memory, and Windows 7 operating system. Tests are performed in 10-, 30-, and 50-dimensional for each test function in the first perspective. All algorithms are executed 30 times for each problem with a total of $D \times 10^4$ and 5×10^4 evaluations of the objective function (Max_FES) in the first and second perspective respectively, where D is the dimension of the problem, and the obtained average results are compared with the other algorithms. The population size $P = 50$ is used for all algorithms. Table II shows the other recommended parameter settings for each algorithm. The values of these parameters are selected based on the recommendation from their original papers or previous related works. Additionally, initial candidate solutions are randomly calculated by uniform distribution between lower and upper limits of benchmark functions.

¹ The first 14 functions are selected from CEC 2005 benchmark set.

TABLE I. FUNCTION TYPES

Test Suite	Unimodal	Simple Multimodal	Hybrid (Multimodal)	Composition (Multimodal)	Expanded (Multimodal)	Hybrid Composition (Multimodal)
CEC 2005	$f_1 - f_5$	$f_6 - f_{12}$	-	-	$f_{13} - f_{14}$	Not used in this study
CEC 2014	$f_1 - f_3$	$f_4 - f_{16}$	$f_{17} - f_{22}$	$f_{23} - f_{30}$	-	-

TABLE II. OPTIMIZATION METHODS PARAMETERS

Optimization Method	Parameters
ABC	Limit = 50D [49]
GSA	$G_0=100, \alpha=20, K_0$ =population size [50]
FOA	Life time = 6, LSC = 2, Area limit = 30, Transfer rate = 10%, GSC = 3 [51]
WOA	$a=2 - 0$ [52]
CA	Acceptance rate = 0.3 [53] Knowledge type = Normative & Situational
DE	$F = 0.9, CR = 0.1$ [54]
Search Manager	Without the need for setting any control parameters

It should be noted that we know better results can be achieved from these algorithms by fine-tuning their control parameters. However, finding the perfect parameter for each problem is expected to be a very time-consuming task. Therefore, fixed parameter settings are adopted for each algorithm and this condition is equal for all algorithms.

B. Results and Discussion on Perspective 1

Experimental results of all algorithms on 14 functions of CEC 2005 and 30 functions of CEC 2014 are presented in Tables III to V and Tables VI to VIII, respectively. In all tables, the mean and standard deviation of 30 runs of the algorithms for each function are adopted to assess the optimization performance of the proposed algorithm. Additionally, a statistical test called Wilcoxon rank-sum test [55], which is a nonparametric statistic test for the independent samples, is conducted on the experimental results at the 5% significance level to judge whether the obtained results from the proposed method are significantly different from the other algorithms and have not occurred by chance [56]. The cases are marked with “+ / \approx / -” when the performance of Search Manager is significantly better than, equal to, and worse than the other test algorithms. For clarity, the best results are marked in bold. Tables IX and X show obtained p -values in the statistical test between the proposed algorithm and each of the remaining algorithms over benchmark functions in 10 and 50 dimensions. Moreover, Fig. 8 and 9 represent the convergence rate of some selected unimodal and multimodal functions for all of the algorithms.

Based on the presented results of CEC 2005 test suite in Tables III to V, Search Manager provided better results over 10 dimensions. Although its performance decreases in 30 and 50 dimensions, it can be seen that its rank is better than the other algorithms overall dimensions. In addition, the proposed algorithm displayed the best performance over expanded functions in all dimensions.

For CEC 2014 test suite, the experimental results in Tables VI to VIII show that the Search Manager presented good results in terms of average rank, but the DE algorithm outperforms Search Manager in most of the functions over 10 and 30 dimensions. Furthermore, Search Manager could not

produce successful results on composition functions as well as unimodal functions as compared to the ABC and DE algorithms. In addition, the GSA showed the best results in 30 and 50 dimensions of hybrid functions. However, Search Manager may get better results by using other types of movement styles.

Based on the convergence curves (Fig. 8 and 9), it can be seen that Search Manager has a good convergence rate and it can reach an optimum solution with less calculation.

To summarize, although Search Manager outperforms most of the compared algorithms in terms of average rank in all benchmark function sets, its performance is lower than the DE and ABC in some complex functions. However, the powerful aspect of the Search Manager is its ability to accept different movement strategies, and hence it may achieve better results for these functions by using other combinations of learn-by-example methods as movement styles. In other words, instead of adjusting parameters to get better results that we have seen in previous optimization algorithms, we can adjust the combination of different learn-by-example methods in Search Manager.

C. Results and Discussion on Perspective 2

Herein, the performance of Search Manager is evaluated through 15 real-world optimization problems. Detailed definitions of these problems can be found in [48].

Considering the experimental results represented in Table XI, Search Manager outperforms all of the compared algorithms in most of the real-world problems. In more details, Search Manager outperforms ABC, FOA, GSA, WOA, CA, and DE on 9, 11, 11, 9, 10, 9, 7 problems respectively. On the contrary, ABC, FOA, GSA, WOA, CA, and DE are better than Search Manager on 2, 2, 2, 2, 2, 1, and 1 problems, respectively. All the compared algorithms obtained similar results in solving the T03BCB and T08TNEP. Therefore, based on the obtained results, it is observable that Search Manager performs better in most of the real-world problems when compared with ABC, FOA, GSA, WOA, CA, and DE.

The conclusion that can be drawn from experiments on real-world problems is that the combination of the used learn-by-example equations as movement styles in Section IV is most suitable for this kind of optimization problems.

VI. CONCLUSION AND FUTURE WORKS

In this paper, we have proposed a framework, called Search Manager that imitates an important part of the management process, which is concerned by managers in any organization. The heart of the Search Manager is applying another movement strategy when the average cost of candidate solutions gets worse. Search Manager has one main operator: Movement. The movement operator implements different learn-by-examples methods as movement styles. Extensive analysis is carried out to reveal the ability of Search Manager

to find the best fitness value over the search space not only in the benchmark functions, but also in the real optimization problems as well. Of course, it should be noted that like the other optimization algorithms, Search Manager is suitable for some sort of problems. However, its important property is that it provides a framework for combining different search methods, so it can be made compatible with any optimization problem by modifying its movement styles.

In the basic Search Manager, we have used four movement strategies as movement styles, but more or other types of movement methods can be formulated to improve its performance.

Search Manager can use the learn-by-example strategies of other optimization algorithms. Therefore, it is capable of combining these strategies in one algorithm and in utilizing their search abilities. From this capability, a possible line for future work would be investigating the behavior of the Search Manager with the combination of different learn-by-example strategies.

Other possible avenues of future research on the Search Manager include: using other mechanisms for rating movement styles because the performance of the proposed framework is highly dependent on this mechanism, applying a mechanism for moving solutions between groups, and finally, management science and its application in organizations have a lot of concepts in the real-world, such as how to fire and hire employees, employee promotion mechanisms, etc. that can be brought into the framework via simulation.

REFERENCES

- [1] A. Nelder, R. Mead, "A simplex method for function minimization," *Comput J.*, vol. 7, no. 4, pp. 308-31, 1965.
- [2] A. V. Levy, A. Montalvo, "The tunneling algorithm for the global minimization of functions," *Siam J Sci Stat Comp*, vol. 6, no. 1, pp. 15-22, 1985.
- [3] D. Shalloway, "Application of renormalization group to deterministic global minimization of molecular conformation energy functions," *J Glob Optim*, vol. 2, pp. 281-311, 1992.
- [4] J. H. Holland, "Genetic algorithms," *Sci Am*, Vol. 267, pp. 66-72, 1992.
- [5] I. Rechenberg, "Evolutions strategien," in: B. Schneider, U. Ranft (editors), *Simulationsmethoden in der Medizin und Biologie*, Berlin: Springer Heidelberg, pp. 83-114, 1978.
- [6] J. R. Koza, "Genetic programming: on the programming of computers by means of natural selection," Cambridge, MIT Press, 1992. ISBN: 978-0-262-11170-6.
- [7] R. Storn, K. Price, "Differential Evolution – A simple and efficient heuristic for global optimization over continuous spaces," *J. Glob Optim*, vol. 11, no. 4, pp. 341-59, 1997.
- [8] G. B. Fogel, D. B. Fogel, "Continuous evolutionary programming: Analysis and experiments," *Cybernet and Sys*, vol. 26, pp. 79-90, 1995.
- [9] D. Simon, "Biogeography-based optimization," *IEEE T Evolut Comput*, vol. 12, pp. 702-713, 2008.
- [10] J. Kennedy, R. C. Eberhart, "Particle swarm optimization," in: *IEEE International Conference on Neural Networks*, Perth, Australia, pp.1942-48, 1995.
- [11] P. Liu, J. Liu, "Multi-leader PSO (MLPSO): A new PSO variant for solving global optimization problems," *Appl Soft Comput*, vol. 61, pp. 256-263, 2017.
- [12] M. S. Nobile, P. Cazzaniga, D. Besozzi, R. Colombo, G. Mauri, G. Pasi, "Fuzzy Self-Tuning PSO: A settings-free algorithm for global optimization," *Swarm Evol Comput*, Inpress, 2017.
- [13] M. Dorigo, T. Stützle, "Ant colony optimization. Cambridge," MIT Press, 2004. ISBN: 978-0-262-04219-2.
- [14] D. Karaboga, "An idea based on honey bee swarm for numerical optimization," Erciyes University, Engineering Faculty, Computer Engineering Department, Technical Report-TR06, 2005.
- [15] K. M. Passino, "Biomimicry of bacterial foraging for distributed optimization and control," *IEEE Contr Sys Mag*, vol. 22, pp. 52-67, 2002.
- [16] S. C. Chu, P. W. Tsai, J. S. Pan, "Cat swarm optimization," in: *9th Pacific Rim International Conference on Artificial Intelligence*, Guilin, China, pp.854-58, 2006.
- [17] G. Wang, S. Deb, L. Coelho, "Elephant herding optimization," in: *3rd International Symposium on Computational and Business Intelligence*, Bali, Indonesia, 2015.
- [18] X. S. Yang, "A new metaheuristic bat-inspired algorithm," in: J. Gonzalez, D. Pelta, C. Cruz, G. Terrazas, N. Krasnogor (editors), *Nature Inspired Cooperative Strategies for Optimization*, vol. 284 of *Studies in Computational Intelligence*, Berlin, Springer Heidelberg, 2010.
- [19] Y. Tan, Y. Zhu, "Fireworks algorithm for optimization," in: Y. Tan, Y. Shi, K. C. Tan (editors), *Advances in swarm intelligence*, Berlin, Springer Heidelberg, pp. 355-64, 2010.
- [20] Z. W. Geem, J. H. Kim, G. Loganathan, "A new heuristic optimization algorithm: harmony search," *Simulation*, vol. 76, pp. 60-8, 2001.
- [21] R. V. Rao, V. J. Savsani, D. P. Vakharia, "Teaching-learning-based optimization: a novel method for constrained mechanical design optimization problems," *Comput Aided Design*, vol. 43, no. 3, pp. 303-5, 2015.
- [22] R. Reynold, "An introduction to cultural algorithms," in: *3rd Annual Conference on Evolutionary Programming*, World Scientific Publishing, pp.131-9, 1994.
- [23] E. Atashpaz-Gargari, C. Lucas, "Imperialist competitive algorithm: an algorithm for optimization inspired by imperialistic competition," in: *the 2007 IEEE congress on evolutionary computation (CEC)*, pp.4661-7, 2007.
- [24] N. Ghorbani, E. Babaei, "Exchange market algorithm," *App Soft Comp*, vol. 19, pp. 177-87, 2014.
- [25] N. Moosavian, B. K. Roodsari, "Soccer league competition algorithm: a new method for solving systems of nonlinear equations," *Int J Intell Sci*, vol. 4, no. 1, pp. 7-16, 2013.
- [26] N. Moosavian N, B. K. Roodsari, "Soccer league competition algorithm: a novel meta-heuristic algorithm for optimal design of water distribution networks," *Swarm Evol Comp*, vol. 17, pp. 14-24, 2014.
- [27] Y. Shi, "Brain Storm Optimization Algorithm," *Adv Swarm Intell*, vol. 6728, pp. 303-309, 2011.
- [28] Y. Masoudi-Sobhanzadeh, H. Motieghader, "World Competitive Contests (WCC) algorithm: A novel intelligent optimization algorithm for biological and non-biological problems," *Informatics in Medicine Unlocked*, vol. 3, pp. 15-28, 2016.
- [29] A. HusseinzadehKashan, "A new metaheuristic for optimization: Optics inspired optimization (OIO)," *Comput Oper Res*, vol. 55, pp. 99-125, 2015.
- [30] A. Kaveh, M. Khayatizad, "A new meta-heuristic method: ray optimization," *Comput Struct*, vol. 112, pp. 283-94, 2012.
- [31] Y. J. Zheng, "Water wave optimization: A new nature-inspired metaheuristic," *Comput Oper Res*, vol. 55, pp. 1-11, 2015.
- [32] O. K. Erol, I. Eksin, "A new optimization method: big bang-big crunch," *Adv Eng Softw*, vol. 37, no. 2, pp. 106-111, 2006.
- [33] A. Kaveh, S. Talatahari, "A novel heuristic optimization method: charged system search," *Acta Mech*, vol. 213, no. 3, pp. 267-89, 2010.
- [34] B. Alatas, "ACROA: Artificial chemical reaction optimization algorithm for global optimization," *Expert Sys Appl*, vol. 38, no. 10, pp. 13170-80, 2011.
- [35] F. F. Moghaddam, R. F. Moghaddam, M. Cheriet, "Curved space optimization: A random search based on general relativity theory," *arXiv preprint*, arXiv:1208.2214, 2012.

[36] R. A. Formato, "Central force optimization: A new metaheuristic with applications in applied electromagnetics," *Prog Electrom Res*, vol. 77, pp. 425–91, 2007.

[37] M. Presuss M, "Multimodal Optimization by Means of Evolutionary Algorithms," 1st ed, Switzerland, Springer, 2015. doi:10.1007/978-3-319-07407-8_1

[38] Y. C. Ho, D. L. Pepyne, "Simple Explanation of the No-Free-Lunch Theorem and Its Implications," *J. Optimiz Theor Appl*, vol. 115, no. 3, pp. 549-70, 2002.

[39] H. J. Leavitt, "Applied organization and readings. Changes in industry: structural, technical and human approach," in: Cooper, WW, et al. *New Perspectives in Organization Research*, New York, Wiley, 1962.

[40] D. Katz, R. L. Kahn, "The social psychology of organizations," New York, Wiley, 1978. ISBN: 978-0-471-02355-5.

[41] T. Manichander, "Educational management," USA, Lulu Publication, 2016. ISBN: 978-1-329-76327-2.

[42] R. B. Rudani, "Principles of management.," New Delhi, McGraw-Hill Education, 2013. ISBN: 978-1-259-02696-6.

[43] R. Benfari, "Understanding your management style: Beyond the Myers-Briggs Type Indicator," New York, Lexington Books, 1991. ISBN: 978-0-669-24814-2.

[44] M. Benefiel, "Using discernment to make better business decisions," in: G. Flynn (editor), *Leadership and Business Ethics*, Netherlands, Springer, pp. 31-38, 2008.

[45] H. Gospel, A. Pendleton, "Corporate governance and labour management.," New York, Oxford university press, 2005. doi:10.1093/acprof:oso/9780199263677.001.0001

[46] P. N. Suganthan, N. Hansen, J. J. Liang, K. Deb, Y. P. Chen, A. Auger, S. Tiwari, "Problem definitions and evaluation criteria for the CEC 2005 special session on real-parameter optimization," Computational Intelligence Laboratory, Zhengzhou University, China and Nanyang Technological University, Singapore, Technical Report, 2005. <http://www.ntu.edu.sg/home/EPNSugan>

[47] J. J. Liang, B. Y. Qu, P. N. Suganthan, "Problem definitions and evaluation criteria for the CEC 2014 special session and competition on single objective real-parameter numerical optimization" Computational Intelligence Laboratory, Zhengzhou University, Zhengzhou, China and Nanyang Technical Report, Nanyang Technological University, Singapore, 2013.

