# INTERNATIONAL JOURNAL OF
# ADVANCED COMPUTER SCIENCE AND APPLICATIONS

# Editorial Preface

## From the Desk of Managing Editor...

It is our pleasure to present to you the February 2014 Issue of International Journal of Advanced Computer Science and Applications.

Today, it is incredible to consider that in 1969 men landed on the moon using a computer with a 32-kilobyte memory that was only programmable by the use of punch cards. In 1973, Astronaut Alan Shepherd participated in the first computer "hack" while orbiting the moon in his landing vehicle, as two programmers back on Earth attempted to "hack" into the duplicate computer, to find a way for Shepherd to convince his computer that a catastrophe requiring a mission abort was not happening; the successful hack took 45 minutes to accomplish, and Shepherd went on to hit his golf ball on the moon. Today, the average computer sitting on the desk of a suburban home office has more computing power than the entire U.S. space program that put humans on another world!!

Computer science has affected the human condition in many radical ways. Throughout its history, its developers have striven to make calculation and computation easier, as well as to offer new means by which the other sciences can be advanced. Modern massively-paralleled super-computers help scientists with previously unfeasible problems such as fluid dynamics, complex function convergence, finite element analysis and real-time weather dynamics.

At IJACSA we believe in spreading the subject knowledge with effectiveness in all classes of audience. Nevertheless, the promise of increased engagement requires that we consider how this might be accomplished, delivering up-to-date and authoritative coverage of advanced computer science and applications.

Throughout our archives, new ideas and technologies have been welcomed, carefully critiqued, and discarded or accepted by qualified reviewers and associate editors. Our efforts to improve the quality of the articles published and expand their reach to the interested audience will continue, and these efforts will require critical minds and careful consideration to assess the quality, relevance, and readability of individual articles.

To summarise, the journal has offered its readership thought provoking theoretical, philosophical, and empirical ideas from some of the finest minds worldwide. We thank all our readers for their continued support and goodwill for IJACSA. We will keep you posted on updates about the new programmes launched in collaboration.

Lastly, 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.

We hope that materials contained in this volume will satisfy your expectations and entice you to submit your own contributions in upcoming issues of IJACSA

**Thank you for Sharing Wisdom!**

# Editorial Board

# Reviewer Board Members

Chhattisgarh Swami Vivekananda Technical University

- **Chao Wang**

- **Chi-Hua Chen**

  National Chiao-Tung University

- **Ciprian Dobre**

  University Politehnica of Bucharest

- **Chien-Pheg Ho**

  Information and Communications Research Laboratories, Industrial Technology Research Institute of Taiwan

- **Charlie Obimbo**

  University of Guelph

- **Chao-Tung Yang**

  Department of Computer Science, Tunghai University

- **Dana PETCU**

  West University of Timisoara

- **Deepak Garg**

  Thapar University

- **Dewi Nasien**

  Universiti Teknologi Malaysia

- **Dheyaa Kadhim**

  University of Baghdad

- **Dong-Han Ham**

  Chonnam National University

- **Dragana Becejski-Vujaklija**

  University of Belgrade, Faculty of organizational sciences

- **Driss EL OUADGHIRI**

- **Duck Hee Lee**

  Medical Engineering R&D Center/Asan Institute for Life Sciences/Asan Medical Center

- **Dr. Santosh Kumar**

  Graphic Era University, Dehradun (UK)

- **Elena Camossi**

  Joint Research Centre

- **Eui Lee**

- **Elena SCUTELNICU**

  "Dunarea de Jos" University of Galati

- **Firkhan Ali Hamid Ali**

  UTHM

- **Fokrul Alom Mazarbhuiya**

  King Khalid University

- **Frank Ibikunle**

  Covenant University

- **Fu-Chien Kao**

  Da-Y eh University

- **Faris Al-Salem**

  GCET

- **gamil Abdel Azim**

  Associate prof - Suez Canal University

- **Ganesh Sahoo**

  RMRIMS

- **Gaurav Kumar**

  Manav Bharti University, Solan Himachal Pradesh

- **Ghalem Belalem**

  University of Oran (Es Senia)

- **Giri Babu**

  Indian Space Research Organisation

- **Giacomo Veneri**

  University of Siena

- **Giri Babu**

  Indian Space Research Organisation

- **Gerard Dumancas**

  Oklahoma Medical Research Foundation

- **Georgios Galatas**

- **George Mastorakis**

  Technological Educational Institute of Crete

- **Gunaseelan Devaraj**

  Jazan University, Kingdom of Saudi Arabia

- **Gavril Grebenisan**

  University of Oradea

- **Hadj Tadjine**

  IAV GmbH

- **Hamid Mukhtar**

  National University of Sciences and Technology

- **Hamid Alinejad-Rokny**

  University of Newcastle

- **Harco Leslie Hendric Spits Warnars**

  Budi LUhur University

- **Harish Garg**

  Thapar University Patiala

- **Hamez l. El Shekh Ahmed**

  Pure mathematics

- **Hesham Ibrahim**

  Chemical Engineering Department, Faculty of Engineering, Al-Mergheb University

- **Dr. Himanshu Aggarwal**

  Punjabi University, India

- **Huda K. AL-Jobori**

  Ahlia University

- **Iwan Setyawan**
  Satya Wacana Christian University
- **Dr. Jamaiah Haji Yahaya**
  Northern University of Malaysia (UUM), Malaysia
- **James Coleman**
  Edge Hill University
- **Jim Wang**
  The State University of New York at Buffalo, Buffalo, NY
- **John Salin**
  George Washington University
- **Jyoti Chaudary**
  High performance computing research lab
- **Jatinderkumar R. Saini**
  S.P.College of Engineering, Gujarat
- **K Ramani**
  K.S.Rangasamy College of Technology, Tiruchengode
- **K V.L.N.Acharyulu**
  Bapatla Engineering college
- **Kashif Nisar**
  Universiti Utara Malaysia
- **Kayhan Zrar Ghafoor**
  University Technology Malaysia
- **Kitimaporn Choochote**
  Prince of Songkla University, Phuket Campus
- **Kunal Patel**
  Ingenuity Systems, USA
- **Krasimir Yordzhev**
  South-West University, Faculty of Mathematics and Natural Sciences, Blagoevgrad, Bulgaria
- **Krassen Stefanov**
  Professor at Sofia University St. Kliment Ohridski
- **Labib Francis Gergis**
  Misr Academy for Engineering and Technology
- **Lai Khin Wee**
  Biomedical Engineering Department, University Malaya
- **Lazar Stosic**
  Collegefor professional studies educators Aleksinac, Serbia
- **Lijian Sun**
  Chinese Academy of Surveying and Mapping, China
- **Leandors Maglaras**
- **Leon Abdillah**
  Bina Darma University
- **Ljubomir Jerinic**

- University of Novi Sad, Faculty of Sciences, Department of Mathematics and Computer Science
- **Lokesh Sharma**
  Indian Council of Medical Research
- **Long Chen**
  Qualcomm Incorporated
- **M. Reza Mashinchi**
- **M. Tariq Banday**
  University of Kashmir
- **MAMTA BAHETI**
  SNJBS KBJ COLLEGE OF ENGINEERING, CHANDWAD, NASHIK, M.S. INDIA
- **Mazin Al-Hakeem**
  Research and Development Directorate - Iraqi Ministry of Higher Education and Research
- **Md Rana**
  University of Sydney
- **Miriampally Venkata Raghavendera**
  Adama Science & Technology University, Ethiopia
- **Mirjana Popvic**
  School of Electrical Engineering, Belgrade University
- **Manas deep**
  Masters in Cyber Law & Information Security
- **Manpreet Singh Manna**
  SLIET University, Govt. of India
- **Manuj Darbari**
  BBD University
- **Md. Zia Ur Rahman**
  Narasaraopeta Engg. College, Narasaraopeta
- **Messaouda AZZOUZI**
  Ziane AChour University of Djelfa
- **Dr. Michael Watts**
  University of Adelaide
- **Milena Bogdanovic**
  University of Nis, Teacher Training Faculty in Vranje
- **Miroslav Baca**
  University of Zagreb, Faculty of organization and informatics / Center for biomet
- **Mohamed Ali Mahjoub**
  Preparatory Institute of Engineer of Monastir
- **Mohamed El-Sayed**
  Faculty of Science, Fayoum University, Egypt
- **Mohammad Yamin**
- **Mohammad Ali Badamchizadeh**
  University of Tabriz

(v)

- **Mohamed Najeh Lakhoua**

  ESTI, University of Carthage

- **Mohammad Alomari**

  Applied Science University

- **Mohammad Kaiser**

  Institute of Information Technology

- **Mohammed Al-Shabi**

  Assistant Prof.

- **Mohammed Sadgal**

- **Mourad Amad**

  Laboratory LAMOS, Bejaia University

- **Mohammed Ali Hussain**

  Sri Sai Madhavi Institute of Science & Technology

- **Mohd Helmy Abd Wahab**

  Universiti Tun Hussein Onn Malaysia

- **Mueen Uddin**

  Universiti Teknologi Malaysia UTM

- **Mona Elshinawy**

  Howard University

- **Maria-Angeles Grado-Caffaro**

  Scientific Consultant

- **Mehdi Bahrami**

  University of California, Merced

- **Miriampally Venkata Raghavendra**

  Adama Science & Technology University, Ethiopia

- **Murthy Dasika**

  SreeNidhi Institute of Science and Technology

- **Mostafa Ezziyyani**

  FSTT

- **Marcellin Julius Nkenlifack**

  University of Dschang

- **Natarajan Subramanyam**

  PES Institute of Technology

- **Noura Aknin**

  University Abdelamlek Essaadi

- **Nidhi Arora**

  M.C.A. Institute, Ganpat University

- **Nazeeruddin Mohammad**

  Prince Mohammad Bin Fahd University

- **Najib Kofahi**

  Yarmouk University

- **NEERAJ SHUKLA**

  ITM UNiversity, Gurgaon, (Haryana) Inida

- **N.Ch. Iyengar**

  VIT University

- **Om Sangwan**

- **Oliviu Matel**

  Technical University of Cluj-Napoca

- **Osama Omer**

  Aswan University

- **Ousmane Thiare**

  Associate Professor University Gaston Berger of Saint-Louis SENEGAL

- **Omaima Al-Allaf**

  Assistant Professor

- **Paresh V Virparia**

  Sardar Patel University

- **Dr. Poonam Garg**

  Institute of Management Technology, Ghaziabad

- **Professor Ajantha Herath**

- **Prabhat K Mahanti**

  UNIVERSITY OF NEW BRUNSWICK

- **Qufeng Qiao**

  University of Virginia

- **Rachid Saadane**

  EE departement EHTP

- **raed Kanaan**

  Amman Arab University

- **Raja boddu**

  LENORA COLLEGE OF ENGINEERNG

- **Ravisankar Hari**

  SENIOR SCIENTIST, CTRI, RAJAHMUNDRY

- **Raghuraj Singh**

- **Rajesh Kumar**

  National University of Singapore

- **Rakesh Balabantaray**

  IIIT Bhubaneswar

- **RashadAl-Jawfi**

  Ibb university

- **Rashid Sheikh**

  Shri Venkteshwar Institute of Technology , Indore

- **Ravi Prakash**

  University of Mumbai

- **Rawya Rizk**

  Port Said University

- **Reshmy Krishnan**

  Muscat College affiliated to stirling University.U

- **Ricardo Vardasca**

  Faculty of Engineering of University of Porto

- **Ritaban Dutta**

  ISSL, CSIRO, Tasmaniia, Australia

- **Rowayda Sadek**

- **Ruchika Malhotra**

  Delhi Technoogical University

- **Saadi Slami**
  University of Djelfa

- **Sachin Kumar Agrawal**
  University of Limerick

- **Dr.Sagarmay Deb**
  University Lecturer, Central Queensland University, Australia

- **Said Ghoniemy**
  Taif University

- **Sasan Adibi**
  Research In Motion (RIM)

- **Sérgio Ferreira**
  School of Education and Psychology, Portuguese Catholic University

- **Sebastian Marius Rosu**
  Special Telecommunications Service

- **Selem charfi**
  University of Valenciennes and Hainaut Cambresis, France.

- **Seema Shah**
  Vidyalankar Institute of Technology Mumbai,

- **Sengottuvelan P**
  Anna University, Chennai

- **Senol Piskin**
  Istanbul Technical University, Informatics Institute

- **Seyed Hamidreza Mohades Kasaei**
  University of Isfahan

- **Shafiqul Abidin**
  G GS I P University

- **Shahanawaj Ahamad**
  The University of Al-Kharj

- **Shawkl Al-Dubaee**
  Assistant Professor

- **Shriram Vasudevan**
  Amrita University

- **Sherif Hussain**
  Mansoura University

- **Siddhartha Jonnalagadda**
  Mayo Clinic

- **Sivakumar Poruran**
  SKP ENGINEERING COLLEGE

- **Sim-Hui Tee**
  Multimedia University

- **Simon Ewedafe**
  Baze University

- **SUKUMAR SENTHILKUMAR**
  Universiti Sains Malaysia

- **Slim Ben Saoud**

- **Sudarson Jena**
  GITAM University, Hyderabad

- **Sumit Goyal**

- **Sumazly Sulaiman**
  Institute of Space Science (ANGKASA), Universiti Kebangsaan Malaysia

- **Sohail Jabb**
  Bahria University

- **Suhas J Manangi**
  Microsoft

- **Suresh Sankaranarayanan**
  Institut Teknologi Brunei

- **Susarla Sastry**
  J.N.T.U., Kakinada

- **Syed Ali**
  SMI University Karachi Pakistan

- **T C. Manjunath**
  HKBK College of Engg

- **T V Narayana Rao**
  Hyderabad Institute of Technology and Management

- **T. V. Prasad**
  Lingaya's University

- **Taiwo Ayodele**
  Infonetmedia/University of Portsmouth

- **Tarek Gharib**

- **THABET SLIMANI**
  College of Computer Science and Information Technology

- **Totok R. Biyanto**
  Engineering Physics, ITS Surabaya

- **TOUATI YOUCEF**
  Computer sce Lab LIASD - University of Paris 8

- **VINAYAK BAIRAGI**
  Sinhgad Academy of engineering, Pune

- **VISHNU MISHRA**
  SVNIT, Surat

- **Vitus S.W. Lam**
  The University of Hong Kong

- **Vuda SREENIVASARAO**
  School of Computing and Electrical Engineering,BAHIR DAR UNIVERSITY, BAHIR DAR,ETHIOPA

- **Vaka MOHAN**
  TRR COLLEGE OF ENGINEERING

- **Wei Wei**

- **Xiaojing Xiang**
  AT&T Labs
- **YASSER ATTIA ALBAGORY**
  College of Computers and Information Technology,
  Taif University, Saudi Arabia
- **YI FEI WANG**
  The University of British Columbia
- **Yilun Shang**
  University of Texas at San Antonio
- **YU QI**
  Mesh Capital LLC
- **Zacchaeus Omogbadegun**
  Covenant University

- **ZAIRI ISMAEL RIZMAN**
  UiTM (Terengganu) Dungun Campus
- **ZENZO POLITE NCUBE**
  North West University
- **ZHAO ZHANG**
  Deptment of EE, City University of Hong Kong
- **ZHIXIN CHEN**
  ILX Lightwave Corporation
- **ZLATKO STAPIC**
  University of Zagreb
- **Ziyue Xu**
- **ZURAINI ISMAIL**
  Universiti Teknologi Malaysia

# CONTENTS

# Applying Cellular Automata for Simulating and Assessing Urban Growth Scenario Based in Nairobi, Kenya

Kenneth Mubea,
Remote Sensing Research Group
(RSRG),
University of Bonn,
Bonn, Germany

Roland Goetzke,
Remote Sensing Research Group
(RSRG),
University of Bonn,
Bonn, Germany

Gunter Menz,
Center for Remote Sensing of Land
Surfaces (ZFL),
University of Bonn,
Bonn, Germany

*Abstract*—This research explores urban growth based scenarios for the city of Nairobi using a cellular automata urban growth model (UGM). African cities have experienced rapid urbanization over the last decade due to increased population growth and high economic activities. We used multi-temporal Landsat imageries for 1976, 1986, 2000 and 2010 to investigate urban land-use changes in Nairobi. Our UGM used data from urban land-use of 1986 and 2010, road data, slope data and exclusion layer. Monte-Carlo technique was used for model calibration and Multi Resolution Validation (MRV) technique for validation. Simulation of urban land-use was done up to the year 2030 when Kenya plans to attain Vision 2030. Three scenarios were explored in the urban modelling process; unmanaged growth with no restriction on environmental areas, managed growth with moderate protection, and a managed growth with maximum protection on forest, agricultural areas, and urban green. Thus alternative scenario development using UGM is useful for planning purposes so as to ensure sustainable development is achieved. UGM provides quantitative, visual, spatial and temporal information which aid policy and decision makers can make informed decisions.

*Keywords—Urban Growth; Scenarios; Nairobi; Cellular automata; Simulation; sustainable development*

## I. INTRODUCTION

Sustainable planning is crucial for future development of cities. Over the last decade there has been large rural urban migration in African cities as people search for employment and better amenities. This has led to a strain on the existing amenities and infrastructure [1]. There has been emergence of slums in cities in Africa such as Nairobi because of unsuitable land-use planning [2]. Undesirable consequences have been noted such as pollution, depletion of natural resources, inadequate transportation systems, urban sprawl among other negative environmental and social effects. Thus there is need for an integrated urban planning paradigm in order to identify and anticipate urban dynamics effectively.

Integration of remote sensing and urban growth modelling has been the frontier edge of urban research. Remote sensing provides spatially consistent data sets that cover large areas with both high spatial detail and high temporal frequency [1]. Such data sets are useful in land-use monitoring and simulation. As urbanisation occurs, changes in land-use increase thus taking up the natural resource base such as forests and agricultural land. This in turn leads to fragmentation and land degradation [3].

Models based on cellular automata (CA) have been used over the last decades in simulating urban development growth and patterns [4]. Early models were based on demographic trends and were not successful in simulating contemporary urban growth [5], [6], [7], and [8]. However, land-use modelling using CA utilise biophysical factors making it possible to simulate various patterns and intensities of urban growth [4]. Land-use change models have been used as decision support tools in urban planning in order to inform planners and decision makers [9]. For an urban model to be used in an area of interest it needs to be localised and this involves calibration. This is done in order to make it adapt to the endogenous characteristics of the particular environment for simulation [4]. Urban models aid in making informed decisions on land-use planning in the context of future development. Sustainable development is thus possible once various simulations of land-use scenarios have been obtained and this helps in understanding the consequences of different driving forces [10].

The "eXtendable Unified Land-use modelling platform" (XULU) was developed as a generic modelling framework at the University of Bonn, Germany [11]. It is able to handle several model types simultaneously such as statistic dynamic or agent based models of urban growth as well land-use change. We adopted the Urban Growth Model (UGM) on the modelling framework XULU for Nairobi. UGM was first developed and applied in the German federal state of North-Rhine Westphalia [12]. UGM is based on the modelling algorithm of the SLEUTH model [13] which uses the concept of cellular automata. Calibration of UGM involved five model parameters similar to SLEUTH model so as to make it adapt to Nairobi [4], [14].

Urban growth modelling based on cellular automata has been used mostly in cities in North America and European cities. Cities in Africa are different to the counterpart cities in the western world in various ways. Major cities in Africa are characterised by high rural urban migration which result on a strain on local urban transport systems, traffic congestion, development of informal settlements [1].

Fig. 1.   Location of Nairobi

Cities in Kenya represent very different environmental and geographic characteristics. Nairobi is the capital city of Kenya and has recently experienced a fast average annual growth rate of 4.9 per cent between the years 1990 and 2006 [15].

In this research, a cellular automata model was used to study land-use change and prediction of future trends in Nairobi as Kenya attains Vision 2030 [16]. The urban land-use data for Nairobi was derived from multi-spectral Landsat imagery captured in 1976, 1986, 2000 and 2010. At the end of the research, calibration and validation of both models were achieved. The models were used to predict the future urban land-use development in the year 2030.

## II.   THE STUDY AREA

Nairobi extends between latitudes 1° 09' and 1° 28' South, and longitude 36° 04' and 37° 10' East in Kenya, with an average altitude of 1,700 meters above sea level, covering an area of 696 km²  (Fig. 1).  Nairobi is the capital city of Kenya. The administratively defined town has land uses divided roughly into urban use, agriculture, rangeland, open/transitional areas, and remnants of evergreen tropical forests. Nairobi has a high growth rate per annum compared to other growth rates in Africa with 75 % of urban population living in informal settlements [17]. From a population of 310,000 in 1960, the population reached 510,000 in 1970 [18], 828,000 in 1979 [19], 1,321,000 in 1989 [20], 2,137,000 in 1999 [21] and 3,138,369 in 2009 [22]. The projected population in the year 2020 will be almost six million [17].

Urban sprawl has a negative impact on infrastructure and the sustainability of cities [15]. This is exhibited for instance in the increase of transport costs, public infrastructure of residential and commercial development. Most African cities show characteristic patterns of urban sprawl where urban development evolves around the nexus of the main transportation routes, with urban growth tending to grow in sectors emanating from city centers [1]. Many urban areas are faced with environmental problems like water pollution, uncontrolled waste disposal, bad air quality and noise.

## III.   MODELLING NAIROBI'S URBAN GROWTH

Urban growth is a complex process which involves the spatio-temporal changes of all socio-economic and physical components at different scales [23]. The process can be demonstrated in a simplified way and be analyzed empirically using urban growth models. Numerical simulation models for land-use change involve highly complex applications that have been developed to solve specific problems in urban areas. Consequently, a majority of these models have been developed at universities and are a result of long-time research [12]. [11] developed XULU (eXtendable Unified Land Use Modelling Platform), a modelling framework that enables model integration and carries out tasks using functionalities such as data storage, input/output methods, editing and visualization. XULU was first used to compute the future land-use for different scenarios with their specific boundary conditions for a watershed in Benin [24]. Hence, the CLUE-s land-use change

model, developed by [25] was implemented in the XULU modelling framework.

A majority of urban growth models are restricted to simulate changes of one land-use category, which is urban, and that typically just in one direction, which is growth [12]. SLEUTH is such a model and is an acronym for "Slope, Land use, Exclusion, Urban, Transport, Hill shade", as its main input parameters. SLEUTH is a cellular automaton (CA) based urban growth and land-use change model [13]. The model was initially applied in the United States of America but has also been applied in other regions of the world such as in Europe [26], South America [27] and Southeast Asia [28]. SLEUTH consists of two components, an urban growth model based on the Clarke Urban Growth Model (UGM) described in [13] and a so-called Land cover Deltatron Model [29] to simulate other land-use changes induces by urban growth. In most of the studies using SLEUTH described in the scientific literature only the UGM component of SLEUTH is applied.

UGM has been implemented in the modelling platform XULU in a modified way [12]. UGM now only needs four spatial input parameters namely a map of urban land-use, transportation, slope and exclusion. The exclusion layer determines, which areas in the research area cannot be changed (e.g. water bodies or protected areas) or, if not excluded, are by a certain degree resistant against urbanization. The transportation layer represents the road network in a research area. While SLEUTH needs at least four urban land-use data sets to calculate a set of calibration coefficients [4], the modified UGM in XULU only needs a map for the starting year of the calibration phase and a reference map at the end year. The simulated urban area of the end year is compared to the reference map with the Multiple Resolution Validation (MRV) as described in [30].

Calibration is the most crucial step in any modelling application [14]. In the calibration phase of UGM a brute-force method is used in order to determine five calibration parameters. These parameters control the transition rules that are implemented in the model and include: dispersion, breed, spread, slope resistance and road gravity. Dispersion determines the dispersiveness of the outward distribution and controls the number of pixels that are selected randomly for possible urbanization. Breed refers to the probability that a newly generated settlement starts its own growth. Spread controls how much existing settlements radiate. Slope resistance influences the likelihood of growth on steep slopes. Road gravity influences the creation of new centers along roads.

A number of Monte-Carlo iterations are performed in the brute-force calibration to obtain the best set of the five calibration parameters. This consequently translates to four different kinds of urban growth: spontaneous growth, diffusive (new spreading centers), organic (infill and edge growth) and road influenced growth. Because testing all possible parameter combinations in Monte-Carlo iterations in a brute-force way would be way too time consuming, calibration is performed in sequential phases ranging from a coarse to a fine calibration [4].

UGM's underlying simulation technique is CA. A CA is a discrete dynamic system in which space is divided into regular spatial cells, and time progresses in discrete steps [31]. Each cell in the system has one of a finite number of states. The state of each cell is updated according to local rules, that is, the state of a cell at a given time depends on its own state and the states of its neighbors at the previous time step [32]. Cellular automata are seen not only as a framework for dynamic spatial modelling but as a paradigm for thinking about complex spatio-temporal phenomena and an experimental laboratory for testing ideas [33]. A cellular automaton consists of five basic elements namely cell space, cell state, cell neighborhood, transition rules and time.

The number and location of the randomly selected cells is controlled by the growth parameters. Depending on the type of growth different properties of the selected cells are investigated. For the diffusive growth (new spreading centers) e.g. this would be the existence of non-built-up cells in the direct proximity of a selected built-up cell and the slope of these cells. Depending on the specified parameters and transition rules the CA computes, if a cell is available for change or not. The CA knows only two states: 1 = urban/built-up and 0 = non- urban/non-built-up.

XULU is a stand-alone JAVA application and serves as a modelling framework whose functionalities include input, output, editing and visualization. XULU offers a model independent graphical user interface. The core program comprises the fields of data management, input/output routines for data import and export, data structure, memory management and data visualization [11]. The user has to load the necessary data objects into the data pool and allocate them to the individual model resources. Several plug-ins of land-use modelling are implemented in XULU include: spatial data types for raster and vector data, I/O routines for shape files and different raster types (e.g. ASCII and GeoTIFF) and a layer-based visualization for raster and vector maps [11]. Additionally land-use change models are loaded as plug-ins. Models that are implemented so far include CLUE-s and the Urban Growth Model UGM [12].

Model calibration is required in order to ensure that a model simulates the reality fully. Diverse model users have different ways of assessing land-use models. Whereas there is a group of model users who wish to make predictions as accurate as possible, another group emphasizes on the ability of a model to support the general knowledge of processes and mechanisms of land-use change [34].

The method of multiple resolution validation (MRV) was used in a comparison of land-use models in which the tests were conducted in seven laboratories with 13 applications, 9 different models and in 12 study areas [35]. Typically maps are compared pixel-wise and every pixel is calculated as an error, where the model map does not exactly fit with the reference map. In MRV method four neighboring pixels are averaged stepwise. The amount of correct pixels increases in every step, until in the last step when the whole research area is inside of one big pixel and both location agreement and location disagreement approach 0 [12]. The MRV technique was incorporated in UGM.

In order to evaluate model results with the technique described above, three datasets are necessary: a reference map of time 1, a reference map of time 2 and a simulation map of time 2. The reference map of time 1 is the initial point for modelling, that is, land-use map of 1986 and at the same time serves as a Null-model, which is the assumption that no change has taken place. Therefore, the reference maps of t1 and t2 are compared. To evaluate the model result, the simulation map of t2 is compared with the reference map of t2.

## IV. Scenarios Of Urban Growth

The use of scenarios to address land-use changes have become useful tools in the assessment of land-use dynamics [36]. This approach is required to anticipate the consequences of various development scenarios. However scenarios are not predictions but rather they are an approach to help manage decisions based on the interpretation of qualitative descriptions of alternative futures translated into quantitative scenarios [37]. There is need for such scenarios to be integrated in land legislation. Several policies and strategies have been formulated by various national and regional governments in order to minimize the negative impacts caused by improper urban developments [38]. However, such policies are not well defined in the context of Kenya. Thus, exploring various scenarios by predicting future urban land-use patterns under different ''what-if'' conditions can help in the management of urban expansion and change as well as in the development of alternative plans before irreversible transformations occur [39]. This paradigm can help Kenya to manage its resources sustainably.