[48] S. Da, P. Suganthan, "Problem definitions and evaluation criteria for CEC2011 competition on testing evolutionary algorithms on real world optimization problems" Jadavpur University and Nanyang Technological University, Technical report, 2011.

[49] D. Karaboga, B. Gorkemli, C. Ozturk, N. Karaboga, "A comprehensive survey: artificial bee colony (ABC) algorithm and applications", *Artif Intell Rev*, vol. 42, no. 1, pp. 21-57, 2014.

[50] E. Rashedi, H. Nezamabadi-Pour, S. Saryazdi, "GSA: a gravitational search algorithm," *Inf Sci*, vol. 179, pp. 2232–48, 2009.

[51] M. Ghaemi, M. R. Feizi-Derakhshi, "Forest Optimization Algorithm," *Expert Sys Appl*, vol. 41, no. 15, pp. :6676-6687, 2014.

[52] S. A. Mirjalili, A. Lewis, "The Whale Optimization Algorithm," *Adv Eng Softw*, 95, pp. 51-67, 2016.

[53] M. Z. Ali, N. H. Award, P. N. Suganthan, R. G. Reynolds, "A modified cultural algorithm with a balanced performance for the differential evolution frameworks," *Knowl-Based Syst*, vol. 111, pp. 73-86, 2016.

[54] M. El-Abd, "Performance assessment of foraging algorithms vs Evolutionary algorithms," *Inf Sci*, vol. 182, pp. 243–63, 2012.

[55] F. Wilcoxon, "Individual comparisons by ranking methods," *Biometrics*, vol. 1, pp. 80-83, 1945.

[56] A. K. Qin, P. N. Suganthan, "Self-adaptive differential evolution algorithm for numerical optimization," in: *IEEE Congress on Evolutionary Computation*, Edinburgh, Scotland, pp. 1785-1791, 2005.

TABLE III. MINIMIZATION RESULTS FOR 14 FUNCTIONS FROM CEC 2005 OVER 30 RUNS AT 10 DIMENSIONS

CEC05	Search method						
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank
f_1	1.44E-19(6.42E-19) 2	1.43E-12(1.50E-12) +3	2.32E-05(8.99E-06) +4	3.07E+00(7.12E-01) +7	2.92E-02(2.32E-03) +5	6.43E-02(2.67E-01) +6	1.64E-24(1.41E-24) -1
f_2	1.59E-05(2.07E-05) 1	2.57E+00(1.03E+00) +2	1.29E-02(5.28E-02) +2	1.13E+01(2.49E+00) +3	4.80E+03(2.18E+03) +7	2.88E+02(3.34E+02) +5	5.05E+01(2.29E+01) +4
f_3	1.53E+05(9.24E+04) 1	3.11E+05(9.44E+04) +2	1.59E+05(1.04E+05) ≈1	1.67E+05(8.77E+04) ≈1	1.19E+06(1.11E+06) +3	1.72E+06(1.45E+06) ≈1	2.50E+06(1.03E+06) +4
f_4	2.69E-01(4.17E-01) 1	5.03E+01(1.91E+01) +3	6.97E-01(4.23E-01) +2	1.37E+01(4.14E+00) +2	1.18E+04(5.93E+03) +6	8.86E+02(8.87E+02) +5	2.63E+02(8.45E+01) +4
f_5	3.90E+01(4.07E+01) 4	1.36E-01(7.43E-01) -1	1.36E+00(3.50E+00) -3	1.52E+03(1.35E+02) +5	1.72E+03(1.86E+03) +5	1.58E+03(9.25E+02) +5	6.14E-01(8.52E-01) -2
f_6	3.91E+01(1.62E+02) 3	7.38E+00(1.58E+00) -1	2.31E+03(2.98E+03) +4	4.62E+02(2.02E+02) +4	6.90E+03(1.44E+04) +4	4.55E+03(8.41E+03) +4	7.21E+00(5.07E+00) -2
f_7	2.22E+00(1.19E+00) 3	1.52E-01(2.87E-01) -1	6.21E-01(3.71E-01) -2	5.00E-01(8.28E-02) -2	5.92E+00(6.13E+00) +4	2.58E+01(2.28E+01) +5	5.38E-01(1.31E-01) -2
f_8	2.01E+01(7.47E-02) 1	2.04E+01(6.62E-02) +5	2.02E+01(4.50E-02) +2	2.04E+01(6.39E-02) +5	2.02E+01(8.92E-02) +4	2.02E+01(1.24E-01) +3	2.04E+01(7.43E-02) +5
f_9	6.63E-02(2.52E-01) 2	2.52E+01(4.30E+00) +4	4.74E+01(1.52E+01) +6	5.51E+01(7.52E+00) +6	3.28E+01(1.26E+01) +5	2.22E+01(1.08E+01) +3	0.00E+00(0.00E+00) -1
f_{10}	3.08E+01(1.18E+01) 2	3.47E+01(4.35E+00) ≈2	7.86E+01(1.86E+01) +4	7.15E+01(9.07E+00) +4	5.59E+01(2.24E+01) +3	2.46E+01(1.28E+01) -1	2.43E+01(4.39E+00) -1
f_{11}	6.66E+00(1.48E+00) 3	7.16E+00(5.39E-01) ≈3	8.10E+00(1.54E+00) +4	9.10E+00(6.67E-01) +5	8.09E+00(1.22E+00) +4	5.34E+00(1.47E+00) -2	6.64E+00(6.77E-01) ≈3
f_{12}	3.46E+02(1.31E+02) 1	1.79E+03(9.87E+02) +2	6.37E+03(8.86E+03) +3	2.75E+04(5.72E+03) +4	5.03E+03(5.42E+03) +3	3.66E+03(4.77E+03) +3	3.75E+02(1.12E+02) ≈1
f_{13}	5.72E-01(1.72E-01) 1	2.14E+00(2.74E-01) +3	5.01E+00(1.51E+00) +5	8.84E+00(1.02E+00) +6	2.79E+00(1.31E+00) +4	1.51E+00(7.43E-01) +2	5.95E-01(9.99E-02) ≈1

f_{14}	3.53E+00(3.30E-01) 1	3.50E+00(1.34E-01) ≈1	3.84E+00(4.04E-01) +/2	3.60E+00(2.72E-01) ≈1	3.68E+00(3.31E-01) +/2	3.60E+00(4.45E+01) ≈1	3.69E+00(1.21E-01) +/2
+/-	8/3/3	11/1/2	11/2/1	14/0/0	10/2/2	5/3/6	
Avg-rank	1.86	2.36	3.14	3.93	4.21	3.28	2.36

"+", "≈", "-" respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

TABLE IV. MINIMIZATION RESULTS FOR 14 FUNCTIONS FROM CEC 2005 OVER 30 RUNS AT 30 DIMENSIONS

CEC05	Search method							
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank	
f_1	3.61E-16(5.09E-17) 2	1.25E-15(6.46E-16) +/3	3.39E-04(8.68E-05) +/4	1.48E+01(1.19E+00) +/6	4.12E-01(1.26E-01) +/5	3.74E+02(5.63E+02) +/7	1.26E-18(7.18E-19) -/1	
f_2	7.13E-01(3.91E-01) 1	2.94E+03(4.59E+02) +/4	2.84E+01(9.88E+00) +/2	1.91E+02(3.03E+01) +/3	5.25E+04(1.21E+04) +/7	2.06E+04(5.44E+03) +/6	1.48E+04(2.49E+03) +/5	
f_3	2.24E+06(1.21E+06) 2	2.87E+07(5.11E+06) +/6	5.40E+06(1.68E+06) +/3	1.07E+06(2.40E+05) -/1	1.90E+07(9.73E+06) +/4	6.22E+07(3.57E+07) +/5	6.04E+07(9.90E+06) +/5	
f_4	3.71E+03(1.44E+03) 3	1.13E+04(1.58E+03) +/4	1.02E+02(3.43E+01) -/1	2.14E+02(2.91E+01) -/2	1.48E+05(3.30E+04) +/7	4.76E+04(1.23E+04) +/6	3.06E+04(4.80E+03) +/5	
f_5	8.11E+03(1.71E+03) 4	3.31E+03(1.58E+02) -/2	7.45E+03(1.93E+03) ≈/4	4.58E+02(3.38E+01) -/1	1.68E+04(5.53E+03) +/6	1.06E+04(2.69E+03) +/5	6.17E+03(9.55E+02) -/3	
f_6	6.77E+02(1.15E+02) 3	3.85E+01(2.19E+01) -/1	2.91E+03(3.36E+03) +/4	2.93E+03(5.19E+02) +/4	1.18E+04(8.20E+03) +/5	1.01E+08(1.29E+08) +/6	4.81E+01(2.48E+01) -/2	
f_7	1.35E-02(7.67E-03) 1	4.55E-01(9.83E-02) +/2	1.29E+00(6.21E-02) +/4	9.43E-01(5.24E-02) +/3	6.37E+01(8.40E+01) +/5	8.47E+02(2.76E+02) +/6	8.82E-01(1.25E-01) +/3	
f_8	2.03E+01(6.32E-02) 1	2.09E+01(6.21E-02) +/4	2.05E+01(4.42E-02) +/2	2.09E+01(4.42E-02) +/4	2.07E+01(9.64E-02) +/3	2.07E+01(1.38E-01) +/3	2.09E+01(4.21E-02) +/5	
f_9	9.86E+00(3.15E+00) 2	1.89E+02(9.20E+00) +/4	2.46E+02(4.04E+01) +/6	2.02E+02(9.41E+00) +/5	2.13E+02(4.66E+01) +/7	1.48E+01(4.16E+01) +/3	1.00E-01(3.06E-01) -/1	
f_{10}	2.99E+02(5.70E+01) 3	2.20E+02(1.22E+01) -/2	5.86E+02(6.02E+01) +/5	2.46E+02(1.59E+01) ≈/3	4.29E+02(9.57E+01) +/4	2.18E+02(8.34E+01) -/2	2.04E+02(1.48E+01) -/1	
f_{11}	3.22E+01(3.44E+00) 2	3.66E+01(1.22E+00) +/5	3.55E+01(3.05E+00) +/4	3.99E+01(9.91E-01) +/6	3.56E+01(3.38E+00) +/4	2.66E+01(3.25E+00) -/1	3.30E+01(1.24E+00) +/3	
f_{12}	1.23E+04(1.05E+04) 2	1.84E+05(2.23E+04) +/5	7.46E+03(5.15E+03) -/1	6.38E+05(8.88E+04) +/6	1.45E+05(9.66E+04) +/4	6.89E+04(2.24E+01) +/3	7.55E+04(1.03E+04) +/3	
f_{13}	3.05E+00(9.80E-01) 1	1.63E+01(8.61E-01) +/4	3.04E+01(5.93E+00) +/6	3.96E+02(1.55E+02) +/6	2.02E+01(7.06E+00) +/5	1.28E+01(5.37E+00) +/3	5.49E+00(5.17E-01) +/2	
f_{14}	1.31E+01(4.34E-01) 1	1.32E+01(1.84E-01) ≈/1	1.32E+01(4.62E-01) ≈/1	1.32E+01(2.24E-01) ≈/1	1.34E+01(3.96E-01) +/2	1.35E+01(2.87E-01) +/3	1.34E+01(1.87E-01) +/2	
+/-	10/1/3	10/2/2	10/2/2	9/2/3	14/0/0	12/0/2	9/0/5	
Avg-rank	2.00	3.36	3.35	3.64	4.78	4.00	2.93	

"+", "≈", "-" respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

TABLE V. MINIMIZATION RESULTS FOR 14 FUNCTIONS FROM CEC 2005 OVER 30 RUNS AT 50 DIMENSIONS

CEC05	Search method							
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank	
f_1	4.62E-16(6.12E-16) 1	4.74E-15(3.35E-15) +/2	9.84E-04(2.01E-04) +/4	1.77E+01(1.57E+00) +/5	1.85E+00(9.14E-01) +/5	4.99E+03(3.98E+03) +/6	2.23E-14(1.11E-14) +/3	
f_2	2.76E+01(9.03E+00) 1	2.43E+04(1.92E+03) +/4	3.66E+02(5.31E+01) +/2	5.11E+02(8.93E+01) +/3	1.32E+05(2.30E+04) +/7	7.45E+04(1.87E+04) +/6	6.91E+04(8.25E+03) +/5	
f_3	2.98E+06(6.56E+05) 2	1.08E+08(1.13E+07) +/5	1.45E+07(3.33E+06) +/3	1.43E+06(2.78E+05) +/1	5.41E+07(1.93E+07) +/4	2.33E+08(9.67E+07) +/6	2.37E+08(4.36E+07) +/6	
f_4	2.42E+04(5.10E+03) 3	4.49E+04(5.15E+03) +/4	1.94E+03(1.16E+03) -/2	5.86E+02(9.07E+01) -/1	4.87E+05(1.68E+05) +/7	1.65E+05(5.18E+04) +/6	1.02E+05(1.31E+04) +/5	
f_5	2.03E+04(2.48E+03) 4	8.63E+03(5.27E+02) -/2	1.97E+04(3.34E+03) ≈/4	2.52E+03(2.17E+02) -/1	2.93E+04(3.27E+03) +/6	2.36E+04(2.50E+03) +/5	1.64E+04(1.33E+03) -/3	
f_6	7.54E+02(1.40E+03) 3	4.72E+01(4.04E+00) -/1	2.64E+03(4.10E+03) +/4	4.43E+03(6.00E+02) +/4	1.62E+04(1.05E+03) +/5	1.25E+09(7.34E+08) +/6	9.68E+01(2.98E+01) -/2	

f_7	5.32E-03(1.01E-02) 1	7.88E-01(1.20E-01) +/2	1.71E+00(1.49E-01) +/5	9.53E-01(2.33E-02) +/3	1.35E+02(1.34E+02) +/6	2.18E+03(4.35E+02) +/7	1.07E+00(2.92E-02) +/4
f_8	2.03E+01(5.42E-02) 1	2.11E+01(4.45E-02) +/4	2.09E+01(4.54E-02) +/3	2.11E+01(3.62E-02) +/4	2.09E+01(1.09E-01) +/3	2.08E+01(9.03E-02) +/2	2.12E+01(3.29E-02) +/5
f_9	3.55E+01(6.80E+00) 2	3.96E+02(1.48E+01) +/4	5.19E+02(5.54E+01) +/7	4.65E+02(1.59E+01) +/6	4.12E+02(6.68E+01) +/5	3.42E+02(6.39E+01) +/3	1.29E+00(1.56E+00) -/1
f_{10}	6.50E+02(7.30E+01) 5	4.36E+02(1.95E+01) -/1	1.13E+03(1.18E+02) +/8	5.85E+02(1.65E+01) -/4	8.31E+02(9.10E+01) +/6	5.03E+02(1.11E+02) -/3	4.44E+02(2.14E+01) -/2
f_{11}	6.35E+01(4.92E+00) 2	6.95E+01(1.34E+00) +/4	6.39E+01(4.45E+00) \approx /2	7.29E+01(1.60E+00) +/6	6.80E+01(3.97E+00) +/5	5.16E+02(5.41E+00) +/7	6.62E+01(1.29E+00) +/3
f_{12}	5.03E+04(3.27E+04) 2	8.52E+05(8.47E+04) +/5	2.19E+05(1.39E+05) +/3	3.17E+06(2.50E+05) +/6	5.14E+05(3.85E+05) +/4	4.93E+05(1.52E+05) +/4	4.00E+05(4.49E+04) +/3
f_{13}	7.03E+00(1.74E+00) 1	3.39E+01(1.14E+00) +/3	6.29E+01(7.09E+00) +/5	2.40E+02(3.51E+01) +/6	4.75E+01(1.19E+01) +/4	5.54E+01(2.24E+01) +/4	1.60E+01(7.57E-01) +/2
f_{14}	2.27E+01(3.59E-01) 1	2.30E+01(1.77E-01) +/2	2.29E+01(4.78E-01) \approx/1	2.30E+01(2.31E-01) +/2	2.29E+01(4.94E-01) \approx/1	2.31E+01(4.27E-01) +/3	2.31E+01(1.68E-01) +/3
	+/ \approx -	11/0/3	10/2/1	11/0/3	13/1/0	13/0/1	10/0/4
	Avg-rank 2.07	3.07	3.71	3.71	4.86	4.85	3.21