The Government of Kenya formulated Kenya Vision 2030 [16]. This was an attempt at maximum protection of natural resources so as to ensure sustainable development is attained in the year 2030. Cities in Kenya have undergone rapid urbanization as people migrate into cities in search of employment and better amenities. Thus this gave us the motivation to investigate scenarios of urban growth in Nairobi. Currently there are a few studies on scenario-based urban growth simulation in Nairobi. Nevertheless, [1] used SLEUTH to model urban growth in Nairobi.

In order to test the usefulness of the urban growth modelling and to provide a coherent and alternate framework for the policy makers, we explored three scenarios in the modelling process. First scenario depicts an unmanaged growth with no restriction on environmental areas, such as forest, agriculture and wetland. Thus urban growth continues with the historical trend of land transition and permits future urban growth allocation without any constraint. The second scenario assumes a managed growth with moderate protection. Here the exclusion layer included government buildings and forest cover. Cities in Kenya have undergone rapid urbanization due to high rural to urban migration as people search for employment and social amenities [2]. There has been significant effect to preserve forest cover in Kenya under the Forest Act, 2005 [40]. The third scenario simulates a managed growth with maximum protection on forest, government

reserved areas, government buildings, military bases, airports, and urban green. Government reserved areas include parks, cemeteries.

## V. Analysis

Scenarios based urban growth modelling of Nairobi involved datasets preparation, land-use change analysis and modelling using UGM. Fig. 2 shows the flow chart of the major steps applied in this research.

### A. Data

Modelling of Nakuru utilized urban extents extracted from land-use maps for 1986 and 2010 as inputs. Other layers used included slope, areas excluded from development and road network. The road layer included three weight values of 100, 50 and 25 [26]. A weight value of 100 was assigned to class A roads (International trunk roads), 50 was assigned to class B and C roads (National Trunk Roads), and 25 was assigned to local streets (Minor roads). The road classification in Kenya is explained in [41]. Thus a road with a value of 100 has the highest potential of attracting urban growth compared to a local street with a value of 25.

### B. Land-use change analysis

Land-use classification of Nairobi consisted of six land-use classes; namely urban, forest, agriculture, open/transitional areas, water and rangeland. Urban land-use included built-up areas within the research area. Forest included evergreen forest, mixed forests with high densities of trees, little or under-storey vegetation. Open/transitional areas included bare land, exposed areas, quarries and transitional areas. Water included rivers and reservoirs. The sewage treatment plant in Ruai was also captured under water class. Rangeland included bush land and ground layer covered by grass and sparsely disturbed scrub species.

Image pre-processing steps for the optical datasets were radiometric correction and geometric correction. Support vector machine (SVM) classification was applied to all the data sets and its performance assessed using error matrices. Recently SVM has been found to perform better compared to maximum likelihood classifier [42]. Post-classification refinements were enforced to diminish categorization errors as a result of the similarities in spectral signatures of certain classes. Spatial modeler and additional rule based procedures were adopted to overcome these classification challenges and differentiate between classes.

### C. Modelling using UGM

Model calibration of UGM involved running the model using default parameters of slope, breed, dispersion, road and spread. The default parameter values were 1, 50 and 100. Model calibration was done iteratively in four sequences from coarse to fine calibration as the parameters were varied using Monte Carlo technique. The MRV method was used to achieve the optimal parameterization for the UGM during the calibration phase as well as for the validation of the model results.

Fig. 2.  Flow chart of urban growth modelling

Three scenarios were explored in the modelling process. This involved varying the exclusion layer so as to achieve three scenarios. In the first scenario there was no restriction on the exclusion layer and thus the exclusion was at zero percentage. In the second scenario we achieved exclusion at 60 % exploring a managed growth with moderate protection. In the third scenario we achieved exclusion at 90 % exploring a managed growth with maximum protection. Thus it was not practical to achieve 100 % exclusion in scenario three since urbanization has already taken place.

## VI.  RESULTS AND DISCUSSION

Land-use summary for Nairobi was performed and results tabulated in Table 1 and Fig. 3. Land-use maps for Nairobi are illustrated on Fig. 4, Fig. 5 and Fig. 6. The urban/built-up areas increased from 35.16 km$^2$ in 1986 to 52.50 km$^2$ in 2000 and 79.38 km$^2$ in 2010. Forest increased from 62.87 km$^2$ in 1986 to 71.14 km$^2$ in 2000 but decreased to 66.86 km$^2$ in 2010. In areas where forest decreased such land was classified as agriculture or urban due encroachment of the forest. Agriculture increased from 144.72 km$^2$ in 1986 to 152.53 km$^2$ in 2000 but decreased to 148.21 km$^2$ in 2010.

Typical agriculture land-use include small-scale crop gardens and peri-urban agriculture for cultivation, and such land-use was converted to urban land-use namely building up of residential and commercial buildings to cater for the increased urban population in Nairobi. Open/Transition areas increased from 99.54 km$^2$ in 1986 to 146.94 km$^2$ in 2000 but decreased to 117.94 km$^2$ in 2010. Rangeland increased from 361.11 km$^2$ in 1986 to 261.74 km$^2$ in 2000 but decreased to 257.61 km$^2$ in 2010. Water increased from 9.60 km$^2$ in 1986 to 11.15 km$^2$ in 2000 and increased further to 26.00 km$^2$ in 2010.

The final model coefficients obtained after successful calibration of UGM for the three scenarios are illustrated in Table 2. We can see the values as follows: slope at 50, spread at 25, dispersion at 1, breed at 50, road at 75, and a weighted value of 0.9449 for scenario one; slope at 52, spread at 25, dispersion at 1, breed at 50, road at 25, and a weighted value of 0.9470 for scenario two; and slope at 52, spread at 27, dispersion at 1, breed at 52, road at 2, and a weighted value of 0.9477 for scenario three. We adopted scenario three since it will ensure sustainable development is met in the future.

TABLE I.        LAND-USE SUMMARY AND ERROR ESTIMATES FOR NAIROBI

| Year | 1986 | | 2000 | | 2010 | |
|------|------|---|------|---|------|---|
| *Land-use classes* | *Area (km²)* | *%* | *Area (km²)* | *%* | *Area (km²)* | *%* |
| Urban | 35.16 | 4.9 | 52.50 | 7.4 | 79.38 | 11.1 |
| Forest | 62.87 | 8.8 | 71.14 | 10.0 | 66.86 | 9.4 |
| Agriculture | 144.72 | 20.3 | 152.53 | 21.4 | 148.21 | 20.8 |
| Open/transition areas | 99.54 | 14.0 | 146.94 | 20.6 | 117.94 | 16.5 |
| Rangeland | 361.11 | 50.6 | 261.74 | 36.7 | 257.61 | 36.1 |
| Water | 9.60 | 1.3 | 11.15 | 1.6 | 26.00 | 3.6 |
| Total | 696 | 100 | 696 | 100 | 696 | 100 |
| | | | | | | |
| Overall Accuracy (%) | 92.64 | | 90.9 | | 91.87 | |



Fig. 3.   Land-use estimates for Nairobi

Fig. 4.   Land-use map for Nairobi in 1986



Fig. 5.   Land-use map for Nairobi in 2000

Fig. 6.   Land-use map for Nairobi in 2010

TABLE II.          BEST MODEL PARAMETERS OBTAINED IN THE THREE SCENARIOS

| Scenario | Model parameters | | | | | Weighted value |
| | Slope | Spread | Dispersion | Breed | Road | |
|---|---|---|---|---|---|---|
| 1 | 50 | 25 | 1 | 50 | 75 | 0.9449 |
| 2 | 52 | 25 | 1 | 50 | 25 | 0.9470 |
| 3 | 52 | 27 | 1 | 52 | 2 | 0.9477 |

An evaluation of the three scenarios was conducted as shown in Table 3. The simulated urban growth values of scenario one, two and three were 82.87 $km^2$, 76.61 $km^2$, and 73.14 $km^2$ in 2010 and 141.72 $km^2$, 127.96 $km^2$, and 118.35 $km^2$ in 2030 respectively. The urban growth simulation maps for all scenarios are illustrated in Fig. 7, Fig. 8, Fig. 9, Fig. 10, Fig. 11 and Fig. 12.

We conducted two map comparisons in Erdas imagine 2011 model maker for scenario three for the city of Nairobi. According to [35] there are three possible two-map comparisons namely observed change, prediction change and prediction error as described above in urban growth modelling of Nairobi. Observed change compares the reference map of time 1 and the reference map of time bearing in mind the dynamics of the landscape. Prediction change compares between the reference map of time 1 and the prediction map of time 2 and thus revealing the behavior of the model. Prediction error compares between the reference map of time 2 and the prediction map of time 2 and thus ascertains the accuracy of the

prediction. In our case time 1 referred to as the year 1986 and time 2 as the year 2010.

The observed change in urban land-use between 1986 and 2010 is illustrated on Fig. 13. Here we have observed built gain of 65.25 $km^2$, observed built persistence 17.80 $km^2$, observed non-built persistence of 607.19 $km^2$ and observed built loss of 3.41 $km^2$ obtained from the observed map of the year 2010. The predicted change in urban land-use between 1986 and 2010 is illustrated on Fig. 14. Here we have predicted built persistence of 73.79 $km^2$ and predicted non-built persistence of 620.51 $km^2$ obtained from the predicted map of the year 2010. The predicted error in urban land-use between 1986 and 2010 is illustrated on Fig. 15. Here we have: non-built observed and built predicted of 57.95 $km^2$; and built observed and built predicted of 40.99 $km^2$ obtained using the observed map of the year 2010 and the predicted map of the year 2010. Our UGM for Nairobi predicts the year 2010 accurately since our gain of built is larger than the loss of built by16.96 $km^2$.

TABLE III.  MODEL EVALUATION FOR NAIROBI

| Year | 2010 | | | 2030 | | |
|---|---|---|---|---|---|---|
| *Scenario* | *1* | *2* | *3* | *1* | *2* | *3* |
| Actual Urban (km$^2$) | 79.38 | 79.38 | 79.38 | | | |
| Simulated Urban (km$^2$) | 82.87 | 76.61 | 73.14 | 141.72 | 127.96 | 118.35 |



Fig. 7.  Urban growth simulation in Nairobi in scenario one (2010)



Fig. 8.  Urban growth simulation in Nairobi in scenario one (2030)

Fig. 9.   Urban growth simulation in Nairobi in scenario two (2010)



Fig. 10. Urban growth simulation in Nairobi in scenario two (2030)

Fig. 11. Urban growth simulation in Nairobi in scenario three (2010)



Fig. 12. Urban growth simulation in Nairobi in scenario three (2030)

In order for an urban growth model to be resourceful to various stakeholders such as policy makers and urban planners, simulation of urban growth has to be performed after calibration. Scenarios three was selected as the best plausible cause for urban planning management with maximum protection on resources. Thus the likelihood of new settlements or built-up areas in Nairobi was obtained at a weighted value of 0.9477 as per scenario three. This indicates that new urban growth is most likely to be caused by breed (at 52), i.e. probability that a newly generated settlement starts its own growth, then followed by slope (at 52) influenced growth and spread (at 27), and finally followed by road and dispersion as least likely factors for new urban growth. Thus, this implies that new areas are developed for residential and commercial uses, which lie in proximity to roads. Such growth could be as a result of high rural urban migration witnessed in Nairobi as new people move immigrate in search for employment, social amenities and business opportunities.

## VII. CONCLUSION

We used Nairobi, Kenya's capital city as an example of a fast expanding African city to analyze the dynamics of land-use changes between 1986 and 2010, and to simulate urban growth into 2030 using cellular automata. Land-use change analyzed demonstrated that substantial changes have taken place as a result of rapid urban growth. Urban land-use maps from image classification were used alongside other datasets in modelling urban growth in Nairobi using UGM. The Monte Carlo iterative method was applied in the UGM calibration. Three scenarios were explored in the urban modelling process; unmanaged growth with no restriction on environmental areas, managed growth with moderate protection, and a managed growth with maximum protection on forest, agricultural areas, and urban green. Scenario three was selected as a plausible paradigm to ensure sustainable development is achieved.

Kenya plans to achieve Vision 2030 in the year 2030 and this can be guided using scenario based urban growth. Thus to achieve the economic and social strategy there is need for land-use scenarios as we conducted in this research in order to cater for anticipated urban growth in the future. Urban growth modelling is vital for guiding decision making for resource management.

Simulated urban growth results for the year 2030 using scenario three indicate that there is the need for tactical planning so as to address rapid urban growth in Nairobi.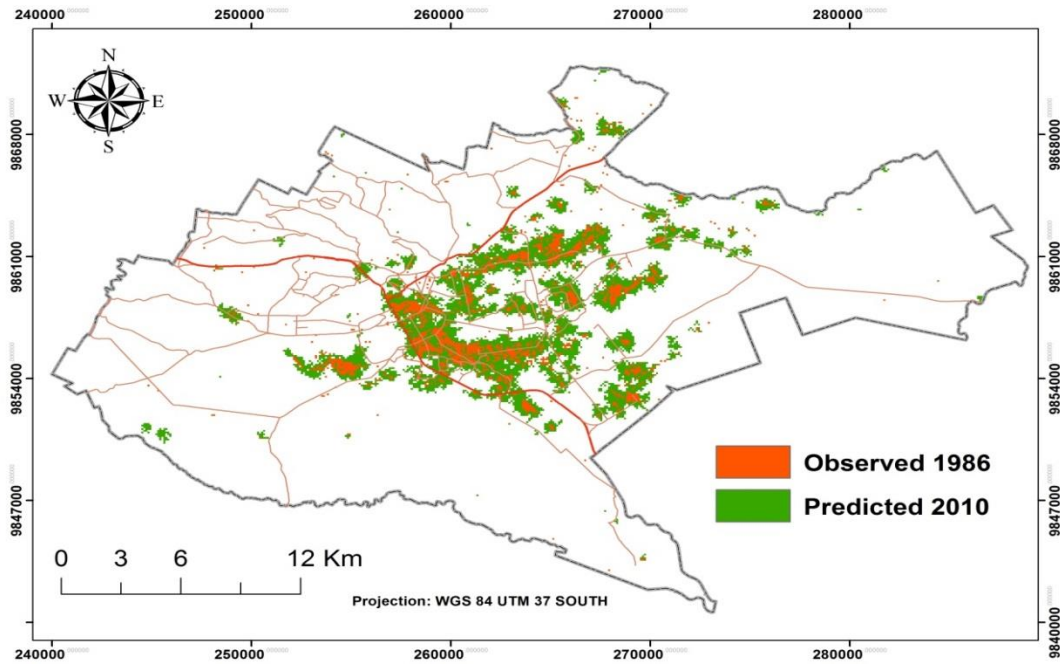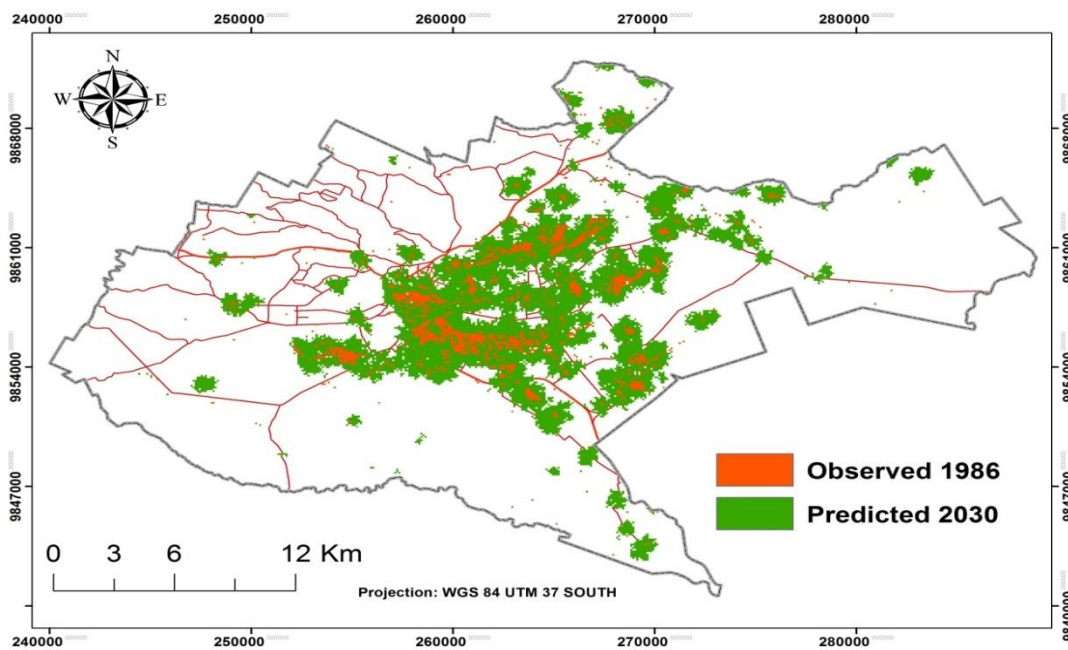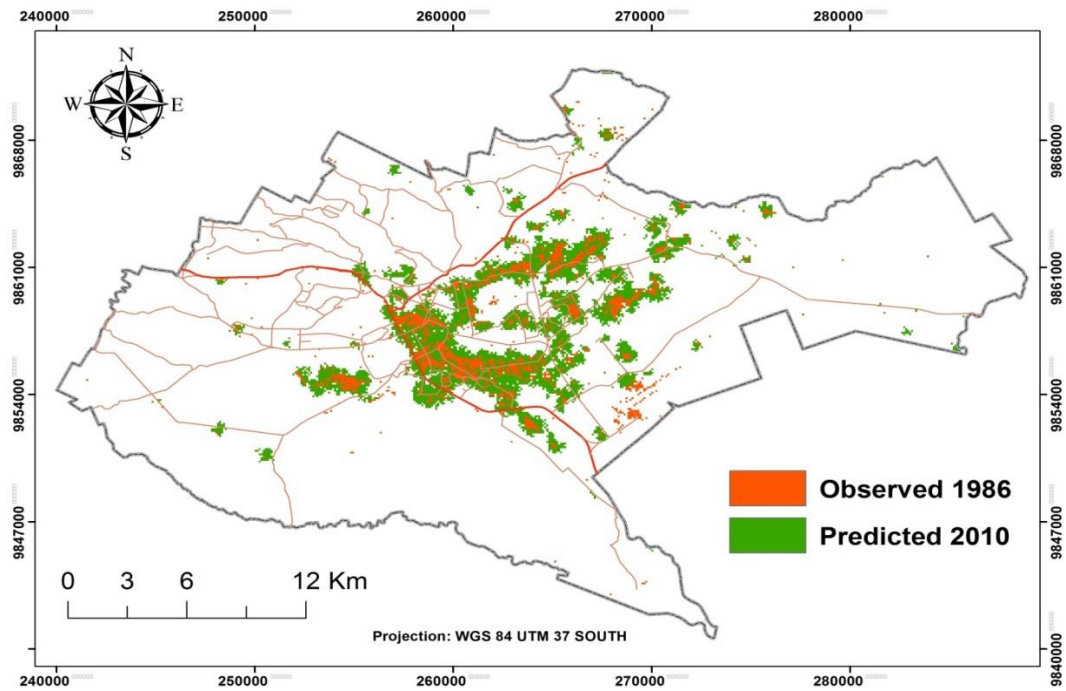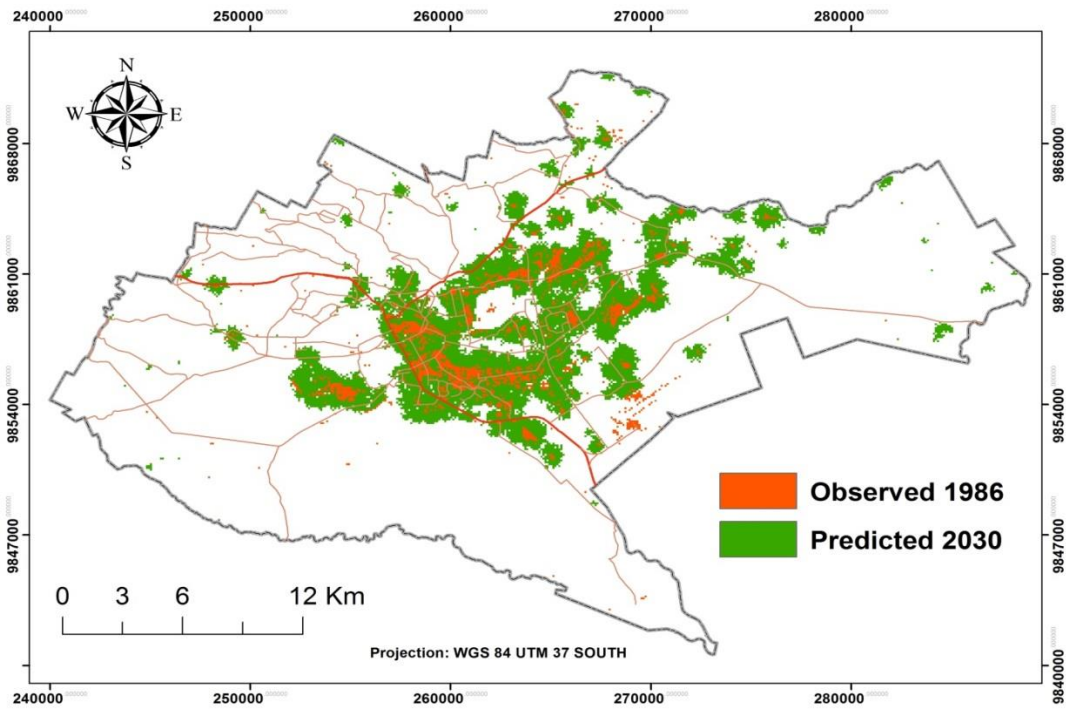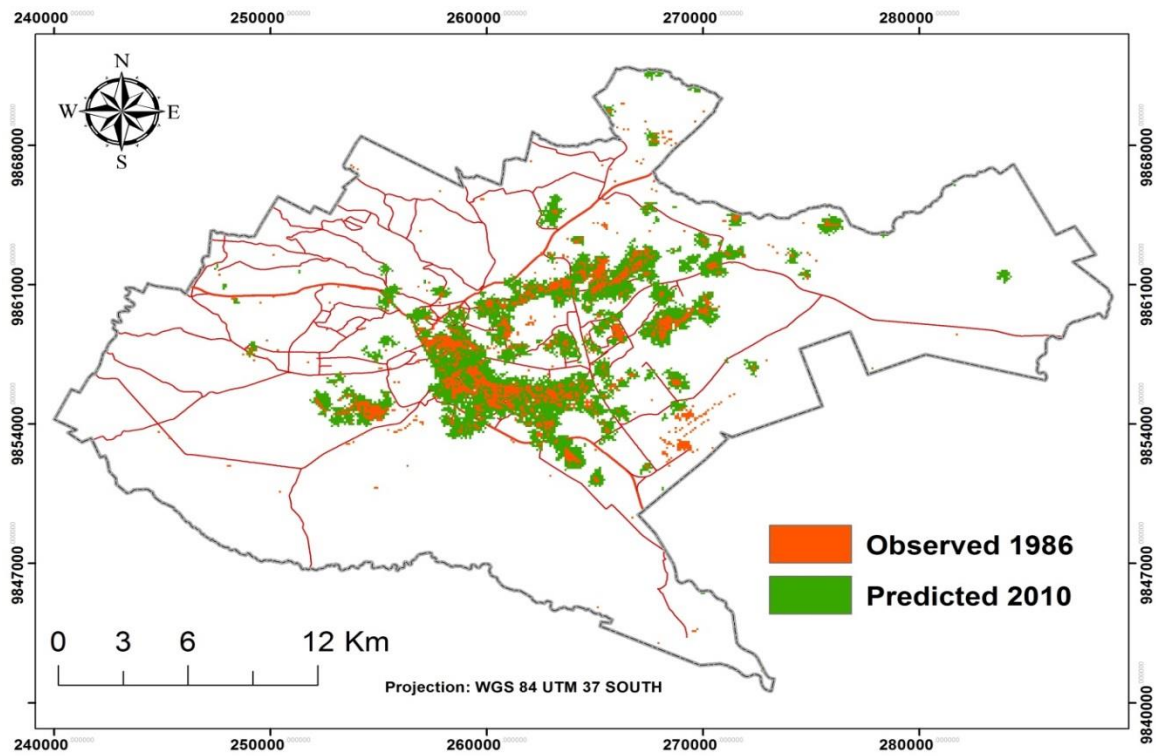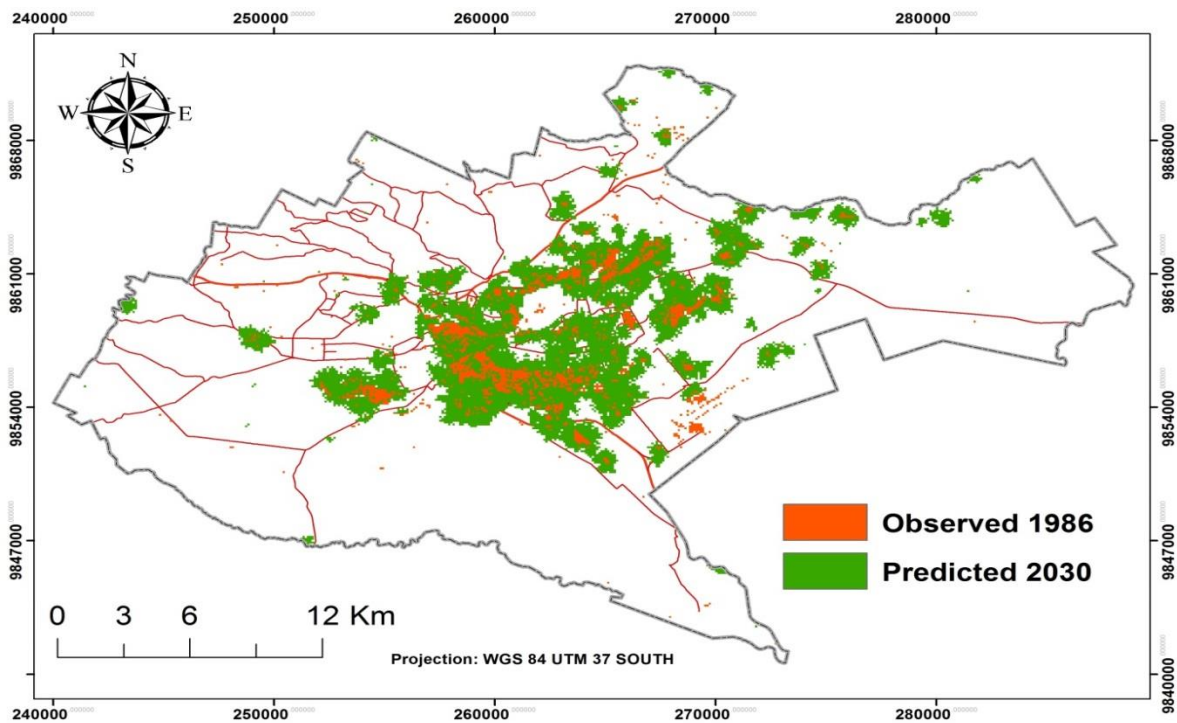 Therefore cellular automata are a valuable approach for regional modelling of big African cities such as Nairobi. Hence it is noble to explore the use of UGM in other cities in Africa and its performance documented accordingly

### References

[1] C. N. Mundia and M. Aniya, "Modeling urban growth of Nairobi city using cellular Automata and Geographical information systems," Geographical Review of Japan, vol. 80, no. 12, pp. 777-788, 2007.