“+”, “ \approx ”, “-” respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

TABLE VI. MINIMIZATION RESULTS FOR 30 FUNCTIONS FROM CEC 2014 OVER 30 RUNS AT 10 DIMENSIONS

CEC14	Search method							
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank	
f_1	3.76E+04(3.61E+04) 1	6.61E+05(2.78E+05) +/4	9.26E+05(5.73E+05) +/5	1.57E+05(9.53E+04) +/3	9.92E+06(3.37E+06) +/6	1.20E+07(1.06E+07) +/7	6.68E+04(3.90E+04) +/2	
f_2	1.20E+03(1.37E+03) 2	4.97E+04(2.35E+04) +/3	3.23E+03(4.36E+03) \approx /2	2.64E+06(6.84E+05) +/5	6.41E+04(4.07E+04) +/4	3.44E+03(3.59E+03) \approx /2	6.98E+01(5.56E+01) -/1	
f_3	1.49E+03(1.81E+03) 3	6.96E+03(2.03E+03) +/5	2.23E+02(1.65E+02) -/2	1.19E+04(4.72E+03) +/6	6.64E+04(2.66E+04) +/7	6.22E+03(6.49E+03) +/4	1.29E+02(1.72E+02) -/1	
f_4	1.93E+01(1.68E+01) 3	1.73E+00(2.70E+00) -/1	2.88E+01(2.01E+01) +/4	3.77E+01(1.00E+01) +/7	3.35E+01(2.00E+01) +/6	3.73E+01(1.88E+01) +/5	4.58E+00(5.72E+00) -/2	
f_5	2.00E+01(3.47E-04) 2	2.03E+01(2.43E-01) +/6	2.00E+01(2.04E-02) +/4	1.92E+01(4.20E+00) -/1	2.03E+01(1.86E-01) +/5	2.00E+01(6.60E-02) +/3	2.00E+01(3.03E-01) +/4	
f_6	4.09E+00(1.57E+00) 3	4.88E+00(7.41E-01) +/4	6.37E+00(1.43E+00) +/5	2.75E+00(1.26E+00) -/2	7.47E+00(1.60E+00) +/6	2.59E+00(1.68E+00) -/2	2.38E+00(5.95E-01) -/1	
f_7	3.77E-01(1.93E-01) 4	2.235E-01(5.52E-02) -/3	1.34E-01(6.63E-02) -/2	8.36E-01(8.71E-02) +/5	9.35E-01(5.15E-01) +/5	2.41E-01(2.18E-01) -/3	6.18E-02(2.17E-02) -/1	
f_8	1.00E-01(3.06E-01) 2	2.37E+01(3.06E+00) +/4	4.25E+01(1.71E+01) +/6	1.35E+01(3.61E+00) +/3	3.74E+01(1.19E+01) +/5	2.38E+01(1.02E+01) +/4	0.00E+00(0.00E+00) -/1	
f_9	2.02E+01(8.68E+00) 2	3.04E+01(4.38E+00) +/3	4.85E+01(1.49E+01) +/4	1.79E+01(4.69E+00) \approx /2	4.76E+01(1.57E+01) +/4	2.45E+01(1.29E+01) \approx /2	5.89E+00(1.49E+00) -/1	
f_{10}	7.12E+01(7.72E+01) 2	1.05E+03(1.23E+02) +/4	1.07E+03(3.12E+02) +/4	5.44E+02(1.34E+02) +/3	5.24E+02(1.73E+02) +/5	5.95E+02(2.64E+02) +/3	9.71E+00(8.23E+00) -/1	
f_{11}	5.78E+02(2.86E+02) 2	1.26E+03(1.42E+02) +/5	1.04E+03(3.23E+02) +/4	5.64E+02(1.39E+02) \approx /2	1.09E+03(3.21E+02) +/4	7.92E+02(3.47E+02) +/3	3.57E+02(1.14E+02) -/1	
f_{12}	1.46E-01(1.10E-01) 1	1.16E+00(2.11E-01) +/6	3.54E-01(1.54E-01) +/3	9.40E-01(1.83E-01) +/5	7.24E-01(3.42E-01) +/4	2.56E-01(1.96E-01) +/2	3.48E-01(6.45E-02) +/3	
f_{13}	3.87E-01(1.53E-01) 3	2.18E-01(3.39E-02) -/2	9.00E-02(2.61E-02) -/1	9.82E-02(1.57E-02) -/1	4.89E-01(1.67E-01) +/4	3.28E-01(1.68E-01) \approx /3	2.09E-01(4.07E-02) -/2	
f_{14}	2.48E-01(6.76E-02) 3	2.17E-01(4.09E-02) \approx /2	5.23E-02(2.70E-02) -/1	4.39E-01(1.66E-02) +/4	3.18E-01(2.18E-01) \approx /3	3.16E-01(9.84E-02) +/3	1.75E-01(2.98E-02) -/2	
f_{15}	2.50E+00(1.12E+00) 5	2.27E+00(2.29E-01) \approx /5	9.12E-01(3.07E-01) -/1	1.97E+00(3.42E-01) +/4	6.26E+00(2.23E+01) +/6	1.75E+00(8.00E-01) -/3	1.09E+00(1.97E-01) -/2	
f_{16}	2.53E+00(5.34E-01) 1	3.31E+00(2.10E-01) +/2	2.87E+00(4.54E-01) +/2	3.24E+00(2.83E-01) +/2	3.26E+00(4.10E-01) +/2	3.21E+00(3.94E-01) +/2	2.50E+00(2.06E-01) \approx/1	
f_{17}	2.84E+03(3.01E+03) 1	8.04E+03(3.92E+03) +/2	4.20E+03(4.00E+03) \approx/1	5.26E+04(8.74E+04) +/3	1.92E+05(3.44E+05) +/6	2.79E+05(3.75E+05) +/6	3.41E+04(1.88E+04) +/4	
f_{18}	7.59E+03(7.23E+03) 3	2.98E+03(1.40E+03) -/2	1.06E+04(1.08E+04) \approx /3	6.67E+03(3.10E+03) +/3	1.18E+04(1.32E+04) \approx /3	8.49E+03(7.50E+03) \approx /3	6.72E+02(5.67E+02) -/1	
f_{19}	1.96E+00(1.07E+00) 2	2.58E+00(2.11E-01) +/3	4.31E+00(1.40E+00) +/4	2.03E+00(5.33E-01) \approx /2	5.01E+00(1.23E+00) +/5	2.77E+00(1.34E+00) +/3	6.29E-01(2.29E-01) -/1	
f_{20}	2.31E+03(3.38E+03) 3	6.81E+02(4.03E+02) \approx /3	9.77E+01(3.15E+01) -/2	4.58E+03(2.83E+03) +/4	3.75E+03(3.47E+03) +/4	1.00E+04(9.68E+03) +/5	4.67E-01(1.77E-01) -/1	

f_{21}	2.02E+03(2.85E+03) 2	2.74E+03(1.10E+03) +/3	4.68E+03(4.37E+03) +/4	5.44E+03(4.74E+03) +/5	1.91E+04(1.71E+04) +/6	5.56E+04(8.41E+04) +/6	1.84E+03(1.81E+03) -/1
f_{22}	7.14E+01(6.55E+01) 2	4.50E+01(8.60E+00) ≈/2	1.57E+02(1.01E+02) +/4	1.55E+02(1.21E+01) +/4	6.68E+01(4.84E+01) ≈/2	1.17E+02(8.49E+01) +/3	2.03E+00(1.73E+00) -/1
f_{23}	3.29E+02(2.40E-09) 2	9.88E+01(1.53E+02) -/1	3.30E+02(6.42E-01) +/4	3.23E+02(2.29E+01) +/5	3.13E+02(4.53E+01) +/4	3.31E+02(3.87E+00) +/4	3.29E+02(7.40E-07) +/3
f_{24}	1.33E+02(9.95E+00) 2	1.34E+02(4.67E+00) ≈/2	1.82E+02(2.44E+01) +/4	1.88E+02(2.77E+01) +/4	1.73E+02(2.62E+01) +/4	1.39E+02(2.19E+01) +/3	1.17E+02(3.00E+00) -/1
f_{25}	1.89E+02(2.35E+01) 3	1.69E+02(6.81E+00) -/2	1.99E+02(7.42E+00) ≈/3	1.99E+02(7.15E-01) ≈/3	1.89E+02(1.44E+01) ≈/3	1.94E+02(1.99E+01) ≈/3	1.49E+02(1.07E+02) -/1
f_{26}	1.00E+02(9.80E-02) 1	1.00E+02(3.36E-02) ≈/1	1.00E+02(2.75E-02) ≈/1	1.01E+02(2.92E+00) ≈/1	1.00E+02(1.87E-01) +/2	1.00E+02(1.21E-01) ≈/1	1.00E+02(4.25E-02) ≈/1
f_{27}	1.41E+02(1.90E+02) 2	6.31E+01(3.45E+01) -/1	4.32E+02(1.89E+02) +/5	3.66E+02(5.02E+01) +/3	3.86E+02(1.38E+02) +/4	3.03E+02(1.69E+02) +/3	1.88E+02(1.76E+02) ≈/2
f_{28}	5.20E+02(1.90E+02) 3	3.07E+02(1.60E-01) -/1	8.78E+02(2.14E+02) +/6	6.53E+02(1.89E+02) +/5	5.45E+02(1.17E+01) ≈/3	5.63E+02(1.30E+02) +/4	3.74E+02(4.80E+00) -/2
f_{29}	1.72E+05(5.26E+05) 4	2.04E+02(4.79E-01) -/1	2.54E+05(8.13E+05) +/5	4.04E+06(1.05E+07) +/6	1.73E+05(5.26E+05) +/5	7.96E+05(1.21E+06) +/5	5.99E+02(7.10E+01) -/2
f_{30}	1.12E+03(4.16E+02) 4	2.54E+02(1.28E+01) -/1	1.92E+03(9.31E+02) +/5	8.52E+02(1.59E+02) -/3	1.67E+03(9.33E+02) +/5	1.67E+03(5.54E+02) +/5	5.33E+02(4.05E+01) -/2
	+/≈/-	14/6/10	19/4/6	20/6/4	25/5/0	21/6/3	5/3/22
	Avg-rank	2.80	3.36	3.53	4.73	3.50	1.63
		2.46					

“+”, “≈”, “-” respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

TABLE VII. MINIMIZATION RESULTS FOR 30 FUNCTIONS FROM CEC 2014 OVER 30 RUNS AT 30 DIMENSIONS

CEC14	Search method						
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank
f_1	1.21E+06(5.53E+05) 3	8.57E+07(1.29E+07) +/7	3.56E+07(1.52E+07) +/5	5.94E+05(1.25E+05) -/1	3.08E+07(1.58E+06) +/5	7.47E+07(4.04E+07) +/6	1.76E+07(5.96E+06) +/4
f_2	1.15E+04(8.29E+03) 2	1.55E+06(1.49E+06) +/3	2.24E+08(1.00E+08) +/6	1.94E+07(2.29E+06) +/5	2.99E+06(3.82E+06) +/4	5.52E+08(4.32E+08) +/7	5.84E-9(0.00E+00) -/1
f_3	9.70E+02(1.21E+03) 2	8.60E+04(1.40E+04) +/6	3.56E+03(1.11E+03) +/3	6.42E+03(2.59E+03) +/4	3.51E+04(2.85E+03) +/4	5.03E+04(2.90E+04) +/5	4.30E+01(4.50E+01) -/1
f_4	7.50E+01(4.18E+01) 3	9.44E+00(1.26E+01) -/1	1.55E+02(4.87E+01) +/5	6.00E+01(3.61E+01) -/2	1.82E+02(5.36E+01) +/5	2.45E+02(8.07E+01) +/6	1.03E+02(1.33E+01) +/4
f_5	2.00E+01(1.87E-04) 1	2.09E+01(5.00E-2) +/5	2.00E+01(6.15E-04) +/2	2.09E+01(6.40E-02) +/5	2.04E+01(1.92E-01) +/4	2.02E+01(1.44E-01) +/3	2.04E+01(4.20E-02) +/4
f_6	2.19E+01(3.61E+00) 4	3.43E+01(1.21E+00) +/5	3.11E+01(3.81E+00) +/5	1.63E+01(2.55E+00) -/1	3.50E+01(3.54E+00) +/5	2.03E+01(4.20E+00) -/3	1.72E+01(1.32E+00) -/2
f_7	2.75E-02(3.91E-02) 3	6.37E-04(5.66E-04) -/2	2.77E+00(5.23E-01) +/6	1.18E+00(2.35E-02) +/5	9.90E-01(7.26E-02) +/4	4.52E+00(5.53E+00) +/6	1.50E-04(2.42E-04) -/1
f_8	6.88E+00(2.91E+00) 2	1.81E+02(1.38E+01) +/4	2.00E+02(4.16E+01) +/4	1.17E+02(1.02E+01) +/3	1.88E+02(3.90E+01) +/4	1.21E+02(2.92E+01) +/3	7.14E-01(1.06E+00) -/1
f_9	1.23E+02(2.64E+01) 2	2.05E+02(9.97E+00) +/4	2.36E+02(4.73E+01) +/5	1.34E+02(1.28E+01) +/3	2.31E+02(6.62E+01) +/5	1.52E+02(4.80E+01) +/3	8.78E+01(9.36E+00) -/1
f_{10}	4.63E+02(1.77E+02) 2	6.54E+03(2.63E+02) +/7	4.13E+03(6.33E+02) +/5	2.50E+03(3.59E+02) +/3	4.00E+03(8.70E+02) +/5	2.98E+03(5.72E+02) +/4	2.55E+01(2.87E+01) -/1
f_{11}	3.36E+03(5.77E+02) 2	7.14E+03(2.50E+02) +/4	4.43E+03(8.18E+02) +/3	2.87E+03(4.48E+02) -/1	5.19E+03(8.10E+02) +/3	3.62E+03(5.88E+02) ≈/2	3.44E+03(2.75E+02) ≈/2
f_{12}	2.92E-01(8.74E-02) 1	2.43E+00(3.08E-01) +/7	1.15E+00(2.52E-01) +/5	6.61E-01(1.12E-01) +/4	1.62E+00(3.90E-01) +/6	3.83E-01(1.79E-01) +/2	5.93E-01(8.21E-02) +/3
f_{13}	3.88E-01(1.07E-01) 3	3.48E-01(3.32E-02) ≈/3	2.80E-01(8.63E-02) -/2	1.70E-01(1.97E-02) -/1	5.30E-01(8.99E-02) +/4	4.88E-01(1.09E-01) +/4	4.50E-01(4.26E-02) +/4
f_{14}	2.66E-01(8.99E-02) 2	2.39E-01(1.86E-02) -/1	2.31E-01(1.07E-01) -/1	3.50E-01(5.11E-02) +/5	2.71E-01(5.26E-02) +/3	2.99E-01(1.03E-01) +/3	3.08E-01(4.71E-02) +/4
f_{15}	3.26E+01(1.14E+01) 3	1.65E+01(9.23E-01) -/2	1.24E+01(6.07E+00) -/1	1.30E+01(7.77E-01) -/1	6.68E+01(2.52E+01) +/4	3.66E+02(5.95E+02) +/5	1.01E+01(8.53E-01) -/1
f_{16}	1.17E+01(4.25E-01) 2	1.30E+01(1.74E-01) +/5	1.20E+01(6.09E-01) +/3	1.28E+01(3.03E-01) +/4	1.25E+01(5.20E-01) +/3	1.27E+01(5.58E-01) +/3	1.08E+01(2.81E-01) -/1
f_{17}	2.22E+05(1.71E+05) 3	1.38E+06(3.71E+05) +/4	8.57E+05(5.87E+05) +/4	3.23E+04(2.06E+04) -/1	4.12E+06(2.54E+06) +/6	4.42E+06(3.40E+06) +/6	3.02E+06(1.51E+06) +/5
f_{18}	2.20E+03(2.76E+03) 3	7.07E+02(5.12E+02) -/2	3.97E+03(4.44E+03) +/6	3.13E+02(2.63E+02) -/1	5.52E+03(4.43E+03) +/6	3.04E+03(3.65E+03) +/4	9.99E+03(4.54E+03) +/5