[2] C. N. Mundia and M Aniya, "Analysis of land use changes and urban expansion of Nairobi city using remote sensing and GIS," International Journal of Remote Sensing, vol. 26, no. 13, pp. 2831-2849, 2005

[3] C. N. Mundia and M. Aniya, "Dynamics of land use/cover changes and degradation of Nairobi city, Kenya," Land Degradation and Development, vol. 17, pp. 97-108, 2006.

[4] E. Silva and K. C. Clarke, "Calibration of the SLEUTH urban growth model for Lisbon and Porto, Portugal," Computers, Environment and Urban Systems, vol. 26, no. 6, pp. 525–552, 2002.

[5] G. Orcutt, M. Greenberger, A. Rivlin, and J. Korbels, Micro analysis of social economic systems: A simulation study. New York: Harper and Row Publishers, 1961.

[6] F. Wu, "GIS based simulation as an exploratory analysis for space: time processes," The Journal of Geographic Systems, vol. 1, no. 3, pp. 199-218, 1999.

[7] N. Oreskes, K. Sharader-Freschete, and K. Belitz, "Verification, validation, and confirmation of numerical models in earth sciences," Science, vol. 263, pp. 641-646, 1994.

[8] M. Wegener, "Operational urban models: State of the art," Journal of the American Planning Association, vol. 60, no. 1, pp. 17-30, 1994.

[9] A. Hill and C. Lindner, "Simulation informal urban growth in Dar es Salaam, Tanzania - A CA-based land-use simulation model supporting strategic urban planning," in Modeling and Simulating Urban Processes. Munster: LIT-Verlag, 2011, pp. 77-98.

[10] P. Verburg, P. Schot, M. Dijst, and A. Veldkamp, "Land use change modelling: current practice and research priorities," GeoJournal, vol. 61, no. 4, pp. 309-324, 2004.

[11] M. Schmitz, T. Bode, H. P. Thamm, and A. B. Cremers, "XULU - A generic JAVA-based platform to simulate land use and land cover change ( LUCC )," in MODSIM 2007 International Congress on Modelling and Simulation, 2007, pp. 2645–2649.

[12] R. Goetzke and M. Judex, "Simulation of urban land-use change in North Rhine- Westphalia ( Germany ) with the Java-based modelling plat- form Xulu," in Modeling and Simulating Urban Processes. Munster: LIT-Verlag, 2011, pp. 99–116.

[13] K. Clarke, S. Hoppen, and L. Gaydos, "A self-modifying cellular automaton model of historical urbanization in the San Francisco Bay area," Environment and Planning B: Planning and Design, vol. 24, no. 2, pp. 247-261, 1997.

[14] K. Clarke, S. Hoppen, and L. Gaydos, "Methods and techniques for rigorous calibration of cellular automaton model of urban growth," in Third International Conference/Workshop on Integrating GIS and Environmental Modeling, Santa Fe, 1996.

[15] UN-HABITAT, "State of the World Cities 2010/2011," Bridging the Urban Divide, 2010.

[16] Government of Kenya, Kenya Vision 2030, 2007.

[17] UN-HABITAT, Regional Urban Sector Profile Study (RUSPS), 2005.

[18] Republic of Kenya, Kenya population census 1969, 1970.

[19] Republic of Kenya, Kenya population census 1979, 1981.

[20] Republic of Kenya, Kenya population census 1989, 1994.

[21] Republic of Kenya, Economic survey 2000, 2000.

[22] Republic of Kenya, Economic survey 2010, 2010.

[23] J. Han, X. Cao, and H. Imura, "Application of an integrated system dynamics and cellular automata model for urban growth assessment: A case study of Shanghai, China," Landscape and Urban Planning, vol. 91, no. 3, pp. 133–141, 2009.

[24] G. Menz et al., "Land use and land cover modeling in Central Benin," in Impacts of Global Change on the Hydrological Cycle in West and Northwest Africa. Heidelberg: Springer, 2010, pp. 70-73.

[25] P. H. Verburg et al., "Modelling the Spatial Dynamics of Regional Land Use: The CLUE-s Model," Environmental Management, pp. 391-405, 2002.

[26] E. Silva and K. Clarke, "Complexity, Emergence and Cellular Urban Models: Lessons Learned from Applying Sleuth to Two Portuguese Metropolitan Areas," European Planning Studies, vol. 13, no. 1, pp. 93-115, 2005.

[27] S. Leão, I. Bishop, and D. Evans, "Spatial-temporal model for demand and allocation of waste landfills in growing urban regions," Computers, Environment and Urban Systems, vol. 28, pp. 353–385, 2004.

[28] L. Lebel, D. Thaitakoo, S. Sangawongse, and D. Huaisai, "Views of Chiang Mai: The Contribution of Remote-Sensing to Urban Governance and Sustainability," in Applied Remote Sensing for Urban Planning, Governance and Sustainability. Berlin: Springer, 2007, pp. 221-247.

[29] J. T. Candau and K. C. Clarke, "Probabalistic Land Cover Transition Modeling Using Deltatrons," in Proceedings of Urban and Regional Information Systems Association (URISA) 38th Annual Conference, Orlando, 2000.

[30] R. G. Pontius Jr, E. Shusas, and M. McEachern, "Detecting important categorical land changes while accounting for persistence," Agriculture, Ecosystems and Environment, vol. 101, no. 2-3, pp. 251-268, 2004.

[31] Y. Liu, Modelling Urban Development with Geographical Information Systems and Cellular Automata, 1st ed. Florida: CRC Press, 2008.

[32] S. Wolfram, "Universality and complexity in cellular automata," Physica, vol. 10D, pp. 1-35, 1984.

[33] R. M. Itami, "Simulating spatial dynamics: cellular automata theory," Landscape and Urban Planning, vol. 30, no. 1-2, pp. 27-47, 1994.

[34] R. G. Pontius Jr and J. Malanson, "Comparison of the structure and accuracy of two land change models," International Journal of Geographical Information Science, vol. 19, no. 2, pp. 243-265, 2005.

[35] R. G. Pontius Jr et al., "Comparing the input, output, and validation maps for several models of land change," Annals of Regional Science, vol. 42, no. 1, pp. 11-47, 2008.

[36] P. H. Verburg, C J Schulp, N Witte, and A Veldkamp, "Downscaling of land use change scenarios to assess the dynamics of European landscapes," Agriculture, Ecosystems and Environment, vol. 114, no. 1, pp. 39-56, 2006.

[37] L. O. Petrov, C. Lavalle, and M. Kasanko, "Urban land use scenarios for a tourist region in Europe: Applying the MOLAND model to Algarve, Portugal," Landscape and Urban Planning, vol. 92, no. 1, pp. 10-23, 2009.

[38] Q. Zhang, Y. Ban, J. Liu, and Y. Hu, "Simulation and analysis of urban growth scenarios for the Greater Shanghai Area, China," Computers, Environment and Urban Systems, vol. 35, no. 2, pp. 126-139, 2011.

[39] T. M. Conway and R. G. Lathrop, "Modeling the ecological consequences of land-use policies in an urbanizing region," Environmental management, vol. 35, no. 3, pp. 278-291, 2005.

[40] Laws of Kenya, Forests Act, 2012.

[41] W. S. Wasike, "Road infrastructure policies in Kenya: historical trends and current challenges.," Development, p. 41, 2001.

[42] K. Mubea and G. Menz, "Monitoring Land-Use Change in Nakuru (Kenya) Using Multi-Sensor Satellite Data," Advances in Remote Sensing, vol. 1, no. 3, pp. 74–84, 2012.

# Construction Strategy of Wireless Sensor Networks with Throughput Stability by Using Mobile Robot

Kei Sawai

Department of Information and Communication
Engineering, Tokyo Denki University
Tokyo, Japan

Yuta Koike

Department of Information and Communication
Engineering, Tokyo Denki University
Tokyo, Japan

Shigeaki Tanabe

Technical Support Department,
Technology Institution of Industrial Safety (TIIS)
Sayama city, Japan

Ryuta Kunimoto

Department of Information and Communication
Engineering, Tokyo Denki University
Tokyo, Japan

Hitoshi Kono

Department of Information Communication and Media
Design Engineering, Tokyo Denki University
Tokyo, Japan

Tsuyoshi Suzuki

Department of Information and Communication
Engineering, Tokyo Denki University
Tokyo, Japan

*Abstract*—We propose a wireless sensor networks deployment strategy for constructing wireless communication infrastructures for a rescue robot with considering a throughput between sensor nodes (SNs). Recent studies for reducing disaster damage focus on a disaster area information gathering in underground spaces. Since information gathering activities in such post disaster underground spaces present a high risk of personal injury by secondary disasters, a lot of rescue workers were injured or killed in the past. Because of this background, gathering information by utilizing the rescue robot is discussed in wide area. However, there are no wireless communication infrastructures for tele-operation of rescue robot in the post-disaster environment such as the underground space. Therefore, we have been discussing the construction method of wireless communication infrastructures for remotely operated the rescue robot by utilizing the rescue robot. In this paper, we evaluated the proposed method in field operation test, and then it is confirmed that maintaining communication connectivity and throughputs between End to End of constructed networks.

*Keywords—Wireless Sensor Networks; Rescue Robot Tele-Operation; Maintaining Throughput*

## I. INTRODUCTION

Gathering information in disaster areas is very important for assessing the situation, avoiding secondary disasters, and managing disaster reduction [1]-[8]. In general, bird's-eye image information gathered by unmanned air vehicles (UAVs) and artificial satellites is useful for understanding post-disaster situation. However, in an underground space in the city part where such UAVs etc. cannot gather information, it is difficult to ascertain the extent of the damage, which is important for avoiding secondary disasters. Also, rescue teams cannot organize a suitable rescue plan for underground spaces because sufficient information is not gathered. Under such a situation, the rescue team must go into the underground spaces directly to gather disaster information, and the information should be shared within the teams by communication between above ground and underground space for efficient and cooperative rescue works. However, when the communication infrastructure is broken due to damage, rescue teams cannot cooperate closely because of communication disconnection. Therefore, the rescue team has to work in the underground space with being unable to know the situation correctly, and they face the added risk of secondary disasters. For example, in the underground disasters in Korea in 2003, a lot of rescue workers were sacrificed because of smoke damage. The rescue teams could not expect the smoke damage because they could not gather enough information about post-disaster situation in underground space, thus many lives were lost. This is a typical case of underground disaster damage due to secondary disasters that has triggered because the rescue teams entered underground areas without adequate information.

From discussions based on past accidents analysis, researchers have recently focused on a disaster information gathering method using a wireless sensor network (WSN) and a rescue robot in closed areas. The WSN consists of spatially distributed sensor nodes (SN) to cooperatively monitor the environmental conditions such as temperature, sound, vibration, pressure, motion, etc. Then, the WSN is enabled to provide the wireless communication function in place without existing infrastructure. The WSN in closed area is constructed by rescue robot. Therefore, an information gathering method by constructing the communication infrastructure to disaster area by using the WSN has been discussed.

One of them, a SN deployment strategy by utilizing the rescue robot is very important to the performance evaluation of adaptability in closed space. Many SN deployment strategies have been discussed in the WSN research field. In these strategies, deployment methods have been proposed based on

evaluation scales that consider factors such as packet routing, energy efficiency, power saving, and coverage area. Several SN deployment methods using mobile SNs and mobile robots to construct the WSN have been developed. Parker et al. proposed the WSN construction method using an autonomous helicopter for environmental monitoring and urban search and rescue [9]. Umeki et al. proposed an ad-hoc network system, Sky Mesh, using a flying balloon for targeted disaster rescue support [10]. Also, deployment methods have been developed based on virtual interaction between the SNs based on several physical models; such as the potential field model and the fluid flow model [11]-[17].

A great deal of effort has been made on SN deployment strategy. However, what seems to be lacking is the strategy that is considered the specification of underground space and the construction method of with concerning a throughput quality for expanding the area where is able to operate remotely the rescue robot. There are a lot of shielding materials of electrical wave. Then a discussion of the construction method of the WSN with concerning specification of underground space is important to prevent a network disconnection. Then to expand the network with stable throughput is important to maintain the information gathering system, it is necessary to prevent the secondly disaster.

Therefore, we have been discussing the information gathering system that is considered these important matters (Fig. 1) [18]-[21]. In this paper, we proposed the novel SN deployment strategy to construct the WSN with the stability of communication connectivity by utilizing the rescue robot. Then we evaluated the availability of the proposed method in field operation test.



Fig. 1. Gathering disaster area information by utilizing wireless sensor networks and rescue robot

## II. Deployment Strategy For Maintaining Communi-Cation Connectivity By Utilizing Rescue Robot

### A. Prior Conditions

In our proposed system, the WSN is constructed by utilizing rescue robot to deploy SNs. In the construction environment that is deployed SN, we assume the place that has entrance stairs and the first basement floor.

First of all, the entrance stairs in under-ground is required to set up at intervals 30 [m], and passage way is built in line in Japanese building standard low. Therefore, in our proposed system gathers this area's information. Then we discussed the WSN construction method by utilizing rescue robot in this area. In the wireless communication of this WSN, IEEE 802.11 series are adopted for wireless communication between SNs including the rescue robot, which has been used as proven communication in many studies of mobile robot and the WSN [22]-[26]. Then in our proposed method, we treat a rescue robot as a SN in the WSN. Heterogeneous networks that are involved some SN and various mobile robots are difficult to manage the system control. Especially, the maintaining the stability of the system control is not easy by occurring secondly disaster in underground spaces. In this environment, to construct the stable system is necessary to simplify the network structure. Therefore, we simplified the network structure by treating a rescue robot as a SN. From here onwards, the communication system of the rescue robot is adopt the IEEE802.11 series as same as the SN.

In our SN deployment method for constructing the WSN, we adopt the method that the rescue robot delivers the previously wireless connected SNs. The rescue robot deploys the SN in the own passageway. Then the WSN is expanded, the operator is able to control the rescue robot by utilizing the communication infrastructure of the WSN. In the network topology of this WSN, it is linearly connected each SNs to prevent the error of routing control. Generally, the WSN is able to decide the routing path of data transfer automatically by utilizing the RSSI between each SN, throughput of End to End or the rate of packet loss. The routing pass of the WSN is reconstructed by changes of these communication qualities.