f_{19}	1.03E+01(1.93E+01) 1	1.93E+01(7.68E-01) +/3	2.55E+01(1.64E+01) +/4	1.76E+01(2.15E+01) +/2	5.96E+01(5.42E+01) +/4	5.44E+01(3.08E+01) +/4	1.13E+01(6.47E-01) +/1
f_{20}	4.98E+02(1.14E+02) 1	2.78E+04(8.29E+03) +/4	5.52E+02(4.25E+02) ≈1	1.68E+04(3.35E+03) +/3	2.72E+04(2.11E+04) +/4	3.54E+04(1.57E+04) +/5	7.34E+03(3.04E+03) +/2
f_{21}	1.95E+05(1.33E+05) 3	3.97E+05(1.22E+05) +/4	1.95E+05(9.13E+04) ≈/3	5.52E+04(1.88E+04) -/1	9.09E+05(8.13E+05) +/7	8.16E+05(7.30E+05) +/6	4.97E+05(1.71E+05) +/5
f_{22}	5.38E+02(1.63E+02) 2	6.13E+02(8.92E+01) +/3	6.57E+02(1.90E+02) +/4	9.52E+02(1.94E+02) +/6	8.15E+02(2.20E+02) +/5	5.55E+02(2.56E+02) ≈/2	1.92E+02(7.98E+01) -/1
f_{23}	3.15E+02(3.09E-04) 1	3.14E+02(9.16E-04) ≈/1	3.51E+02(1.48E+01) +/6	3.16E+02(1.76E-01) +/4	3.33E+02(1.08E+01) +/5	3.30E+02(1.03E+01) +/5	3.17E+02(1.86E+00) +/4
f_{24}	2.29E+02(7.53E+00) 4	2.27E+02(5.23E-01) ≈/4	2.64E+02(2.73E+01) +/5	2.08E+02(5.37E-01) -/1	2.07E+02(4.55E+00) -/2	2.49E+02(7.43E+00) +/5	2.25E+02(1.75E+00) -/3
f_{25}	2.21E+02(3.49E+00) 4	2.15E+02(2.78E+00) -/3	2.32E+02(1.83E+01) +/5	2.02E+02(7.92E-02) -/1	2.19E+02(1.91E+01) ≈/4	2.19E+02(6.64E+00) ≈/4	2.07E+02(1.08E+00) -/2
f_{26}	1.24E+02(4.28E+01) 3	1.00E+02(6.05E-02) -/1	1.65E+02(6.64E+01) +/6	1.78E+02(4.12E+01) +/6	1.00E+02(1.24E-01) -/1	1.27E+02(4.52E+01) +/4	1.00E+02(5.41E-02) -/2
f_{27}	5.83E+02(3.42E+02) 2	7.58E+02(1.46E+02) +/4	8.36E+02(4.27E+02) +/4	8.39E+02(3.89E+02) +/5	9.66E+02(4.05E+02) +/5	6.66E+02(2.08E+02) +/3	4.43E+02(2.29E+01) -/1
f_{28}	2.32E+03(3.39E+02) 5	4.25E+02(5.55E+00) -/1	3.95E+03(1.05E+03) +/6	2.40E+03(5.72E+02) ≈/5	2.33E+03(5.56E+02) -/3	2.09E+03(5.65E+02) -/2	8.65E+02(2.64E+01) -/4
f_{29}	2.79E+05(2.19E+06) 4	2.31E+02(3.84E+00) -/1	1.67E+04(1.94E+04) -/3	4.08E+07(5.88E+07) +/7	5.34E+06(4.77E+06) +/6	2.61E+06(3.34E+06) +/5	4.62E+04(1.01E+04) +/5
f_{30}	3.61E+03(8.28E+02) 3	7.60E+02(1.25E+02) -/1	2.66E+04(2.80E+04) +/5	1.97E+03(3.86E+02) -/2	6.86E+04(3.13E+04) +/6	3.08E+04(3.58E+04) +/5	5.60E+03(1.88E+03) +/4
	+/-/≈-	17/3/10	23/3/4	17/1/12	26/1/3	25/3/2	13/1/16
	-						
	Avg-rank 2.53	3.40	4.10	3.06	4.33	4.40	2.56

“+”, “≈”, “-” respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

TABLE VIII. MINIMIZATION RESULTS FOR 30 FUNCTIONS FROM CEC 2014 OVER 30 RUNS AT 50 DIMENSIONS

CEC14	Search method						
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank
f_1	2.36E+06(9.42E+05) 2	4.64E+08(5.76E+07) +/6	1.03E+08(2.55E+07) +/5	1.21E+06(1.87E+05) -/1	3.13E+07(8.67E+06) +/3	7.57E+07(4.18E+07) +/4	8.57E+07(1.29E+07) +/7
f_2	5.90E+03(5.79E+03) 2	3.06E+01(1.43E+02) -/1	7.11E+08(1.93E+08) +/7	2.75E+07(2.15E+06) +/5	2.16E+07(1.66E+07) +/4	5.04E+08(1.31E+09) +/6	1.53E+04(9.99E+03) +/3
f_3	1.95E+03(9.38E+02) 1	2.52E+05(3.22E+04) +/7	5.88E+03(1.09E+03) +/2	1.33E+04(2.62E+03) +/4	4.37E+04(1.21E+04) +/5	5.23E+04(2.51E+04) +/6	9.39E+03(3.23E+03) +/3
f_4	1.04E+02(4.74E+01) 4	4.38E+01(1.26E+00) -/2	4.47E+02(1.17E+02) +/6	6.28E+01(2.87E+01) -/3	2.37E+02(6.19E+01) +/5	2.38E+02(5.97E+01) +/5	1.14E+02(1.02E+01) ≈/4
f_5	2.00E+01(2.40E-05) 1	2.11E+01(2.45E-02) +/6	2.00E+01(2.38E-04) +/2	2.11E+01(3.30E-02) +/6	2.04E+01(1.66E-01) +/4	2.02E+01(9.20E-02) +/3	2.06E+01(3.65E-02) +/5
f_6	4.66E+01(5.02E+00) 3	6.61E+01(1.52E+00) +/4	5.94E+01(4.62E+00) +/4	3.44E+01(3.50E+00) -/2	6.51E+01(4.17E+00) +/4	1.93E+01(3.24E+00) -/1	4.53E+01(1.78E+00) +/3
f_7	1.03E-02(9.58E-03) 2	3.17E-03(3.95E-03) -/1	8.99E+00(2.38E+00) +/6	1.26E+00(2.29E-02) +/4	1.17E+00(1.25E-01) +/3	6.39E+00(4.97E+00) +/5	5.40E-03(4.10E-03) -/1
f_8	3.12E+01(7.48E+00) 2	3.95E+02(8.76E+00) +/4	3.93E+02(6.16E+01) +/7	2.54E+02(1.54E+01) +/5	3.30E+02(5.27E+01) +/6	1.26E+02(3.78E+01) +/3	1.36E+00(1.33E+00) -/1
f_9	3.05E+02(5.00E+01) 3	4.16E+02(1.22E+01) +/4	5.34E+02(7.37E+01) +/5	2.56E+02(1.92E+01) -/2	4.52E+02(9.23E+01) +/4	1.39E+02(3.64E+01) -/1	2.88E+02(1.26E+01) -/2
f_{10}	9.65E+02(2.44E+02) 2	1.28E+04(3.16E+02) +/6	7.70E+03(8.65E+02) +/5	5.29E+03(5.13E+02) +/4	7.45E+03(1.17E+03) +/5	3.07E+03(6.11E+02) +/3	7.7110E+01(6.65E+01) -/1
f_{11}	6.34E+03(9.18E+02) 2	1.37E+04(3.61E+02) +/5	7.60E+03(8.43E+02) +/3	6.22E+03(4.55E+02) ≈/2	8.94E+03(1.55E+03) +/3	3.73E+03(6.42E+02) -/1	9.13E+03(3.88E+02) +/4
f_{12}	4.07E-01(9.21E-02) 2	3.29E+00(3.21E-01) +/7	1.77E+00(2.73E-01) +/5	6.50E-01(7.85E-02) +/3	2.60E+00(5.28E-01) +/6	3.54E-01(1.77E-01) -/1	1.01E+00(1.32E-01) +/4
f_{13}	4.97E-01(1.03E-01) 2	5.07E-01(4.33E-02) ≈/2	4.35E-01(8.45E-02) ≈/2	2.29E-01(2.29E-02) -/1	5.88E-01(1.31E-01) +/3	4.50E-01(1.17E-01) ≈/2	5.67E-01(6.11E-02) +/3
f_{14}	3.08E-01(3.84E-02) 2	2.81E-01(2.31E-02) -/1	3.40E-01(1.45E-01) +/3	3.80E-01(3.72E-02) ≈/2	4.00E-01(1.50E-01) +/4	3.70E-01(3.98E-01) ≈/2	3.38E-01(4.06E-02) +/3
f_{15}	9.70E+01(2.42E+01) 4	3.45E+01(9.53E-01) -/3	1.02E+02(7.29E+01) ≈/4	2.43E+01(1.15E+00) -/1	2.43E+02(4.89E+01) +/6	2.41E+02(4.15E+02) ≈/4	3.03E+01(1.22E+00) -/2
f_{16}	2.08E+01(6.99E-01) 3	2.28E+01(1.67E-01) +/5	2.11E+01(7.16E-01) +/4	2.20E+01(6.20E-01) +/4	2.22E+01(5.45E-01) +/4	1.27E+01(6.14E-01) -/1	2.06E+01(2.55E-01) -/2
f_{17}	3.24E+05(1.82E+05) 3	1.06E+07(2.41E+06) +/5	3.83E+06(1.36E+06) +/4	3.53E+04(1.89E+04) -/1	1.68E+07(9.44E+06) +/6	3.84E+06(2.73E+06) +/4	1.06E+07(3.74E+06) +/5

f_{18}	1.62E+03(1.30E+03)) 1	2.44E+03(1.43E+03)) +/2	2.62E+05(1.18E+06)) +/4	1.79E+03(5.82E+02)) ≈/1	7.87E+03(7.02E+03)) +/3	1.06E+04(3.94E+04)) +/4	8.13E+03(4.61E+03)) +/3
f_{19}	5.54E+01(3.06E+01)) 3	3.57E+01(6.18E-01)) ≈/3	7.44E+01(2.08E+01)) +/4	2.00E+01(6.14E+00)) -/2	7.65E+01(2.80E+01)) +/4	6.20E+01(3.51E+01)) ≈/3	3.33E+01(6.47E+00)) ≈/3
f_{20}	4.82E+02(2.00E+02)) 1	9.15E+04(2.80E+04)) +/5	1.12E+03(2.56E+02)) +/2	1.96E+04(3.96E+03)) +/3	1.14E+05(1.96E+05)) +/5	3.33E+04(1.88E+04)) +/4	2.16E+04(6.66E+03)) +/3
f_{21}	2.73E+05(1.59E+05)) 3	4.39E+06(9.70E+05)) +/5	1.94E+06(1.15E+06)) +/5	4.15E+04(1.01E+04)) -/1	6.36E+06(3.69E+06)) +/6	9.33E+05(8.40E+05)) +/4	5.14E+06(2.44E+06)) +/5
f_{22}	1.48E+03(3.24E+02)) 3	1.89E+03(1.51E+02)) +/4	1.38E+03(3.65E+02)) ≈/3	1.76E+03(3.76E+02)) +/4	1.97E+03(4.89E+02)) +/5	6.41E+02(1.76E+02)) -/1	8.20E+02(1.58E+02)) -/2
f_{23}	3.44E+02(1.34E-04)) 5	3.37E+02(5.14E-03)) -/2	4.29E+02(2.62E+01)) +/7	3.43E+02(6.01E+00)) -/4	3.73E+02(4.80E+01)) +/6	3.28E+02(7.79E+00)) -/1	3.44E+02(8.89E-01)) ≈/5
f_{24}	2.71E+02(8.22E+00)) 6	2.53E+02(1.99E+01)) -/3	3.30E+02(2.04E+01)) +/7	2.19E+02(6.76E+00)) -/1	2.00E+02(3.56E-01)) -/2	2.49E+02(7.32E+00)) -/4	2.59E+02(2.75E+00)) -/5
f_{25}	2.51E+02(1.19E+01)) 4	2.48E+02(4.58E+00)) ≈/4	2.67E+02(2.25E+01)) +/5	2.03E+02(1.13E-01)) -/1	2.08E+02(1.96E+01)) -/2	2.22E+02(5.24E+00)) -/3	2.22E+02(2.65E+00)) -/3
f_{26}	1.94E+02(2.53E+01)) 5	1.00E+02(4.44E-02)) -/1	2.05E+02(1.98E+01)) +/6	1.78E+02(4.08E+01)) -/6	1.04E+02(1.82E+01)) -/3	1.42E+02(6.15E+01)) ≈/5	1.01E+02(1.13E-01)) -/2
f_{27}	1.56E+03(2.42E+02)) 3	2.04E+03(4.93E+01)) +/4	2.04E+03(1.58E+02)) +/4	3.31E+03(8.88E+02)) +/6	2.07E+03(1.17E+02)) +/4	7.51E+02(2.07E+02)) -/1	1.36E+03(1.32E+02)) -/2
f_{28}	7.02E+03(1.01E+03)) 6	4.37E+02(1.77E+01)) -/2	8.63E+03(1.48E+03)) +/7	6.07E+03(9.00E+02)) -/5	4.89E+03(1.66E+03)) -/4	2.00E+03(6.21E+02)) -/3	2.94E+01(2.49E+01)) -/1
f_{29}	1.90E+03(7.11E+02)) 2	2.35E+02(1.16E+00)) -/1	4.87E+06(2.61E+07)) +/6	1.77E+08(1.18E+08)) +/7	3.58E+07(1.95E+07)) +/5	1.13E+06(2.87E+06)) +/4	3.48E+06(8.29E+06)) +/5
f_{30}	1.33E+04(2.06E+03)) 4	2.02E+03(1.46E+02)) -/1	1.63E+05(9.26E+04)) +/7	3.01E+04(2.38E+03)) +/5	1.09E+05(9.17E+04)) +/6	4.03E+04(4.21E+04)) +/5	1.18E+04(1.17E+03)) -/3
	+/-/≈	16/3/11	27/3/0	13/3/14	26/0/4	14/5/11	14/3/13
	Avg-rank 2.86	3.56	3.07	3.30	4.30	3.13	3.16