However, the reconstruction of the routing pass repeatedly occurs the situation that is the disconnection and reconnection in between SNs. This situation is a problem for the system with tele operating the mobile robots. The tele-operating with abeyance of wireless communication degrades an operability of rescue robot and the performance of the gathering disaster area information. Then the change of the communication qualities is often occurred in disaster area by the damage of secondly disaster, the routing path is repeatedly reconstructed in the WSN. Therefore, we adopt the network topology that is linearly connected SNs, and it refers to the previously determined routing path to prevent the lowering of mobile robot activity.

### B. Requested Specifications

IEEE 802.11 series, it is necessary to keep the throughput that is more than 1.0 [Mbps] in between the operator and the rescue robot (End-to-End communications) (Add the references). Then in the construction of the WSN by utilizing the rescue robot, the throughput between End-to-End communications has to be maintained in the environment that is constructed the WSN. The construction length of the WSN is required 50 [m] by concerning the distance of first basement floor 30 [m] and entrance stairs 20 [m]. However, the communication connectivity of IEEE802.11 series is

characterized by decreasing in turn area covered with concrete material such as the underground space. Thus in our proposed system for constructing WSN, we should consider this communication characteristic that has a risk of network disconnection.

In the communication system of the SN and the rescue robot, it adopt IEEE802.11b as the system that has the high connectivity in the environment where has a lot of obstacles. The theoretical values of throughput by utilizing IEEE802.11b are 11.0 [Mbps], and then the actual measurement values are lowered around 7.0 [Mbps] by the efficiency of the various factors in the real environment. Then this wireless LAN protocol provides the communication distance that is 100 [m] as on the straight line. The throughput is satisfied with the required specification that is more than 1.0 [Mbps] to operate the mobile robot and the high connectivity. Therefore, we adopt IEEE802.11b to our proposed system. And then the throughput we defined is the amount of packet transferred per unit time in the networks.

When IEEE802.11b is adopted in the ad-hoc networks, the number of SNs that is able to maintain is more than 1.0 [Mbps] is 4 nodes in the situation that the entire throughput between each SN is more than 6.0 [Mbps]. In the function of ad-hoc networks constructing the WSN, the delay of data transfer is occurred with an increasing amount of the hop number in between the source and the destination of the network. Whence to linearly connect the SNs for expanding the WSN with maintaining the throughput, the number of the SN is required to decide for constructing the network in advance. Then this method provides the high connectivity in turn area covered with concrete material such as the underground space by deploying the SN as communication relay device to construct WSN. Also the constructed network consists of a source SN, three SNs that are deployed and a rescue robot regarded as the SN (Fig. 3).

## III. SN Deployment Method To Construct Lineally Networks

Our proposed SN deployment strategy is required the communication quality parameters to construct lineally networks. The communication qualities are measured the electrical field density (RSSI) and the throughput for the decision of SN deployment place. The lower of the packets throughput is occurred by the efficiency of various factors that is the multipath facing of the radio waves, the packet collisions, changing RSSI …etc. Thus it is difficult to construct the WSN with maintaining the throughput over 1.0 [Mbps] in the section of End to End. Therefore, our proposed method constantly monitors these parameters for the construction of stable communication infrastructure.

The rescue robot mounts three SNs, and then it expands the WSN by deploying these SNs. The entire SNs mounting on the rescue robot is linearly connected by reference the routing plan in advance. Then, SN is deployed by the state of network connection. The deployment with connecting the adjacent SNs has no steps that involving the reconstruction of connection by adding to the WSN. The reconstruction of WSN need to momentary disconnect the adjacent SNs of deploying SN, whence the operator cannot control the rescue robot in that split

second. The rescue robot with the status of uncontrolled tele-operating has the risk which is the occurring the losing of the rescue robot and the secondly disaster. Therefore, we adopted the method that connecting the entire SNs in advance.



Fig. 2.   Existing approach of wireless tele-operation



Fig. 3.   Constructing method of wireless sensor networks by utilizing rescue robot

To keep the throughput to over 1.0 [Mbps] in End to End, it requires maintaining the two communications qualities of between each adjacent SN. The RSSI in between two adjacent SNs (1 [Hop]) requires over -86 [dBm]. A wireless LAN module that controlling the throughput speed constantly refers the RSSI for stability of the network connection. If the RSSI value get down to under -86 [dBm], the wireless LAN module controls the throughput speed to under 6.0 [Mbps]. Whence our proposed method also should measure the RSSI to predict the throughput speed control of the wireless LAN module. Throughput requires the value more than 6.0 [Mbps] in between the deploying SN and the adjacent SNs on the condition that maintaining the throughput over 1.0 [Mbps] in between end to end. Moreover in the decision of deployment

position, measuring the End-to-End throughput is required to evaluate the communication quality between the operator and the rescue robot.

Therefore, our proposed algorithm requires the repetitive measurement of communication quality and the movement of rescue robot. The rescue robot in decided place deploys the SN. Fig. 4 shows the workflow of this deployment strategy. In the workflow, $N$ is parameter of previously deployed SN ID, $M$ is the next deployment SN. The workflow is outlined below.

*1)    Rescue robot deploys the first SN in the point of 0 [m], the deployed SN is numbered the ID "N" (initial value = 1). The secondary SN is numbered the ID "M" (initial value = 2).*

*2)    After the moving of the rescue robot, the operator constantly observes the RSSI between SN "N" and "M" in interval at 1.0 [m].*

*3)    If the RSSI of between "N" and "M" is higher than -86 [dBm], the operator measures the throughput of between the End-to-End communications. Moreover if the throughput is more than 1.0 [Mbps], the rescue robot keeps task to construct the WSN.*

*4)    If the deployed SN ID is "N= 1" in the situation that the RSSI is lower than -86 [dBm] or the throughput is not enough 1.0 [Mbps], the rescue robot goes back the place where is kept the communication qualities.*

*5)    After the movement, the rescue robot evaluates the throughput between End to End. If the throughput is stable, the rescue robot deploys the SN of ID "M". After the action of deployment, the value of "N" and "M" are changed to "N"=2 and "M"=3. The number of "N" and "M" are incremented a value after deployment of SN. (N=N+1, M=M+1)*

*6)    Then the rescue robot repeats above deployment action (2) - (5) until the number of "M" is incremented 5. Our proposed SN deployment strategy constructs WSN by utilizing above workflow.*

IV.    COMMUNICATION QUALITY EVALUATION OF CONSTRUCTED WIRELESS COMMUNICATION INFRASTRUCTURE BY USING PROPOSED MODEL

*A.  Experimental condition*

actually constructing the WSN composing the SNs and the rescue robot. Then, the rescue robot deploying the SNs constructed the WSN in this experiment. In the evaluation, evaluation item was targeted at the extended distance and the throughput in between End to End of the WSN.

To construct the WSN, we adopt the developed SN device shown in Fig. 5 in our previous studies. This SN mounts the CPU board, memory device, CompactFlash disc, IEEE 802.11b/g wireless LAN module, a digital camera, an A/D converter, and a battery. Then these devices of the SN are controlled by Linux OS (Debian). It enables to construct the WSN by utilizing the "AODV-uu" of the application connecting Ad-Hoc networks. Table 1 shows the specification of our developed SN.



Fig. 4.    Workflow of SN deployment method



Fig. 5.    Developed wireless sensor node

TABLE I.　　　SPECIFICATION OF WIRELESS SENSOR NODE

| Sensor Node | |
|---|---|
| Operating system | Linux Kernel 2.6 (Debian) |
| CPU board | Armadillo-300 (ARM 200[MHz]) |
| Web camera | Axis 207MW |
| Fish eye lens | Nissin 4CH190 (AOV 190[deg]) |
| Weight | 1.5 [kg] |
| Height×Width×Length | 225 [mm] × 180 [mm] × 380 [mm] |
| Battery No. 1 | Output : 5 [V], 1.8 [A] |
| Battery No. 2 | Output : 12 [V], 2.1 [A] |
| Operating time | 3 [hour] |

The crawler-type mobile robot,"S-90LWX" (TOPY INDUSTRIES, LIMITED), in Fig. 7 is adopted as the rescue robot in this experiment. The SN deployment mechanism was developed for the WSN construction and installed to the rescue robot, which can mount up to five SNs using five solenoid-operated locks. Figure 6 shows the framework of this mobile robot with the SN deployment mechanism. The mobile robot and the entire SNs are named IP address in this system. Then, the operator can operate the crawler robot and the deployment mechanism remotely by utilizing the TCP/IP and UDP.



Fig. 6.　Configuration of Mobile Robot



Fig. 7.　Rescue robot mounting sensor node

## B. Measurement Method

The operator linearly advances the mobile robot in a straight way. In this experiment, the RSSI and the throughput are measured to 10 times at every 1.0 [m] interval, and calculated the average at each measuring point. The throughput is measured in the between deploying SN and the adjacent SNs, and the End to End. Generally, a wireless communication quality in physical layer level is measured using the spectrum analyzer in anechoic chamber. For the RSSI measurement in this experiment, however, we used "iwlist" command contained in the Linux wireless tools package because we aimed to evaluate the transport layer level communication. To measure the packet throughput, "utest" (NTTPC Communications Ltd.) was used.

The experiment is performed in the passageway with a length of 300 [m] or more in Tokyo Denki University, and the rescue robot constructed WSN by utilizing our proposed algorithm in this environment (Fig. 8 and 9). Then we evaluated the distance that the mobile robot moved the without the SN deployment for the cooperative evaluation. In the experiment without the SN deployment, the mobile robot is controlled in area that keeping the throughput to over 1.0 [Mbps].



Fig. 8.　Experimental environment



(a) Experimental place　　　(b) S-90LWX with SNs

Fig. 9.　Overview of experimental environment

## C. Experimental Results

Figure 10 shows the results of the average of through-put between End to End communications, the value of the RSSI and the extending distance by utilizing our proposed algorithm. Arrowed lines on the graph indicate the SN deployment point and the extended distance. In the results, the extended distance with maintaining the throughput over 1.0 [Mbps] was 252 [m].

Then the distance without the extending method was 140 [m] that the remit keeping the throughput to 1.0 [Mbps].

Therefore, we confirmed the extended distance for the WSN construction with keeping the throughput of 1.0 [Mbps] is over the theoretical distance of the networks at this experimental situation.



Fig. 10. Experimental results of measured RSSI and throughput

## V. DISCUSSION

In this communication quality evaluation, we confirmed the WSN that has the maintenance capability of throughput between End to End communications was constructed by utilizing our proposed method. Throughput between the End to End was stable over 1.0 [Mbps] to the point of 252 [m] from the point of 0 [m]. The deployment point of SNs was 130 [m] (SN2), 180 [m] (SN3), 192 [m] (SN4) and 252 [m] (Rescue robot as SN5). Decreasing throughput was 0.7 [Mbps] at most in between SN1 and SN2, however, it was maintained over 1.2 [Mbps] in entire measurement point. It is assumed that this reduction of the throughput was occurred as a result of increase of communication distance.

Thus there was no significant decrease of throughput in between SN2 and SN3, SN3 and SN4, SN4 and Rescue robot (SN5). Then it was stable in between End to End. These experimental results is caused by the SN deployment point interval was short, thus the throughput was not affected by the attenuation of radio wave. In the reason that SN deployment intervals was short, it was caused by the environment there are various noises. However, the throughput was stable to more than 1.0 [Mbps] in this environment, it was confirmed the availability of our proposed method in this field test.

## VI. CONCLUSION

This paper proposed the WSN deployment strategy that maintains throughput more than 1.0 [Mbps]. The proposed strategy maintained communication conditions such that the throughput between End to End communications in the WSN enables smooth tele-operation of the mobile rescue robot in a post-disaster underground space. Experimental results showed the effectiveness of the proposed strategy that is enable to construct the WSN in the field test.The rapid implementation of actions to reduce secondary disasters in disaster areas requires the stable referral of disaster information. Therefore, this strategy which constructs WSN that maintains the throughput stable by utilizing rescue robot is effective for gathering

disaster area information in actual disaster scenarios. We will apply the proposed strategy to WSN deployment in practical underground space in the future.

## REFERENCES

[1] CHI Hao-yuan, LIU Xu, XU Xiao-dong, "A Framework for Earthquake Disaster Mitigation System," Proceedings of 2011 China located International Conference on Information Systems for Crisis Response and Management (ISCRAM), pp.490-495, 2011.

[2] Huang AN, "China's Emergency Management Mechanisms for Disaster Prevention and Mitigation," Proceedings of International Conference on E-Business and E-Government (ICEBEG), pp.2403-2407, 2010.

[3] Yoshiaki KANAEDA, Kazushige MAGATANI, "Development of the device to detect SPO2 in the Field,"31st Annual International Conference of the IEEE EMBS, pp.412-415, September 2009.

[4] Y. Kawata, "Disaster Mitigation due to next Nankai earthquake tsunami occurring in around 2035," Proceedings of International Tsunami Symposium 2001, session 1, pp. 315-329, 2001.

[5] Y. Kawata, "The great Hanshin-Awaji earthquake disaster: damage, social response, and recovery," Journal of Natural Disaster Science, Vol. 17, No. 2, pp.1-12, 1995.

[6] H. Kawakata, Y. Kawata, H. Hayashi, T. Tanaka, K. C. Topping, K. Yamori, P. Yoshitomi, G. Urakawa and T. Kugai, "Building an integrated database management system of information on disaster hazard, risk, and recovery process," Annuals of Disas. Prev. Res. Inst., Kyoto Univ., No.47 C, 2004.

[7] Abishek T K, Chithra K R and Maneesha V. Ramesh, "ADEN:Adaptive Energy Efficient Network of Flying Robots Monitoring over Disaster Hit Area," Proceedings of 8th IEEE International Conference on Distributed Computing in Sensor Systems (IEEE DCOSS), pp.306-310, 2012.