“+”, “≈”, “-” respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

TABLE IX. P-VALUES OBTAINED WITH WILCOXON’S RANK SUM TEST OVER 14 FUNCTIONS FROM CEC 2005 BENCHMARK IN 10D

CE C05	Search Manager vs					
	ABC	FOA	GSA	WOA	CA	DE
f_1	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11
f_2	3.02E-11	1.01E-08	3.01E-08	3.02E-11	3.02E-11	3.02E-11
f_3	3.80E-07	8.77E-01	9.23E-01	4.68E-08	8.10E-10	3.02E-11
f_4	3.02E-11	7.22E-06	3.02E-11	3.02E-11	3.02E-11	3.02E-11
f_5	3.77E-12	6.51E-09	3.02E-11	8.84E-07	3.34E-11	1.69E-09
f_6	2.07E-02	1.60E-06	5.57E-10	5.57E-10	7.38E-10	7.95E-01
f_7	6.30E-11	8.35E-08	1.20E-08	5.26E-04	1.33E-10	1.70E-08
f_8	7.39E-11	1.44E-02	6.69E-11	5.09E-06	4.71E-04	4.97E-11
f_9	1.94E-11	1.94E-11	1.94E-11	1.94E-11	1.94E-11	6.81E-13
f_{10}	1.30E-01	6.06E-11	4.08E-11	9.51E-06	3.78E-02	5.10E-03
f_{11}	7.01E-02	1.11E-03	1.56E-08	2.39E-04	3.00E-03	8.53E-01
f_{12}	7.66E-05	2.27E-03	3.02E-11	2.15E-06	9.03E-04	1.54E-01
f_{13}	3.02E-11	3.02E-11	3.02E-11	1.61E-10	2.44E-09	2.00E-01
f_{14}	6.10E-01	1.08E-02	2.22E-01	4.84E-02	3.63E-01	3.00E-03

TABLE X. P-VALUES OBTAINED WITH WILCOXON’S RANK SUM TEST OVER CEC 2014 BENCHMARK FUNCTIONS IN 50D

CE C14	Search Manager vs					
	ABC	FOA	GSA	WOA	CA	DE
f_1	3.02E-11	3.02E-11	1.69E-08	3.02E-11	3.02E-11	3.02E-11
f_2	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11	2.77E-10
f_3	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11	2.87E-10
f_4	1.07E-07	3.02E-11	6.55E-04	2.92E-09	1.46E-10	1.41E-01
f_5	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11
f_6	3.02E-11	6.06E-11	6.72E-10	3.01E-11	3.02E-11	1.43E-05
f_7	3.03E-02	3.02E-11	3.02E-11	3.02E-11	3.02E-11	8.00E-03
f_8	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11	2.84E-11
f_9	1.55E-09	7.39E-11	1.63E-05	2.03E-09	3.02E-11	1.25E-04
f_{10}	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.34E-11	3.02E-11
f_{11}	3.02E-11	1.86E-06	5.99E-01	8.48E-09	1.09E-10	1.95E-10
f_{12}	3.02E-11	3.02E-11	2.37E-10	3.01E-11	1.50E-02	3.02E-11
f_{13}	4.20E-01	2.51E-02	3.02E-11	8.68E-03	9.63E-02	3.00E-03

f_{14}	3.03E-03	8.41E-01	4.68E-08	3.83E-05	1.86E-01	5.80E-03
f_{15}	3.02E-11	2.11E-01	3.02E-11	4.07E-11	5.30E-01	3.02E-11
f_{16}	3.01E-11	3.51E-02	2.83E-08	1.28E-09	3.02E-11	3.77E-04
f_{17}	3.02E-11	4.61E-10	3.68E-11	3.02E-11	8.99E-11	3.02E-11
f_{18}	1.50E-02	4.71E-04	1.12E-01	6.73E-06	3.51E-02	3.16E-05
f_{19}	1.85E-01	3.91E-02	6.01E-08	3.50E-03	7.73E-01	1.58E-01
f_{20}	3.02E-11	4.20E-10	3.02E-11	3.02E-11	3.02E-11	3.02E-11
f_{21}	3.02E-11	4.50E-11	4.97E-11	3.34E-11	1.86E-06	3.02E-11
f_{22}	3.09E-06	2.40E-01	8.31E-03	4.08E-05	8.15E-11	1.95E-10
f_{23}	3.02E-11	3.02E-11	3.99E-04	8.48E-09	5.57E-10	1.84E-01
f_{24}	1.99E-05	3.02E-11	3.02E-11	3.02E-11	8.48E-09	1.70E-08
f_{25}	2.77E-01	8.12E-04	3.02E-11	2.74E-08	2.87E-10	3.02E-11
f_{26}	2.23E-09	5.07E-10	3.82E-09	5.07E-10	8.89E-10	5.97E-09
f_{27}	3.02E-11	4.97E-11	8.35E-08	4.97E-11	3.02E-11	9.75E-10
f_{28}	3.01E-11	3.37E-05	6.91E-04	3.37E-05	4.80E-07	3.02E-11
f_{29}	3.02E-11	3.02E-11	5.26E-04	3.02E-11	3.02E-11	3.83E-05
f_{30}	3.02E-11	3.02E-11	3.02E-11	3.02E-11	3.02E-11	7.74E-06

TABLE XI. MINIMIZATION RESULTS FOR REAL-WORLD PROBLEMS

Problem	Search method						
	Search Manager Mean(Std.) Rank	ABC Mean(Std.) Comp/Rank	FOA Mean(Std.) Comp/Rank	GSA Mean(Std.) Comp/Rank	WOA Mean(Std.) Comp/Rank	CA Mean(Std.) Comp/Rank	DE Mean(Std.) Comp/Rank
T01 FM	1.79E+01(4.88E+00) 1	1.64E+01(7.05E+00) ≈/1	2.01E+01(3.24E+00) +/2	2.31E+01(1.95E+00) +/3	2.11E+01(3.70E+00) +/2	2.43E+01(1.85E+00) +/4	1.62E+01(3.68E+00) ≈/1
T02L-J	- 1.29E+01(2.97E+00) 2	-4.71E+00(5.15E-01) +/6	- 5.24E+00(1.58E+00) +/6	-5.58E+00(6.45E-01) +/4	- 1.55E+01(5.46E+00) -/1	- 8.41E+00(1.95E+00) +/5	- 1.01E+01(1.41E+00) +/3
T03BCB	1.15E-05(0.00E+00) 1	1.15E-05(0.00E+00) ≈/1	1.15E-05(0.00E+00) ≈/1	1.15E-05(0.00E+00) ≈/1	1.15E-05(0.00E+00) ≈/1	1.15E-05(0.00E+00) ≈/1	1.15E-05(0.00E+00) ≈/1
T04STR	1.42E+01(2.27E-01) 1	1.43E+01(3.51E-01) +/2	1.61E+01(3.17E+00) +/4	2.24E+01(9.64E-01) +/7	2.09E+01(6.88E+00) +/6	1.82E+01(3.33E+00) +/4	2.00E+01(2.31E+00) +/5
T05Si(B)	- 2.88E+01(1.76E+00) 1	- 2.18E+01(1.08E+00) +/4	- 2.81E+01(7.88E+00) +/2	1.13E+01(3.03E+01) +/6	- 2.87E+01(1.76E+00) ≈/1	- 2.33E+01(3.89E+00) +/3	- 2.46E+01(2.51E+00) +/2
T06Si(C)	- 1.83E+01(2.71E+00) 1	- 1.03E+01(1.27E+00) +/5	- 2.20E+01(3.91E+00) -/1	1.36E+01(7.56E+00) +/7	- 1.19E+01(4.44E+00) +/4	6.54E+00(2.06E+00) +/6	- 1.57E+01(2.66E+00) +/2
T07SPR	1.45E+00(2.11E-01) 3	1.67E+00(9.62E-02) +/4	1.35E+00(2.29E-01) -/2	1.90E+00(8.78E-02) +/6	1.86E+00(2.41E-01) +/5	1.30E+00(1.92E-01) -/2	1.57E+00(7.47E-02) +/4
T08TNE	2.20E+02(0.00E+00) 1	2.20E+02(0.00E+00) ≈/1	2.20E+02(0.00E+00) ≈/1	2.20E+02(0.00E+00) ≈/1	2.20E+02(0.00E+00) ≈/1	2.20E+02(0.00E+00) ≈/1	2.20E+02(0.00E+00) ≈/1
T09 LSTP	1.58E+06(3.32E+05) 5	1.25E+06(6.94E+04) -/4	2.52E+06(6.60E+05) +/6	2.07E+07(1.81E+06) +/6	9.05E+05(4.11E+05) -/3	3.25E+07(8.44E+06) +/7	1.03E+05(8.16E+03) -/2
T10CAA	- 1.37E+01(2.78+00) 2	- 1.63E+01(1.36E+00) -/1	- 1.38E+01(2.80E+00) ≈/2	- 1.68E+01(2.99E-01) -/1	-1.05E+01(6.81E-01) +/4	- 1.21E+01(2.25E+00) +/3	- 1.41E+01(1.45E+00) ≈/2
T11.IDE	9.85E+04(5.68E+04) 2	3.22E+06(4.37E+05) +/4	1.85E+06(7.82E+05) +/4	6.41E+04(1.29E+03) -/1	1.31E+06(1.21E+05) +/3	2.96E+06(1.77E+06) +/4	1.03E+08(1.21E+07) +/6
T11.3EL	1.54E+04(3.20E+01) 1	1.54E+04(1.28E+01) ≈/1	3.91E+04(4.80E+04) +/3	9.30E+04(5.87E+04) +/4	1.56E+04(6.07E+01) +/2	1.55E+04(2.20E+01) +/2	1.55E+04(1.37E+01) ≈/1
T11.8HS	1.16E+06(3.25E+05) 1	2.07E+06(3.48E+05) +/3	1.16E+06(3.90E+05) ≈/1	1.72E+06(1.30E+05) +/2	5.21E+06(2.95E+06) +/5	1.38E+06(7.59E+05) ≈/1	2.30E+06(4.27E+05) +/4
T12(me)	2.25E+01(5.64E+00) 1	3.09E+01(3.71E+00) +/3	2.98E+01(5.41E+00) +/3	2.99E+01(4.38E+00) +/3	3.69E+01(7.58E+00) +/4	2.64E+01(6.92E+00) +/2	2.28E+01(2.56E+00) ≈/1
T13(Ca)	2.71E+01(4.11E+00) 2	3.10E+01(2.52E+00) +/3	3.35E+01(6.94E+00) +/3	4.44E+01(4.30E+00) +/4	4.19E+01(5.69E+00) +/4	3.14E+01(7.51E+00) +/3	2.15E+01(1.46E+00) -/1
+/-	9/4/2	9/4/2	11/2/2	10/3/2	11/3/1	7/6/2	
Avg-rank	1.66	2.86	2.73	3.73	3.06	3.20	2.40

“+”, “≈”, “-” respectively denote that the performance of Search Manager is better than, similar to, and worse than the corresponding algorithm.

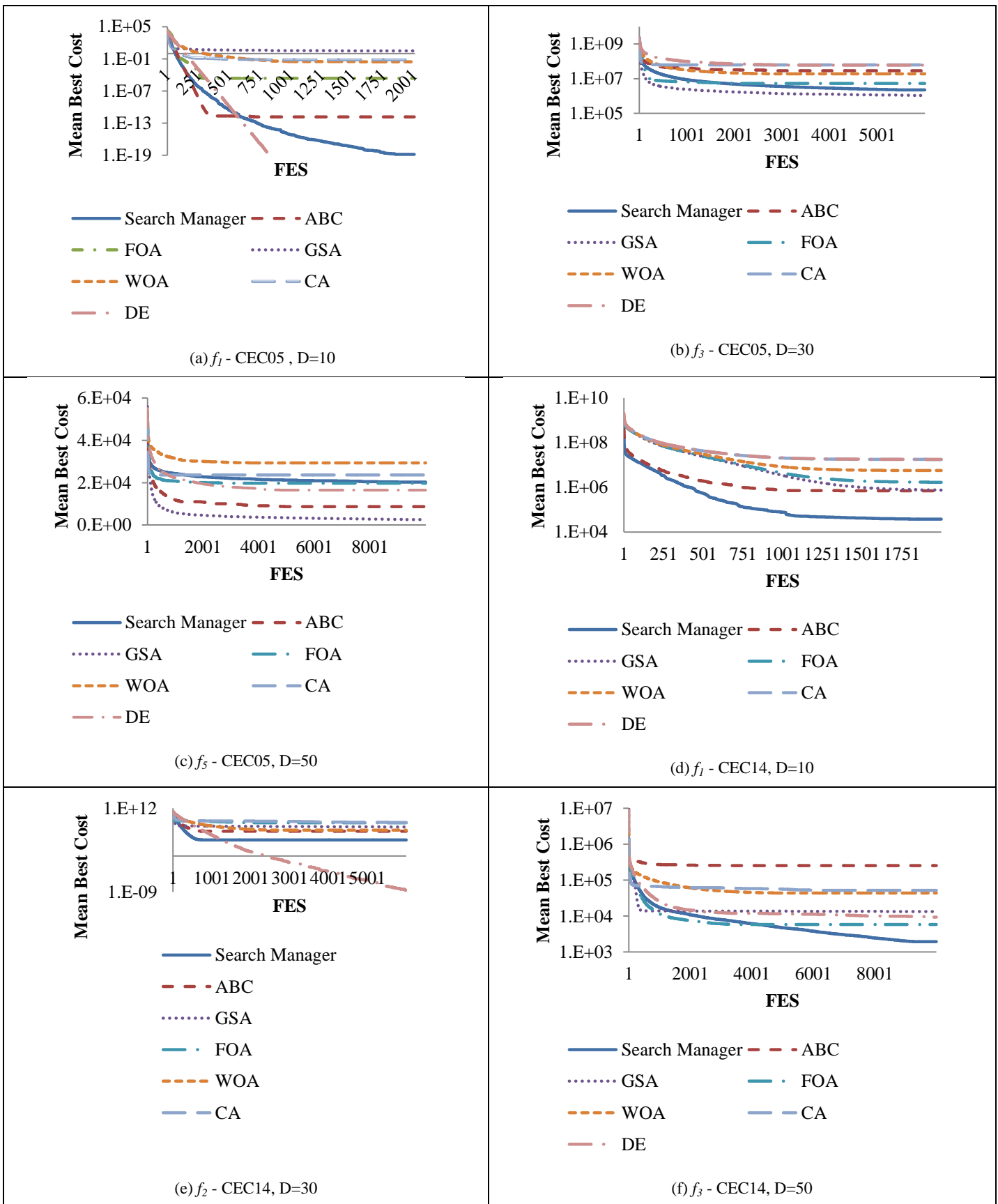


Fig. 8. Convergence curves of the algorithms on unimodal benchmark functions.

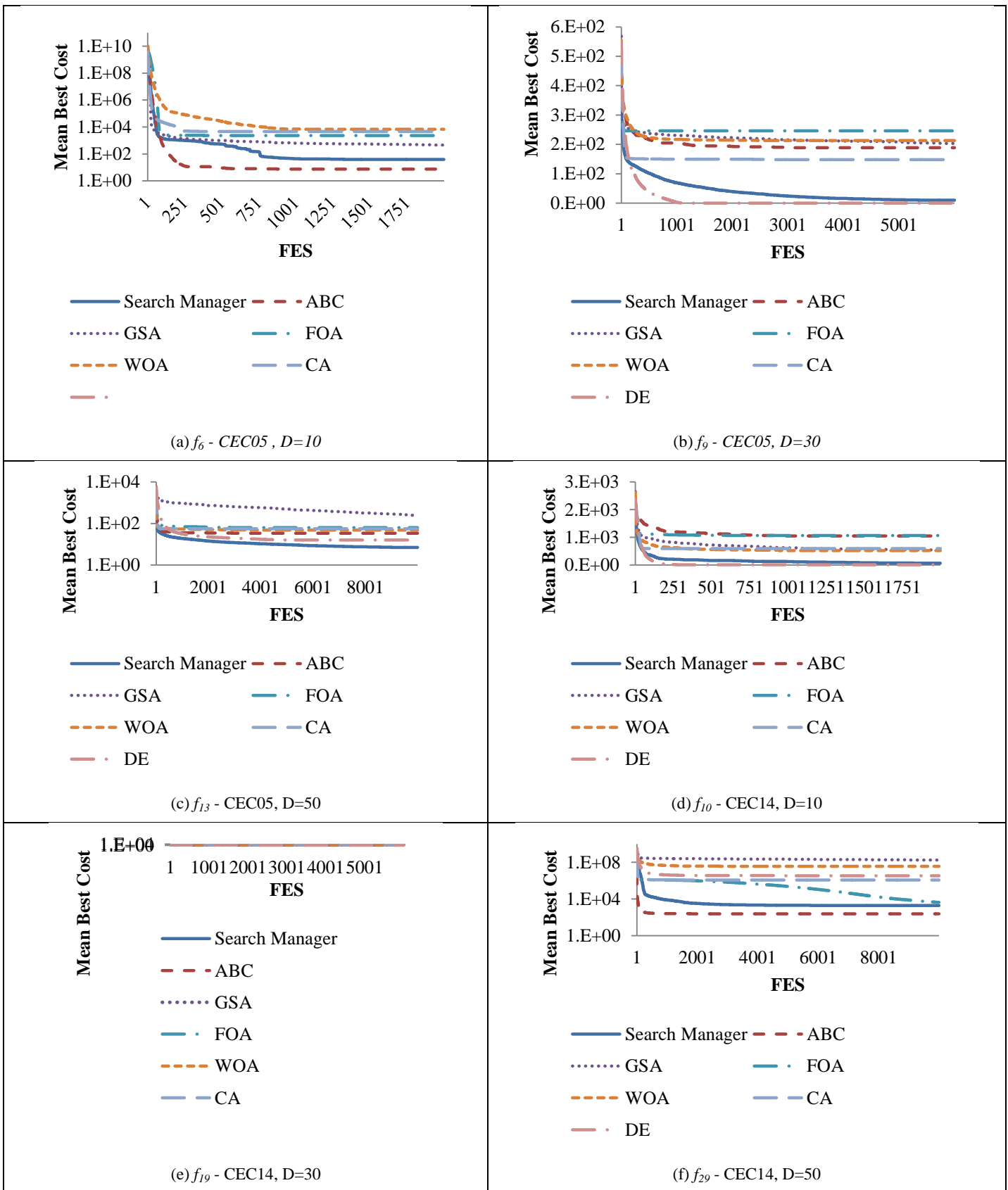


Fig. 9. Convergence curves of the algorithms on multimodal benchmark functions.

A Study of Feature Selection Algorithms for Predicting Students Academic Performance

Maryam Zaffar

Department of Computer and Information Sciences
Universiti Teknologi PETRONAS,
32610 Seri Iskandar, Malaysia

Manzoor Ahmed Hashmani

High Performance Cloud Computing Center
Department of Computer and Information Sciences
Universiti Teknologi PETRONAS,
32610 Seri Iskandar, Malaysia

K.S. Savita

High Performance Cloud Computing Center
Department of Computer and Information Sciences
Universiti Teknologi PETRONAS,
32610 Seri Iskandar, Malaysia

Syed Sajjad Hussain Rizvi

Department of Telecommunication
Hamdard University, Karachi,
Pakistan

Abstract—The main aim of all the educational organizations is to improve the quality of education and elevate the academic performance of students. Educational Data Mining (EDM) is a growing research field which helps academic institutions to improve the performance of their students. The academic institutions are most often judged by the grades achieved by the students in examination. EDM offers different practices to predict the academic performance of students. In EDM, Feature Selection (FS) plays a vital role in improving the quality of prediction models for educational datasets. FS algorithms eliminate unrelated data from the educational repositories and hence increase the performance of classifier accuracy used in different EDM practices to support decision making for educational settings. The good quality of educational dataset can produce better results and hence the decisions based on such quality dataset can increase the quality of education by predicting the performance of students. In the light of this mentioned fact, it is necessary to choose a feature selection algorithm carefully. This paper presents an analysis of the performance of filter feature selection algorithms and classification algorithms on two different student datasets. The results obtained from different FS algorithms and classifiers on two student datasets with different number of features will also help researchers to find the best combinations of filter feature selection algorithms and classifiers. It is very necessary to put light on the relevancy of feature selection for student performance prediction, as the constructive educational strategies can be derived through the relevant set of features. The results of our study depict that there is a 10% difference of prediction accuracies between the results of datasets with different number of features.

Keywords—Educational data mining; feature selection algorithms; classifiers; CFS; relief feature selection algorithm

I. INTRODUCTION

Education is a prime factor for the development of a nation. The quality of education is one of the most needed ingredients in creating remarkable members of society. The data kept in academic institution databases plays noteworthy role for the improvement of educational process by exploring the hidden

information [1]. Many techniques are being used to evaluate the performance of students' academics. Data Mining techniques are being broadly used on student data these days [2], [3] and is playing a positive role in the area of Educational Data Mining (EDM). EDM discovers the educational data to comprehend the issues in student's academic performance using the fundamental nature of data mining techniques [4]. Student performance prediction is considered as an important topic in EDM. As the performance of student not only effect the organization reputation, but also the future of the student itself, therefore the student performance prediction models are in spot light in front of educational stakeholders. EDM deploys data to help academic organizations in planning educational strategies and in turn enhancing the quality of education.

Student academic progress can be monitored through the prediction models. These prediction models use different EDM techniques to analyze the students' academic performance. It is very hard to distinguish the features affecting the student academic performance [5]. Student academic performance prediction can be helpful for institutions to identify students in need of financial assistance [6], [7], improve institution enrolment quality [7], [8], help students to plan better for future, and also to overcome their struggle with studies. The students' performance prediction model depends on the selected features from the dataset. The most suitable features can be selected by applying feature selection algorithm [9]. These algorithms can refine the prediction results [10]. However, the Feature selection algorithms are best to extract the relevant features and avoid redundancy, without cost of data loss [11], therefore it is very suitable to use FS algorithms in EDM to avoid loss of important data to build strategies with the help of such a quality data.

Feature Selection algorithms are used in in pre-processing step of data. It supports to select the appropriate subset of features to construct a model for data mining. However, Feature Selection algorithms are utilized to improve the predictive accuracy and lower the computational complexity [4], [12], [13]. The feature selection algorithms can increase

the performance of student performance prediction models. There are three main types of feature selection algorithms three main categories: filter, wrapper, and hybrid models. Filter method is performed on pre-processing step, and are not depended on any learning algorithm, but they depend on overall features of the training data. Wrapper method uses learning algorithms to estimate the features. Whereas Hybrid Feature selectin combines the properties of both filter and wrapper method [12]. In this study we focus mainly on the filter feature selection algorithm.

Feature selection has been used in EDM in different research works [5], [9], [14]. Researchers in EDM use different feature selection algorithms to yield effective results in predicting academic performance of students. But still a lot of attention is required to construct student performance prediction models with the help of feature selection algorithms. Our paper is a step towards detecting the best amalgamations of feature section algorithms and classification algorithms on student datasets.

The outline of the paper is as follows: Section II provides the literature related to the feature selection algorithm used in the field of EDM. Section III provides the research methodology followed by the paper. Section IV illustrates the results and discussions. Conclusion of the study is described in Section V.

II. RELATED LITERATURE

This section gives a brief literature review on the feature selection algorithms used in the field of EDM and the different combinations of feature selection along with classification algorithms used in the other studies. The study in [15] proposed an improved decision tree to predict the indicators of student dropouts. The study collects the dataset of 240 students through a survey and applies Correlation based Feature Selection (CFS) algorithm (Filter feature selection algorithm) in pre-processing step. The classification accuracy of the model shows more than 90%. However, the study took only one dataset into consideration. The investigation in [4] evaluated six feature selection algorithms to predict the performance of higher secondary students. The results of the study conclude that Voted Perceptron, and One Rule (OneR) shows high predictive performance with all the feature subsets gained through feature selection algorithms. Furthermore, Information Gain (IG) and CFS shows better ROC value and F-measure values on higher secondary school dataset.

A study to predict the performance of student in secondary school at Tuzla was presented in [1]. The study used Gain Ratio (GR) feature selection algorithm on the dataset with 19 features. The results with Random Forest classification (RF) algorithm reveals best results in terms of prediction accuracy.

The investigation in [16] was conducted to predict the enrolment of students in Science, Technology, Engineering and Mathematics (STEM) in higher educational institutions in Kenya. Almost 18 features were collected through a questionnaire. The CART decision tree shows better prediction

accuracy results with Chi-Square and IG feature selection algorithms.

A study to predict the grades of student was conducted in [14], Principal Component Analysis (PCA) was performed on the dataset of students enrolled in the computer science bachelor's degree . The study uses PCA to build decision trees from the features extracted through the Moodle Logs, to predict student grades.

A comparison between Greedy, IG-ratio ,Chi-Square and mRMR feature selections, was conducted in the study of [17]. The study collected first year students' record with 15 attributes, from the database of University of Technology, Thailand. The study proposed that Greedy Forward selection can give better prediction accuracy result with artificial neural network (ANN) as compared to Naïve Bayes, decision tree and k-NN.

The existing studies in educational data mining have used different filter feature selection algorithms on student datasets. In this study, we used two different datasets, with different number of features. This study is an extension of our previous work [18].

III. RESEARCH METHODOLOGY

This research article is an extended version of the paper [18]. One dataset was used in previous study to check the performance of different feature selection algorithms. The foremost objective of this research is to estimate the performance of different FS algorithms along with different classification algorithms using different students' datasets with dissimilar number of features. The comparison between the results of FS algorithms is based on two datasets to provide to new educational data mining for the performance of various feature selection algorithms with different number of features. This study will answer two research questions that are:

RQ1. What are the important feature selection algorithms to predict the academic performance of students (Whether they pass or fail)?

RQ2. What are the best possible combinations of feature selection algorithms and classification algorithms to predict the performance of students (Whether they pass or fail)?

To achieve the research objective and to answer the above-mentioned research questions, two student datasets are taken from valid sources, after which different FS algorithms are applied which was not used earlier on this dataset in the previous studies. As in this paper we try to evaluate different feature selection algorithms to check their performance. Various classification algorithms are applied using different FS algorithms. It is evaluated to check the performance among all the combinations applied on students' dataset. Fig. 1 describes a basic flow of our study. Two student datasets were taken in this study. In the second step feature selection algorithms are applied separately on both datasets, in combination of different classification algorithms. Results of precision and correctly classified instances were compared in the final step.

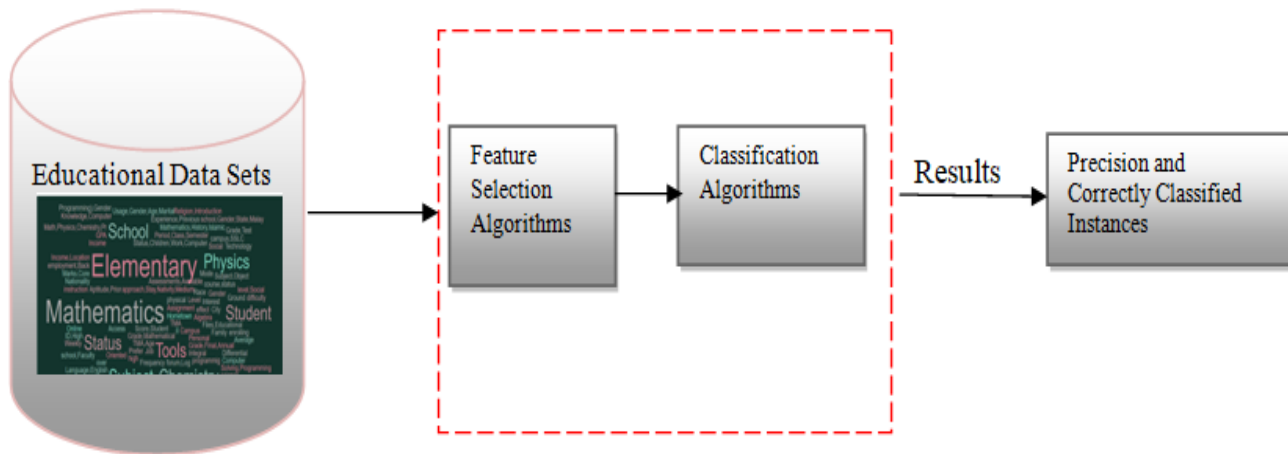


Fig. 1. Flow of methodology.

A. Dataset Description

In this study we have taken two student datasets with different number of features to check the performance of feature selection algorithm on different number of features. The details of two datasets used in this study are given below.

1) *Dataset 1*: The dataset 1 is comprised of 500 students records with 16 features. This dataset has been used in the study [19], and is available publicly even on Kaggle dataset repository. It is being previously used to check the learner's interactivity with e-learning management system. However, only information gain based feature selection algorithm is used previously. There are three categories of attributes in this dataset demographic, academic and behavioral. The dataset is being used by our previous version of this study.

2) *Dataset 2*: The dataset 2 is comprised of 300 students with 24 features. It was collected from the three different collages of India. This dataset is used in the study [20]. The dataset is being used in this paper to analyze the student's academic performance.

B. Experimental Setup

Waikato Environment for Knowledge Analysis (WEKA) is developed by University of Waikato in New Zealand as data mining tool. It is built in Java language, and a rich source of data mining algorithms. WEKA offers skill for developing machine learning techniques for different data mining tasks [21], [22]. In this experiment we have used Weka version 3.9, and explorer application.

C. Feature Selection Algorithm and Classifiers

Feature selection is one of the most recurrent and significant technique in data pre-processing and is said to be an essential element of machine learning process [23]. The main focus of our research in this paper is on six important FS algorithm CfsSubsetEval, ChiSquared-AttributeEval, FilteredAttributeEval, GainRatioAttribute-Eval, Principal Components, and ReliefAttributeEval feature selection algorithms are evaluated.

1) *CfsSubsetEval*: This approach identify the predictive capability of every feature. However, the redundancy factor also plays a critical role in this approach [24], [25]. CFS algorithm uses homogeneous feature in selection process along with discretization preprocessing steps [26].

2) *ChiSquaredAttributeEval*: Chi-Square is used to compare the tests of independence and the test of goodness of fit. Test of independence estimates whether the class label is dependent or independent of a feature. ChiSquared-AttributeEval estimates an attribute by calculating the value of the chi-squared statistic relating to the class [17], [25].

3) *FilteredAttributeEval*: This filter feature selection algorithm is available in Weka plate form.

4) *GainRatioAttributeEval*: The Gain Ratio is the non-symmetrical measure that is introduced to compensate for the bias of the information gain [27]. It is a filter feature selection algorithm that measures how common a feature in a class associated to all other classes.

5) *Principal Components*: Principal Component analysis reduces the dimensionality of space, without reducing the number of features [28].

6) *ReliefAttributeEval*: Relief is a simple weight-based algorithm which depends totally on a statistical method. It evaluates the significance of an attribute by sampling an instance repeatedly [25]. It detects those features which are statistically related to the target concept. It has a limitation of non-optimal feature set size [29].

Prediction accuracy of the features selected from the feature selection algorithms can be evaluated through classification algorithms. In our previous work we have used fifteen classification algorithms that are: Bayesian Network (BN), Naïve Bayes (NB), NaiveBayesUpdateable (NBU), MLP, Simple Logistic (SL), SMO, Decision Table (DT), OneR J rip, Decsion Stump (DS), J48, Random Forest (RF), RandomTree (RT), REpTree (RepT). However due the limitation of space we have selected six classification algorithms in this paper.