[8] Abishek T K, Chithra K R, Maneesha V Ramesh, "AER: Adaptive Energy Efficient Routing Protocol For Network of Flying Robots Monitoring over Disaster Hit Area," Proceedings of 21st Annual Wireless and Optical Communi-cations Conference (WOCC), pp.166-169, 2012.

[9] Parker, E., L., Kannan, B., Xiaoquan, F., Yifan, T. (2003). Heterogeneous Mobile Sensor Net Deployment Using Robot Herding and Line of Sight Fromations, Proceedings of 2003 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS2003), Volume 3. pp.2488-2493, 2003.

[10] T. Umeki, H. Okada, K. Mase, "Evaluation of Wireless Channel Quality for an Ad Hoc Network in the Sky SKYMESH," Proceedings of Sixth International Symposium on Wireless Communication Systems 2009 (ISWCS'09). pp.585-589, 2009.

[11] Helge-Bjorn Kuntze, Christian W. Frey, Igor Tchouchenkov, Barbara Staehle, Erich Rome, Kai Pfeiffer, Andreas Wenzel and Jurgen Wollenstein, "SENEKA - Sensor Network with Mobile Robots for Disaster Management," Homeland Security (HST), pp.406-410, 2012.

[12] E. Budianto, M.S. Alvissalim, A. Hafidh, A. Wibowo, W. Jatmiko, B. Hardian, P. Mursanto and A. Muis, "Telecommunication Networks Coverage Area Expansion in Disaster Area using Autonomous Mobile Robots : Hardware and Software Implementation," Proceedings of International Conference on Advanced Computer Science and Information Systems (ICACSIS), pp.113-118, 2011.

[13] Andrew Chiou, and Carol Wynn, "Urban Search and Rescue Robots in Test Arenas: Scaled Modeling of Disasters to Test Intelligent Robot Prototyping," Proceedings of International Conference on Autonomic and Trusted Computing (ATC), pp.200-205, 2009.

[14] A. Howard, J. Matric, and J. Sukhatme, "Mobile Sensor Network Deployment using Potential Fields," A Distributed Scalable Solution to the Area Coverage Problem, Distributed Autonomous Robotics Systems 5, Springer-Verlag. pp.299-308, 2002.

[15] Wing-Yue Geoffrey Louie, and Goldie Nejat, "A victim identification methodology for rescue robots operating in cluttered USAR environments," Advanced Robotics, vol. 27, issue. 5, pp. 373-384, 2013.

[16] Andrew Markham and Niki Trigoni, "Magneto-Inductive NEtworked Rescue System (MINERS):Taking Sensor Networks Underground," Proceedings of the 11th international conference on Information Processing in Sensor Networks (IPSN '12), pp. 317-328, 2012.

[17] Josh D. Freeman, Vinu Omanan, and Maneesha V. Ramesh, "Wireless Integrated Robots for Effective Search and Guidance of Rescue Teams," Proceedings of 8th International Conference on Wireless and Optical Communications Networks (WOCN 2011), pp. 1-5, 2011.

[18] K. Sawai, T. Suzuki, H. Kono, Y. Hada and K. Kawabata, "Development of a SN with impact-resistance capability for gathering disaster area information," 2008 Internatio-nal Symposium on Nonlinear Theory and its Applications (NOLTA2008), pp.17-20, 2008.

[19] Tsuyoshi Suzuki, Kei Sawai, Hitoshi Kono and Shigeaki Tanabe, "Sensor Network Deployment by Dropping and Throwing Sensor Node to Gather Information Underground Spaces in a Post-Disaster Environment," Discrete Event Robot, iConcept PRESS, in Press. 2012.

[20] K. Sawai, H. Kono, S. Tanabe, K. Kawabata, T. Suzuki, "Design and Development of Impact Resistance Sensor Node for Launch Deployment into Closed Area," In international journal of sensing for industry(Sensor Review), Emerald Group Publishing Ltd., Vol. 32, pp.318 – 326, 2012.

[21] S. Tanabe, K. Sawai and T. Suzuki, "Sensor Node Deployment Strategy for Maintaining Wireless Sensor Network Communication Connectivity," International Journal of Advanced Computer Science and Applications ( IJACSA ) , The Science and Information organization, Vol.2, No. 12, pp.140 – 146, 2011.

[22] H. Sato, K. Kawabata and T. Suzuki, "Information Gathering by wireless camera node with Passive Pendulum Mechanism," International Conference on Control, Automation and Systems 2008 (ICCAS2008), pp.137-140, 2008.

[23] T. Yoshida, K. Nagatani, E. Koyanagi, Y. Hada, K. Ohno, S. Maeyama, H. Akiyama, K. Yoshida and Satoshi Tadokoro, "Field Experiment on Multiple Mobile Robots Conducted in an Underground Mall," Field and Service Robotics Springer Tracts in Advanced Robotics, vol. 62, pp365-375, 2010.

[24] H. Jiang, J. Qian, and W. Peng, "Energy Efficient Sensor Placement for Tunnel Wireless Sensor Network in Underground Mine," Proceedings of 2nd International Conference on Power Electronics and Intelligent Transportation System(PEITS 2009), pp. 219-222, 2009.

[25] J. Xu, S. Duan and M. Li, "The Research of New Type Emergency Rescue Communication System in Mine Based on Wi-Fi Technology," Proceedings of IEEE 3rd International Conference on Communication Software and Networks (ICCSN), pp. 8-11, 2011.

[26] K. Nagatani, S. Kiribayashi, Y. Okada, K. Otake, K. Yoshida, S. Tadokoro, T. Nishimura, T. Yoshida, E. Koyanagi, M. Fukushima and S. Kawatsuma, "Emergency Response to the Nuclear Accident at the Fukushima Daiichi Nuclear Power Plants using Mobile Rescue Robots," Journal of Field Robotics, vol. 30, no. 1, pp. 44-63, 2013.

# A Secured Framework for Geographical Information Applications on Web

Mennatallah H. Ibrahim, Hesham A. Hefny
Computer Science and Information
Institute of Statistical Studies and Research, Cairo University
Cairo, Egypt

*Abstract*—Current geographical information applications increasingly require managing spatial data through the Web. Users of geographical information application need not only to display the spatial data but also to interactively modify them. As a result, the security risks that face geographical information applications are also increasing. In this paper, a secured framework is proposed. The proposed framework's goal is, providing a fine grained access control to web-based geographic information applications. A case study is finally applied to prove the proposed framework feasibility and effectiveness.

*Keywords—spatial data; geographic information systems; access control; authorization*

## I. Introduction

Geographic information applications on web strongly need to be secured for several factors: (1) most of the geographical data contain sensitive information, so data cannot be freely disclosed to or altered by illegitimate users, (2) geographical information application's users have different roles and expertise, so they need to be assigned different rights for operating on data, and (3) the power of geographical information applications comes from its ability to relate different types of data in a spatial context, in which these related data are supplied by different data providers such as governments, private companies, academic organizations, and so forth. Each data provider need to ensure the protection of its own data when published on the Web. However, security of geographic information applications on web is an issue that has not been much investigated by Geographic Information Systems (GIS) community.

Applying a controlled access to geographic information applications on web is one of the most important security aspects to such applications. Controlled access ensures the information confidentiality and integrity. Since the geographic information applications are in most cases critical and complex applications (e.g., health applications) contextual permissions are needed , such permissions give concrete rights to users for operating on data in specific situations (e.g., doctors may be given concrete permission to operate on patient's data in emergency situations).

In general, there are two main models used for digital spatial data: (1) vector data models that use discrete elements such as points that represent small objects, lines that represent linear objects and boundaries, and polygons that represent areas. These three elements used to represent the geometry of real world entities, (2) raster data models which identifies and represents grid cells for a given region of interest, raster cells are arrayed in a row/column pattern cell values represent type or quality of mapped variables used with values that may change continuously across a region: elevation, mean temperature, average rainfall [8].

Fine-grained access control and spatial or non-spatial access control are two important requirements for spatial data access control. Fine-grained access control is important where; the spatial data in database usually have different granularities, which are organized in hierarchical architecture. The hierarchy from top to down has two representations in which one is using terms of database, namely, tables, records and cells, while the other is using terms of geospatial domain, namely, map layer, geospatial objects, geometric or descriptive properties. Spatial or non-spatial access control is also very important where; restricting access to some spatial objects, whose descriptive properties meet some conditions, is frequently needed. For example, those spatial objects, whose type is military unit, cannot be accessed by ordinary users.

The aim of this paper is to propose a framework that provides access control to web-based geographic information applications depending on Organization Based Access Control (ORBAC) model [9]. The proposed framework supports new concepts that were not addressed before in order to provide more security to the geographic information applications on Web. These concepts are: (1) contextual permissions, and (2) supporting various security polices in a unique framework. The proposed model also has an important advantage which is achieving fine grained access control through views. Thus, users of the proposed framework are not allowed to access database tables, instead, they are only allowed to access views of these tables. The proposed framework deals with vector data models where, vector data models are more adequate for usage in current GIS applications and spatial database management systems, also vector data models are more adequate for dynamic applications that require data modifications.

## II. Related Work

The pioneer access control model for vector-based spatial data on web is proposed in [1, 2]. This pioneer model is based on Role-Based Access Control (RBAC). It extends the classical discretionary access control model in which it adds a spatial dimension to the authorization rules by assigning a geographical scope in which this geographical scope defines the spatial region in which the authorization is valid. When an

access request is issued for an object, the system checks if the requested object lies in the authorization space and if so, it grants the access. A similar architecture but differs in focusing on XML-based representation of spatial data, has been proposed in [3]. The main limitation in such models is represented in, not addressing the issue of multi granularity of spatial data. This limitation has been addressed in [4] where, a more complex spatial data model has been proposed in which, the specification of authorization rules to access complex structured spatial data stored in a DBMS is allowed and organized according to multiple spatial representation levels and at multiple granularities. The model proposed in [4], however, does not deal with geographically bounded roles.

A fine-grained access control model based on RBAC for grid environment is proposed in [5]. The proposed model is based on Globus Security Infrastructure, in such model, every user is mapped to a given role, and every role is given a unique digital certificate to distinguish its identification, then every role had the given permission to access the resources. The main advantage of the model proposed in [5] is represented in, providing strong access control through digital certifications. In [6], also a fine-grained access control model based on RBAC to spatial data in grid environment is proposed. The proposed model adopts a double authorization mechanism: the first authorization authorizes the role, similarly to the RBAC model, and the second authorization authorizes the specific user based on the user's attribution. The limitations of such model are: (1) the role authorization is achieved through Access Control List which is a time consuming method when the resources are massive, (2) the fine-grained authorization method is complex, and (3) conflicts may occur between both the role and the fine-grained authorizations.

As mentioned before, the power of geographical information applications comes from its ability to relate different types of data in a spatial context, such related data are provided by different data providers, and as a result each data provider needs to apply its own security police to ensure its data security while sharing them over the web. Furthermore geographic information applications are complex and critical; such applications strongly need contextual permissions. All of the previously proposed models are based on RBAC model. RBAC models are not fully satisfactory for web-based geographic information applications, as they are not supporting the authorizations rules that specify contextual permissions, or the rules that are specific to particular organization. In another word, none of the previous access control models is able to model security policies that are not restricted to static permissions or to support various security polices in a unique framework.

## III. ORGANIZATION BASED ACCESS CONTROL MODEL

ORBAC defines permissions that are applied within an organization to control the activities performed by roles on views under specific conditions. In ORBAC, the subject must be assigned to a given role, the object must be used in a given view and the action must partake in some activities. There are eight basic sets of entities in ORBAC which are: Org (a set of organization), S (a set of subjects), α (a set of actions), O (a set of objects), R (a set of roles), a (a set of activities), V (a set of views) and C (a set of contexts). The following are the basic entities of ORBAC model:

### A. Organizations

The Organization is the most important entity in ORBAC model. An organization can be seen as an organized group of subjects who agreed to form an organization. Subject must plays specific roles in the organization according to their agreement.

### B. Subjects and Roles

A subject is an active entity (e.g., a user). A role is used to create a link between subjects and organizations. If org is an organization, s is a subject and r is a role, then Employ (org; s; r) means that org employs subject s in role r.

### C. Objects and Views

The entity Object covers inactive entities (e.g., database tables). As in relational databases, a view corresponds to a set of objects that satisfy a common property. If org is an organization, o is an object and v is a view, then Use (org; o; v) means that org uses object o in view v.

### D. Actions and Activities

The entity Action represents computer actions such as read, write, send, and so forth. The entity Activity is used to abstract actions. If org is an organization, α is an action and a is an activity, then Consider (org; α; a) means that org considers that action α falls within the activity a.

### E. Context

Context is used to specify the concrete circumstances where organizations grant roles permissions to perform activities on views. If org is an organization, s is a subject, o is an object, α is an action and c a context, then Define (org; s; o; α; c) means that within organization org, context c is true between subject s, object o and action α. The conditions required for a given context to be linked, within a given organization, to subjects, objects and actions is formally specified by logical rules.

In ORBAC Security Policy, the relationship "Permission" corresponds to a relation between organizations, roles, views, activities and contexts, also the relationships Prohibition, Obligation and Recommendation are defined similarly. If org is an organization, r is a role, v is a view, a is an activity and c a context then Permission (org; r; v; a; c) means that organization org grants role r permission to perform activity a on view v within context c.

## IV. VIEWS AND FINE GRAINED ACCESS CONTROL

A view is a logical representation of database table(s). In essence, a view is a stored query with no physical storage. A view derives its data from database tables on which it is based; such tables are called base tables. All operations performed on a view affect the base tables. Views provide an additional level of table security by restricting access to a predetermined set of rows or columns of a table, it enable one to tailor the presentation of data to different types of users [7].

With respect to security, we usually want to let specific users access some columns and rows of base tables while hiding other sensitive ones. In our proposed framework, views provide a solution for realizing fine-grained access control for spatial database. By creating views, table level (map layer level), record level (feature level), field level (property level) or even spatial context access control can be easily implemented.