IV. RESULTS AND DISCUSSIONS

This research reported focuses on the performance evaluation of six Feature Selection algorithms using two different student’s datasets. The effectiveness of these algorithms is measured through Precision, Recall, F-measure and prediction accuracy (Correctly classified instances). F-measure is defined as the harmonic mean of precision and recall [30]. The results presented in our previous study [18] and which is then compared with the results obtained using dataset 1 and dataset 2. The outcomes of the six Feature Selection techniques using dataset 1 are reported in Tables I to VI by applying 15 classifiers. These tables illustrate results obtained by each of the Feature Selection (FS) algorithms. Furthermore, each table of results contains four columns that are FS-Classification Algorithm, Precision, Recall and F-measure values.

A. Results on Dataset 1

The results in Table I shows the different values of accuracy measures for fifteen classifiers with CfsSubsetEval feature selection algorithm using dataset 1. Fig. 2 graphically illustrates the results obtained with ChiSquared-AttributeEval feature selection algorithms. The results presented in Table II and Fig. 3 depicts that the classifier Decision Stump (DS) has the lowest performance on educational dataset 1 with ChiSquaredAttributeEval; however, MLP classifier shows comparatively better results than other classifiers with the same FS technique.

The results presented in Table III and Fig. 4 indicates that the accuracy of classifiers used on educational data with FilteredAttributeEval feature selection algorithm. The results demonstrate that the values of Precision, Recall and F-measure are comparatively low when Decision Stump and Jip classifiers are applied. While MLP performance is relatively improved than other classifiers using FilteredAttributeEval.

TABLE I. RESULTS OF CFSUBSETEVAL ON DATASET 1 USING DIFFERENT CLASSIFIERS [18]

FS-Classification Algorithm	Precision	Recall	F-Measure
Cfs-BN	0.724	0.743	0.742
Cfs-NB	0.73	0.729	0.728
Cfs-NBU	0.73	0.729	0.729
Cfs-MLP	0.736	0.729	0.729
Cfs-SL	0.724	0.722	0.723
Cfs-SMO	0.668	0.667	0.667
Cfs-DT	0.693	0.688	0.688
Cfs-Jrip	0.659	0.66	0.658
Cfs-OneR	0.611	0.583	0.571
Cfs-PART	0.713	0.708	0.71
Cfs-DS	0.373	0.528	0.437
Cfs-J48	0.708	0.701	0.702
Cfs-RF	0.64	0.632	0.633
Cfs-RT	0.627	0.618	0.621
Cfs-RepT	0.667	0.66	0.655

TABLE II. RESULTS OF CHISQUAREDATTRIBUTEVAL ON DATASET 1 USING DIFFERENT CLASSIFIERS [18]

FS-Classification Algorithm	Precision	Recall	F-Measure
Chi-BN	0.716	0.715	0.716
Chi-NB	0.66	0.66	0.654
Chi-NBU	0.66	0.66	0.654
Chi-MLP	0.769	0.764	0.764
Chi-SL	0.715	0.708	0.709
Chi-SMO	0.741	0.736	0.737
Chi-DT	0.71	0.701	0.702
Chi-Jrip	0.698	0.694	0.692
Chi-OneR	0.611	0.583	0.571
Chi-PART	0.64	0.639	0.639
Chi-DS	0.373	0.528	0.437
Chi-J48	0.709	0.708	0.708
Chi-RF	0.718	0.715	0.716
Chi-RT	0.674	0.674	0.674
Chi-RepT	0.651	0.653	0.651

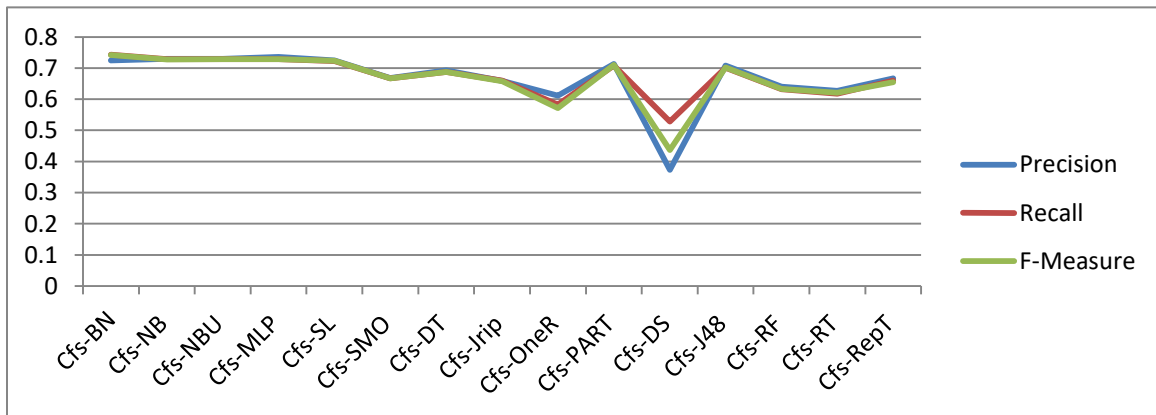


Fig. 2. Performance of CfsSubsetEval using Dataset 1.

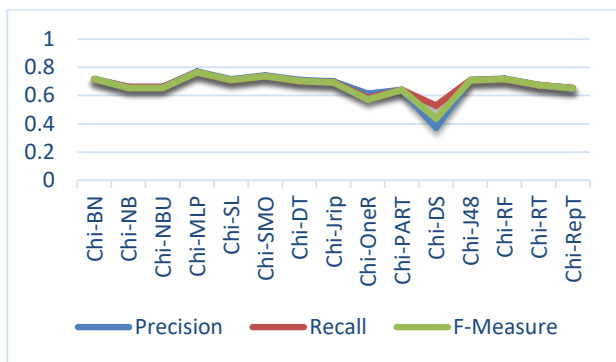


Fig. 3. Performance of ChiSquaredAttributeEval using Dataset 1.

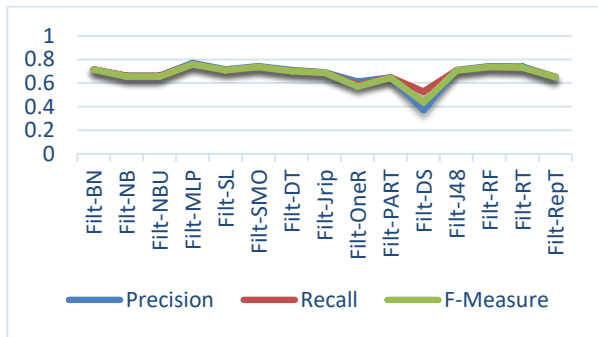


Fig. 4. Performance of FilteredAttributeEval using Dataset 1.

TABLE III. PERFORMANCE EVALUATION OF FILTERED ATTRIBUTE EVAL USING PRECISION RECALL AND F-MEASURE ON DATASET 1 [18]

FS-Classification Algorithm	Precision	Recall	F-Measure
Filt-BN	0.716	0.715	0.716
Filt-NB	0.66	0.66	0.654
Filt-NBU	0.66	0.66	0.654
Filt-MLP	0.768	0.757	0.758
Filt-SL	0.715	0.708	0.709
Filt-SMO	0.741	0.736	0.737
Filt-DT	0.71	0.701	0.702
Filt-Jrip	0.691	0.688	0.688
Filt-OneR	0.611	0.583	0.571
Filt-PART	0.646	0.646	0.645
Filt-DS	0.373	0.528	0.437
Filt-J48	0.709	0.708	0.707
Filt-RF	0.741	0.736	0.737
Filt-RT	0.738	0.729	0.73
Filt-RepT	0.651	0.653	0.651

The results reported in Table IV and Fig. 5 are exhibiting the identical performance details as illustrated earlier in Table III and Fig. 4. The results show that the decrease in performance by applying GainRatioAttributeEval Jrip classifier, however, MLP and SMO performed comparatively better than other classifiers.

The results in Table V present the performance of Principal Components using fifteen selected classification algorithms. Fig. 6 is the graphical representation of the performance of Principal Components. The result in Table V depicts that SMO classifier performed relatively better, while the performance of Jrip and Decision Stump classifiers is contradictory to the expected with Principal component.

TABLE IV. PERFORMANCE EVALUATION OF GAINRATIOATTRIBUTEVAL USING PRECISION RECALL AND F-MEASURE ON DATASET 1 [18]

FS-Classification Algorithm	Precision	Recall	F-Measure
GR-BN	0.716	0.715	0.716
GR-NB	0.66	0.66	0.654
GR-NBU	0.66	0.66	0.654
GR-MLP	0.768	0.757	0.758
GR-SL	0.715	0.708	0.709
GR-SMO	0.741	0.736	0.737
GR-DT	0.71	0.701	0.702
GR-Jrip	0.691	0.688	0.688
GR-OneR	0.611	0.583	0.571
GR-PART	0.646	0.646	0.645
GR-DS	0.373	0.528	0.437
GR-J48	0.709	0.708	0.707
GR-RF	0.741	0.736	0.737
GR-RT	0.738	0.729	0.73
GR-RepT	0.651	0.653	0.651

TABLE V. RESULTS OF PRINCIPALCOMPONENTS ON DATASET 1 USING DIFFERENT CLASSIFIERS [18]

FS-Classification Algorithm	Precision	Recall	F-Measure
PC-BN	0.643	0.632	0.633
PC-NB	0.508	0.507	0.506
PC-NBU	0.508	0.507	0.506
PC-MLP	0.694	0.694	0.693
PC-SL	0.692	0.688	0.688
PC-SMO	0.745	0.736	0.737
PC-DT	0.633	0.618	0.617
PC-Jrip	0.57	0.549	0.545
PC-OneR	0.445	0.444	0.445
PC-PART	0.591	0.59	0.591
PC-DS	0.345	0.486	0.403
PC-J48	0.674	0.667	0.668
PC-RF	0.701	0.694	0.695
PC-RT	0.585	0.576	0.576
PC-RepT	0.659	0.66	0.659

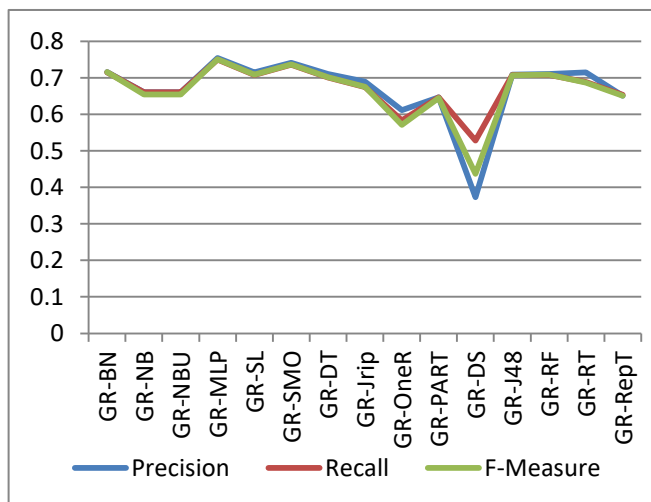


Fig. 5. Performance of GainRatioAttributeEval using dataset 1.

Table VI and Fig. 7 presents the result of ReliefAttributeEval(Rel) using different classifiers. It is observed through the results analysis that Random Forest classifiers shows better results with ReliefAttributeEval, however, the Decision Stump (DS) classifier depicts poor performance with ReliefAttributeEval using data set1 of students records.

B. Comparison of Results on Dataset 1 and Dataset 2

The comparison between the correctly classified instances using dataset 1 and dataset 2 are illustrated in Table VII. In this table six classifiers are presented only which performed better as compared to the other classifiers. The results indicate significant difference in the performance using both the datasets. There is approximately 10 to 20% performance and accuracy difference with each of the FS algorithm.

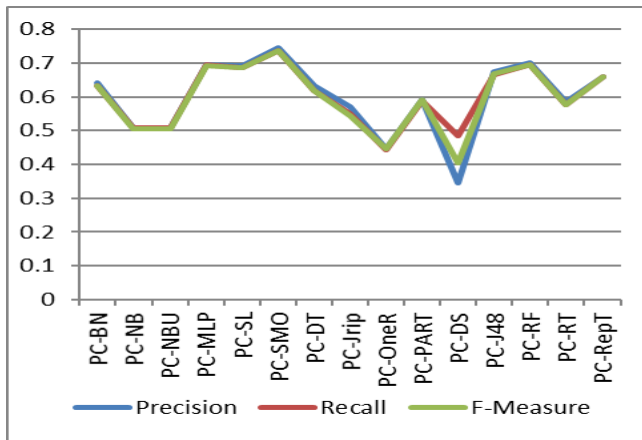


Fig. 6. Precision, recall and F-measure of principal components.

TABLE VI. RESULTS OF RELIEF ATTRIBUTE ON DATASET 1 USING DIFFERENT CLASSIFIERS [18]

FS-Classification Algorithm	Precision	Recall	F-Measure
Rel-BN	0.716	0.715	0.716
Rel-NB	0.66	0.66	0.654
Rel-NBU	0.66	0.66	0.654
Rel-MLP	0.767	0.764	0.764
Rel-SL	0.715	0.708	0.709
Rel-SMO	0.741	0.736	0.737
Rel-DT	0.71	0.701	0.702
Rel-Jip	0.713	0.708	0.708
Rel-OneR	0.611	0.583	0.571
Rel-PART	0.646	0.646	0.645
Rel-DS	0.373	0.528	0.437
Rel-J48	0.709	0.708	0.707
Rel-RF	0.756	0.75	0.873
Rel-RT	0.665	0.66	0.657
Rel-RepT	0.651	0.653	0.651

1) *Feature Selection Algorithms Accuracy:* Relief feature selection and Chi-Square algorithm with MLP classifier provides maximum accuracy using the dataset 1. While dataset 2 is used with chi feature selection technique in combination with Bayes Net (BN) classification algorithm offers the maximum accuracy. Principal component feature reduction technique in combination with Naïve Bayes (NB), provides least accuracy on dataset 1. Though other selected FS

techniques in combination with decision tree algorithm exhibits the least accuracy. Hence, the overall performance degrades for dataset 1 with the combination of FS technique and Decision Tree (DT) classifiers. Likewise, the Chi-square FS algorithm with Decision tree results in least performance on the dataset 2. It is concluded from the accuracy measures illustrated in Table VII that performance is better with 16 features of dataset 1 than the 24 features of dataset 2.

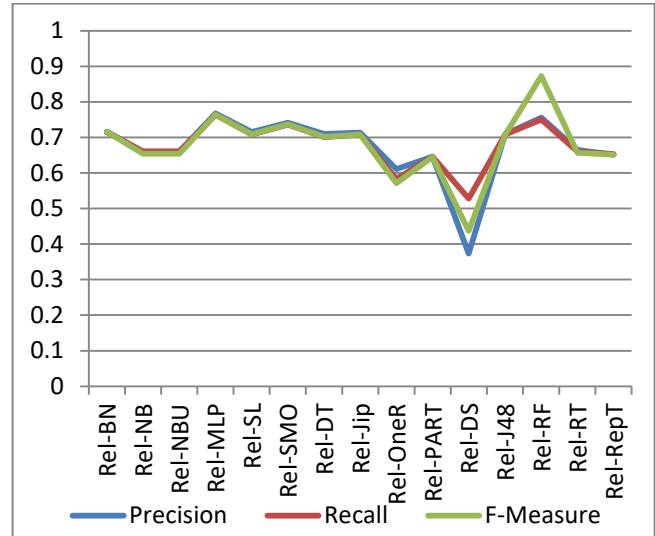


Fig. 7. Performance of ReliefAttributeEval using Dataset 1.

A comparative analysis based on the number of features selected in the dataset 1 and dataset 2 with respect to the precision is presented in Table VIII. The chi-square FS technique with Mlp classifiers results in maximum precision using the dataset 1 whereas Cfs algorithm along with the Bayes Net and Naïve Bayes provides maximum precision using the dataset 2. However, the performance of FS techniques with decision tree classification algorithm degrades using the dataset 1 and 2. The performance analysis discussed answer the two research questions discussed in Section III. These results give the answer of two research questions.

RQ1. What are the important feature selection techniques to predict the performance of students?

It is concluded from Tables VII and VIII that the performance of FS techniques has been improved using dataset 1 as compared to the dataset 2. Relief feature selection technique and Chi square algorithm perform better on dataset 1. Whereas Chi square and Cfs feature selection techniques perform better on dataset 2. Hence, these techniques must be considered in predicting the performance of students. According to the analysis, relief, chi-square and cfs are important FS techniques to predict the performance of student.

RQ2. What are the best possible combinations of feature selection techniques and Classification algorithms to predict the performance of students?

Fig. 8 and 9 shows that there is an evident difference in the results of dataset 1 and dataset 2. The results with the dataset 1 are much better than the results with the dataset 2. Both the figures are presenting a clear picture of the results.

TABLE VII. PERFORMANCE EVALUATION OF FEATURE SELECTION ALGORITHMS ON DATASET 1 & 2 IN CONTEXT WITH % OF CORRECTLY CLASSIFIED INSTANCES

FS-Classification Technique	Data set1	Dataset2
Cfs-BN	0.724	0.625
Cfs-NB	0.73	0.625
Cfs-MLP	0.736	0.561
Cfs-SMO	0.668	0.523
Cfs-DS	0.373	0.287
Cfs-RF	0.64	0.614
Chi-BN	0.716	0.616
Chi-NB	0.66	0.597
Chi-MLP	0.769	0.441
Chi-SMO	0.741	0.548
Chi-DS	0.373	0.367
Chi-RF	0.718	0.452
Filt-BN	0.716	0.61
Filt-NB	0.66	0.614
Filt-MLP	0.768	0.496
Filt-SMO	0.741	0.534
Filt-DS	0.373	0.287
Filt-RF	0.741	0.438
GR-BN	0.716	0.559
GR-NB	0.66	0.555
GR-MLP	0.754	0.506
GR-SMO	0.741	0.519
GR-DS	0.373	0.287
GR-RF	0.71	0.565
PC-BN	0.643	0.367
PC-NB	0.508	0.488
PC-MLP	0.694	0.436
PC-SMO	0.745	0.495
PC-DS	0.345	0.28
PC-RF	0.701	0.363
Rel-BN	0.716	0.58
Rel-NB	0.66	0.596
Rel-MLP	0.767	0.439
Rel-SMO	0.741	0.444
Rel-DS	0.373	0.287
Rel-RF	0.756	0.499

TABLE VIII. PERFORMANCE EVALUATION OF FEATURE SELECTION ALGORITHMS ON DATASET 1 & 2 IN CONTEXT WITH % OF CORRECTLY CLASSIFIED INSTANCES

FS-Classification Technique	Dataset1	Dataset2
Cfs-BN	74.31	57.84
Cfs-NB	72.08	55.88
Cfs-MLP	72.92	57.84
Cfs-SMO	66.67	55.88
Cfs-DS	52.78	42.51
Cfs-RF	63.19	59.8
Chi-BN	71.52	61.33
Chi-NB	65.97	59.33
Chi-MLP	76.39	44.33
Chi-SMO	73.61	55
Chi-DS	52.78	42
Chi-RF	71.53	45.33
Filt-BN	71.53	59.8
Filt-NB	65.97	59.8
Filt-MLP	75.69	48.03
Filt-SMO	73.61	51.96
Filt-DS	52.78	42.15
Filt-RF	73.61	42.15
GR-BN	71.53	56.33
GR-NB	65.97	55.66
GR-MLP	75	51
GR-SMO	65.97	54.3
GR-DS	52.78	42.15
GR-RF	70.83	55.88
PC-BN	63.19	45.09
PC-NB	50.69	51.96
PC-MLP	69.44	45.09
PC-SMO	73.61	49.01
PC-DS	48.61	43.13
PC-RF	69.44	47.05
Rel-BN	71.53	55.88
Rel-NB	65.97	53.92
Rel-MLP	76.39	46.07
Rel-SMO	73.61	48.03
Rel-DS	52.78	42.15

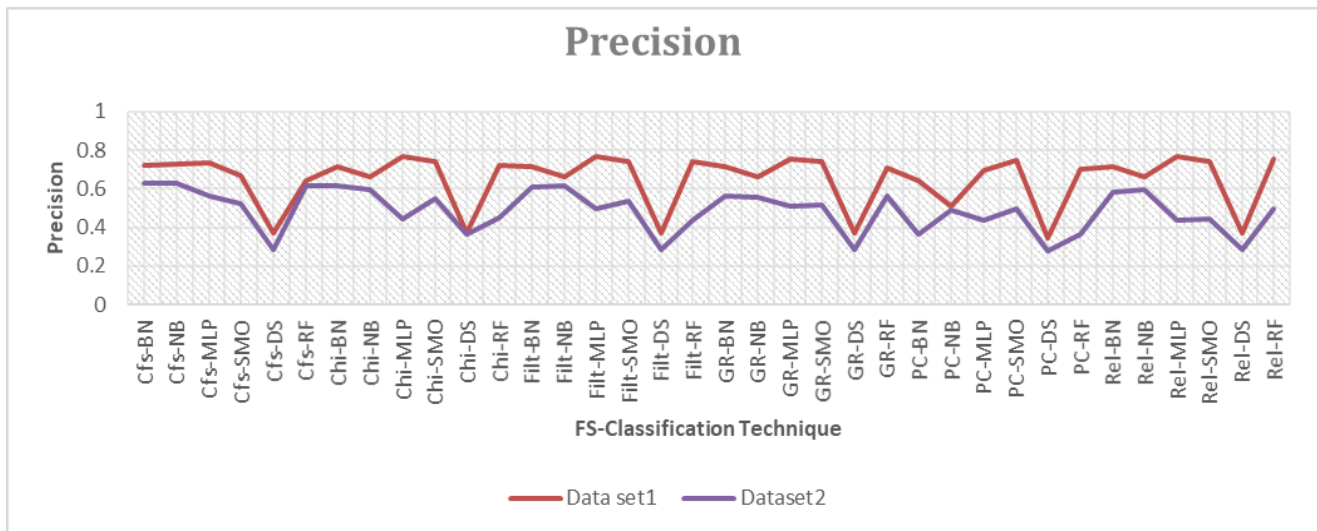


Fig. 8. Comparison of precision accuracy using dataset 1 & 2.

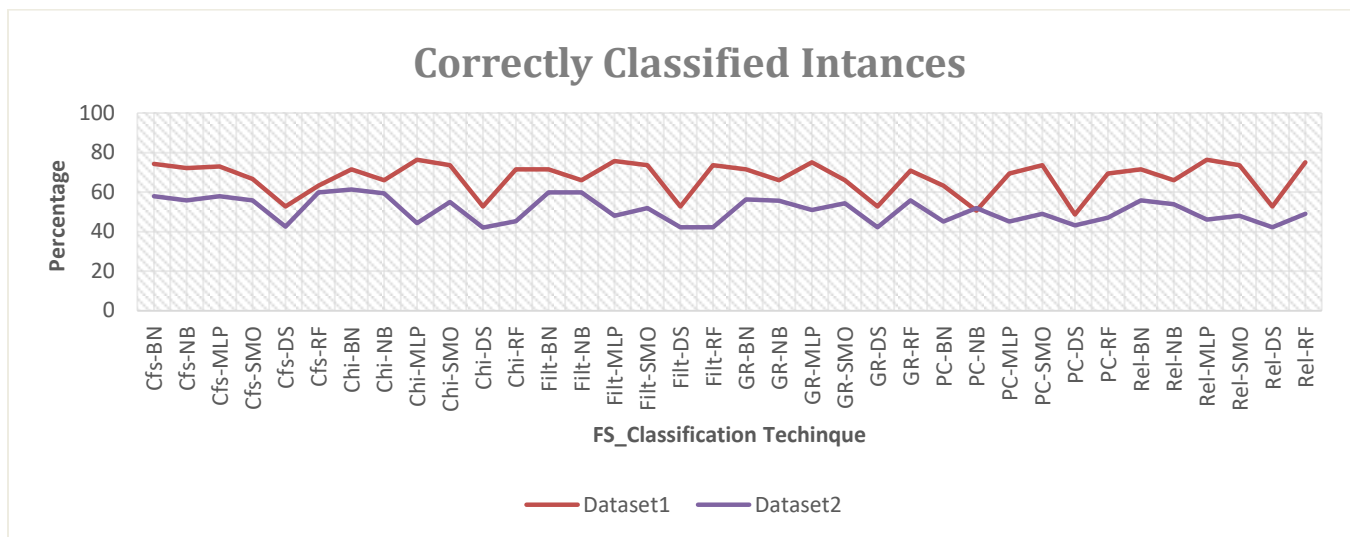


Fig. 9. Comparison of correctly classified instances using dataset 1 & 2.

V. CONCLUSION

This paper presents the study of various feature selection algorithms and analysed their performance using two different datasets. The results indicated that there is significant performance difference of feature selection algorithms using the datasets with different numbers of features; shows 10 to 20 per cent difference in accuracy percentages. The performance of the filter feature selection techniques reduces as the number of feature increases. To predict the academic performance of the student, having a large number of feature sets, wrappers feature selection techniques can also be evaluated. In future we will also evaluate the feature selection results through confusion m. Furthermore, we cannot neglect the advantages of filter feature selection techniques. In future, the study can be enhanced by applying few hybrid feature selection algorithms on student datasets in order to predict the performance of the student.

REFERENCES

- [1] E. Osmanbegović, M. Suljić, and H. Agić, "DETERMINING DOMINANT FACTOR FOR STUDENTS PERFORMANCE PREDICTION BY USING DATA MINING CLASSIFICATION ALGORITHMS," *Tranzicija*, vol. 16, pp. 147-158, 2015.
- [2] A. M. Shahiri and W. Husain, "A review on predicting student's performance using data mining techniques," *Procedia Computer Science*, vol. 72, pp. 414-422, 2015.
- [3] C. Romero and S. Ventura, "Educational data mining: a review of the state of the art," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 40, pp. 601-618, 2010.
- [4] M. Ramaswami and R. Bhaskaran, "A study on feature selection techniques in educational data mining," *arXiv preprint arXiv:0912.3924*, 2009.
- [5] A. Mueen, B. Zafar, and U. Manzoor, "Modeling and Predicting Students' Academic Performance Using Data Mining Techniques," *International Journal of Modern Education and Computer Science*, vol. 8, p. 36, 2016.
- [6] M. Ramaswami and R. Rathinasabapathy, "Student Performance Prediction," *International Journal of Computational Intelligence and Informatics*, vol. 1, 2012.
- [7] N. T. Nghe, P. Janeczek, and P. Haddawy, "A comparative analysis of techniques for predicting academic performance," in *Frontiers In Education Conference-Global Engineering: Knowledge Without Borders, Opportunities Without Passports*, 2007. FIE'07. 37th Annual, 2007, pp. T2G-7-T2G-12.
- [8] P. Golding and O. Donaldson, "Predicting academic performance," in *Frontiers in education conference*, 36th Annual, 2006, pp. 21-26.
- [9] H. M. Harb and M. A. Moustafa, "Selecting optimal subset of features for student performance model," *Int J Comput Sci*, p. 5, 2012.
- [10] M. Doshi, "Correlation Based Feature Selection (Cfs) Technique To Predict Student Performance," *International Journal of Computer Networks & Communications*, vol. 6, p. 197, 2014.
- [11] W. Punlumjeak and N. Rachburee, "A comparative study of feature selection techniques for classify student performance," in *Information Technology and Electrical Engineering (ICITEE)*, 2015 7th International Conference on, 2015, pp. 425-429.
- [12] D. Koller and M. Sahami, "Toward optimal feature selection," *Stanford InfoLab1996*.
- [13] P. Mitra, C. Murthy, and S. K. Pal, "Unsupervised feature selection using feature similarity," *IEEE transactions on pattern analysis and machine intelligence*, vol. 24, pp. 301-312, 2002.
- [14] A. Figueira, "Predicting Grades by Principal Component Analysis: A Data Mining Approach to Learning Analytics," in *Advanced Learning Technologies (ICALT)*, 2016 IEEE 16th International Conference on, 2016, pp. 465-467.
- [15] S. Sivakumar, S. Venkataraman, and R. Selvaraj, "Predictive Modeling of Student Dropout Indicators in Educational Data Mining using Improved Decision Tree," *Indian Journal of Science and Technology*, vol. 9, 2016.
- [16] K. W. Stephen, "Data Mining Model for Predicting Student Enrolment in STEM Courses in Higher Education Institutions," 2016.
- [17] N. Rachburee and W. Punlumjeak, "A comparison of feature selection approach between greedy, IG-ratio, Chi-square, and mRMR in educational mining," in *Information Technology and Electrical Engineering (ICITEE)*, 2015 7th International Conference on, 2015, pp. 420-424.
- [18] M. Zaffar, M. A. Hashmani, and K. Savita, "Performance analysis of feature selection algorithm for educational data mining," in *Big Data and Analytics (ICBDA)*, 2017 IEEE Conference on, 2017, pp. 7-12.
- [19] E. A. Amrieh, T. Hamtini, and I. Aljarah, "Mining Educational Data to Predict Student's academic Performance using Ensemble Methods," *International Journal of Database Theory and Application*, vol. 9, pp. 119-136, 2016.
- [20] S. Hussain, N. A. Dahan, F. M. Ba-Alwi, and N. RIBATA, "Educational Data Mining and Analysis of Students' Academic Performance Using WEKA," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 9, 2018.
- [21] K. Patel, J. Vala, and J. Pandya, "Comparison of various classification algorithms on iris datasets using WEKA," *Int. J. Adv. Eng. Res. Dev.(IAERD)*, vol. 1, 2014.

- [22] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA data mining software: an update," ACM SIGKDD explorations newsletter, vol. 11, pp. 10-18, 2009.
- [23] A. Kalousis, J. Prados, and M. Hilario, "Stability of feature selection algorithms: a study on high-dimensional spaces," Knowledge and information systems, vol. 12, pp. 95-116, 2007.
- [24] M. A. Hall and L. A. Smith, "Practical feature subset selection for machine learning," 1998.
- [25] C. Anuradha and T. Velmurugan, "Feature Selection Techniques to Analyse Student Academic Performance using Naïve Bayes Classifier," in The 3rd International Conference on Small & Medium Business, 2016, pp. 345-350.
- [26] C. Huertas and R. Juárez-Ramírez, "Filter feature selection performance comparison in high-dimensional data: A theoretical and empirical analysis of most popular algorithms," in Information Fusion (FUSION), 2014 17th International Conference on, 2014, pp. 1-8.
- [27] J. Novaković, "Toward optimal feature selection using ranking methods and classification algorithms," Yugoslav Journal of Operations Research, vol. 21, 2016.
- [28] Q. Guo, W. Wu, D. Massart, C. Boucon, and S. De Jong, "Feature selection in principal component analysis of analytical data," Chemometrics and Intelligent Laboratory Systems, vol. 61, pp. 123-132, 2002.
- [29] K. Kira and L. A. Rendell, "The feature selection problem: Traditional methods and a new algorithm," in Aaai, 1992, pp. 129-134.
- [30] T. Velmurugan and C. Anuradha, "Performance Evaluation of Feature Selection Algorithms in Educational Data Mining," Performance Evaluation, vol. 5, 2016.