## V. PROPOSED FRAMEWORK

The proposed framework aims to: provide fine-grained access control to geographic information applications on web. The proposed framework based on ORBAC model which provides mean to specify different security policies within a unique framework, so that each organization providing data to build the geographic application will be able to define its own security policy, furthermore organizations will be able to specify contextual permissions.



Fig. 1.   "Proposed Framework"

As shown in Fig.1, the proposed framework depends on the well known three-tier architecture which consists of three layers: (1) Presentation layer; it resides on the users side and consists of either html pages or specialized programs, such as Java code, and plugs in. Users interact with the web-based geographic application through the presentation layer by sending requests and receiving responses, (2) Application layer consists of three services which are: (a) Access Control Service that exposes and implements the operations for both authorization rules checking and administration, (b) Application Service that exposes and implements the application logic and access the application data, and (c) Authentication Service. (3) The data Storage layer consists of database servers.

The proposed framework consists of the following components:

1) *User accounts with usernames and passwords;*
2) *Organizations that users belonging to;*
3) *Roles that users plays in their organizations;*

4) *Actions that is assigned to users according to their roles;*
5) *Views which includes spatial objects that have common properties;*
6) *Contexts in which the roles permitted to perform actions*
7) *Database tables that contain spatial and non-spatial data about the spatial objects.*

In our proposed framework, the views are created from the base tables in the spatial database according to different access control requirements: either spatial, non-spatial or their combination. These created views are granted to different roles with corresponding authorized actions. The views are granted to users according to their roles and organizations. The proposed framework will support the fine grained access control through these created views.

The typical interaction between the user and the system is as follows: (1) the user will connect to the system through the Authentication Service, (2) an organization and a role will be assigned to the authenticated user, (3) The user will be allowed to issue request(s) to the system through the presentation layer. Each request from the user is then mapped onto one or more operations of the Application Service. The application service in turn interacts with the Access Control Service to verify whether the operation can be performed or not, and (4) the response is sent back to the user from the application service through the presentation layer.

## VI. CASE STUDY

A case study has been carried in order to prove the feasibility and effectiveness of the framework proposed in this paper. This case study is represented in creating a web-based geographic information application for two business organizations dealing with each other. These two organizations need to collaborate together in order to enhance their performance level and such enhancement will be achieved by creating a web-based geographic information application. The first organization owns a number of warehouses for electronic products, while the second organization owns a number of stores that sell the products of the first organization. The goal of the application is to maintain the warehouses, the stores and their products over the Web. Data in such application can be queried, inserted and modified using a Web browser.



Fig. 2.   "AllWarehouses" View on map that retrieves all warehouses

Fig. 3.   "AllStores" View on map that retrieves all stores



Fig. 4.   "MidAmericaWarehouse" View on map that retrieves only Mid America Warehouse.

The oversimplified information which defines the access control policy is assuming that there are two roles which are Manager and Coordinator. These roles belong to two organizations which are, Organization1 and Organization2. The features to be secured are the warehouses and the stores. The privileges are RetrieveData, InsertData, UpdateData and DeleteData. The contexts in which privileges will be granted to roles are normal and emergency.

Three views are created: (1) "AllWarehouses" view which is created from three base tables. Such view contains data about all warehouses owned by Organization2, (2) "AllStores" view which is also created from three base tables and contains data about all stores owned by Organization1, and (3) "MidAmericaWarehouse" view which is a subset of "AllWarehouses" view that contains data about only one warehouse called Mid America warehouse.

Let O be the set of organizations, R be the set of roles, V be the set of views, A be the set of actions and C be the set of contexts:

O= {Organization1, Organization2}
R= {Manager, Coordinator}
V= {AllWarehouses, AllStores, MidAmericaWarehouse}
A= {RetrieveData, InsertData, UpdateData, DeleteData}
C= {Normal, Emergency}
Note: The keyword ALL stands for all possible values for the field.

The authorization rules are defined as follows:
r1= <Organization1, Manager, AllStores, ALL, ALL >
r2= < Organization2, Manager, AllWarehouses, ALL, ALL>

r3= < Organization2, Coordinator, MidAmericaWarehouse, ALL, ALL>
r4= < Organization2, Coordinator, AllWarehouses, RetrieveData, Emergency >

**Rule r1**: states that the role "Manager" in "Organization 1" is authorized to retrieve, insert, update and delete data in the "AllStores" view in all contexts.
**Rule r2:** states that the role "Manager" in "Organization 2" is authorized to retrieve, insert, update and delete data in the "AllWarehouses" view in all contexts.
**Rule r3**: states that role "Coordinator" in "Organization 2" is authorized to retrieve, insert, update and delete data in the view "MidAmericaWarehouse" in all contexts.
**Rule r4:** states that the role "Coordinator" in "Organization 2" is authorized **only** to retrieve data in the "AllWarehouses" view and this is **only** in Emergency context.

From the previous rules, it is clear that by using views fine granularity and is achieved. Fine granularity is clearly shown in the "role Coordinator" in "Organization1" which is operating on data of a certain warehouse (i.e., Mid America Warehouse), but not operating on data of all other warehouses. The previous rules also show that using the entity context enable us to control when to allow a specific role (i.e., Coordinator) to perform a specific privilege (i.e., RetrieveData) on specific view (i.e., AllWarehouses) in specific context (i.e., Emergency) In normal cases such role is allowed only to deal with "MidAmericaWarehouse" view.

The effect of authorization rules and views on different user's interactions are illustrated through a number of screen shots. "Fig. 2" shows the "AllWarehouses" view which is displayed to the role "Manager" in "Organization2" as this role is allowed to retrieve all warehouses. "Fig. 3" shows the "AllStores" view which is displayed to the role "Manager" in "Organization 1" as this role is allowed to retrieve all stores. While "Fig. 4" shows the "MidAmericaWarehouse" view which is displayed to the role "Coordinator" as this role is allowed to retrieve only Mid America Warehouse except in Emergency context in which such role will be allowed to retrieve all warehouses.

Such results cannot be achieved with the previously proposed models as such models don't allow any organization as a data provider to specify its own security policy as a result all users who assigned the role "Manager" will be able to access same views regardless the organizations that these users belong to. Also in the previous models, contextual permissions are not allowed as a result the role "Coordinator" will be allowed either to retrieve all warehouses data or not without considering any contexts. Furthermore, the previously proposed models allow all authenticated users to access base tables of database while in our proposed framework all users are prevented from accessing those tables.

By comparing the proposed framework and the previously proposed ones, it is clear that, the proposed framework supports two new important concepts that were not supported previously, although these new concepts are very important to the web-based geographic information applications. The first

concept presented in, the ability of the proposed framework to provide a mean to specify different security policies within a unique framework. So each data provider that collaborated in building the application can define its security policy to protect its own data. The second concept presented in, the ability to specify contextual permissions. Furthermore, the proposed framework supports the fine-grained access control through views. Furthermore, an important difference between the proposed model in this paper and the previously proposed models is that the model proposed in this paper prevents authenticated users from accessing base tables, as users are allowed only to access specific views according to their roles and organizations. TABLE I summarizes the differences between the model proposed in this paper and the previously proposed ones.

TABLE I.        PROPOSED MODEL AND PREVIOUS MODELS COMPARISON

| Criteria | Frameworks | |
|---|---|---|
| | *Previous Models* | *Proposed Model* |
| Allowing each data provider to define its security policy to protect its own data | – | ✓ |
| Preventing access to base tables | – | ✓ |
| Allowing Contextual Permissions | – | ✓ |

The proposed framework advantages:

*1)    The framework depends on ORBAC model which provides means to specify different security policies within a unique framework.*

*2)    The framework supports fine granularity access control to data through views.*

*3)    The framework provides authorization rules that specify contextual permissions.*

The proposed framework limitations:

*1)    Concurrent control:* The multi views mode to a single base table may cause the concurrent control problem and this happens when many users access those views that based on the same base table concurrently. Fortunately, this problem has been solved to some extent by the database internal mechanism.

*2)    Redundancy and information leakage:* The abusive usage of views can result in redundancy of the access control predicates, and the potential of information leakage through exceptions and errors that are caused by user-defined functions.

## VII.    CONCLUSION

Security is very important and critical for web-based geographic information applications. Access control is required to keep data confidentiality and integrity. In this paper, we proposed a framework for secured web-based geographic information applications based on ORBAC model. The goal of this framework is to ensure the security of data in such applications.

REFERENCES

[1]    E. Bertino & M. L. Damiani. *"A controlled access to spatial data on Web"*. 7th AGILE Conference on Geographic Information Science, Heraklion, Greece. 29 April-1May 2004. Pages 369-377.J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.

[2]    E. Bertino, M.L. Damiani, & D. Momini. *"An Access Control System for a Web Map Management Service"*. In Proc. of the 14th International Workshop on Research Issues in Data Engineering (RIDE-WS-ECEG), Boston, USA. March 2004. Pages 33–39.

[3]    B. Purevjii, T. Amagasa, S. Imai & Y. Kanamori. *"An access control model for geographic data in an XML-based framework"*. In Proc of the 2ⁿᵈ International Workshop on Information Systems Security (WOSIS). 2004. Pages 251-260.

[4]    A. Belussi, E. Bertin, B. Catania, M.L. Damiani & A. Nucita. *"An authorization model for geographical maps"*. Proceedings of the 12th annual ACM international workshop on Geographic information systems, New York, NY, USA. 2004. Pages 82-91.

[5]    M. Yan, Y. Gao, L. Wu, P.Wu, & Y. Zhao. *"Spatial data access control in grid environment"*. Geoinformatics, 17th International Conference. August 2009. Pages 1-6

[6]    F. Ma, Y. Gao, M. Yan, F. Xu & D. Liu. "The fine-grained security access control *of spatial data"*. Geoinformatics 18th International Conference. June 2010.  Pages 1-4.

[7]    http://docs.oracle.com/cd/E18283_01/server.112/e16508.pdf

[8]    P. Bolstad. April 2012.*"GIS fundamentals: A firat text on geographic information systems"*.4ᵗʰ Edition.U.S. *State of* of Minnesota: Eider Press.

[9]    A. Abou El Kalam, S. Benferhat, R. El Baida, A. Miege & et al. "*Organization based*  access control". Proceeding of the 4ᵗʰ IEEE International Workshop on Policies for Distributed Systems and Networks. 2003. Page 120.

# Bimodal Emotion Recognition from Speech and Text

Weilin Ye
Institute of Computer Science and Technology
Chongqing University of Posts and Telecommunications
Chongqing, China

Xinghua Fan
Institute of Computer Science and Technology
Chongqing University of Posts and Telecommunications
Chongqing, China

*Abstract*—**This paper presents an approach to emotion recognition from speech signals and textual content. In the analysis of speech signals, thirty-seven acoustic features are extracted from the speech input. Two different classifiers Support Vector Machines (SVMs) and BP neural network are adopted to classify the emotional states. In text analysis, we use the two-step classification method to recognize the emotional states. The final emotional state is determined based on the emotion outputs from the acoustic and textual analyses. In this paper we have two parallel classifiers for acoustic information and two serial classifiers for textual information, and a final decision is made by combing these classifiers in decision level fusion. Experimental results show that the emotion recognition accuracy of the integrated system is better than that of either of the two individual approaches.**

*Keywords—emotion recognition; acoustic features; textual features; decision level fusion*

## I. Introduction

With the advent of information age and the popularity of Internet, more and more kinds of information come to our life. Phoning has become the main means of daily communication and follow-up contacting. We often play some customer service phone to ask for information about some products, after the call we always asked to evaluate the service attitude of the telephone operator, so that the businesses can know the service quality of the staff. However, manual evaluation often has a problem of objectivity and authenticity. Automatic emotion recognition is one of the key techniques of human-computer interaction [1].

In recent years, several research works have focused on emotion recognition. Hoch et al[2] presented a method to recognize three kinds of emotional states in the automotive environment from speech and expression information. Busso et al[3] analyzed the complementarity of speech emotion recognition and facial expression recognition, presented a multi-modal emotion recognition method from feature level fusion and decision level fusion. Wangner et al[4] combined electromyogram, ECG, skin resistance and breathing these four kinds of physiological parameters to recognize emotional state and got a recognition rate of 92%.

However, few approaches have focused on emotion recognition from textual input. Textual information is another important communication medium and can be retrieved from many sources, such as books, newspapers, web pages, e-mail messages, etc. It is not only the most popular communication medium, but also rich in emotion. With the help of natural language processing techniques, emotions can be extracted

from textual input. In this paper, a bimodal emotion recognition method is used to extract emotion information from both speech and text input. In this paper, the classifiers recognize emotions according to two simple types: positive and non-positive. This paper designed two parallel classifiers for acoustic information and two serial classifiers for textual information, and a final decision is made by combing these classifiers in decision level fusion. This method can be applied to a telephone service center dialogue system to recognize customers' negative emotions, such as anger, impatience etc. so that to turn the answering service to manual service automatically to avoid losing customers.

## II. Emotional Speech Corpus

At present, there still not have a public database for Chinese speech emotion recognition research. Generally there are two ways to get emotional speech corpus: a) Recording; b) Clipping. Recording method has better customization, and can record emotional speech which meets the speaker, text, emotion categories and other requirements. According to the general rules of building corpus, four college students around the age of 20 with higher emotional expression ability are invited to participate in recording (2 females, 2 males). After five non-recording people's perception experiments, we removed nearly 40% corpus which are not sure which kind of emotion. Finally we picked out a total of 600 available corpuses, including positive and non-positive each 300, where non-positive include anger, sadness, fear and other negative emotions.

## III. Preprocessing

The purpose of voice and text preprocessing are different. Voice preprocessing is to get pure voice by eliminating the interference of various factors. Text preprocessing is to get relatively clean data sets by filtering noise data.

### A. Speech Signal Preprocessing

#### 1) Pre-emphasis
Since speech signal are affected by the glottis excitation and snout radiation, the high frequency part of the speech signal falls down. Pre-emphasis enhance the high frequency part, make the signal spectrum flat over the entire frequency band.

#### 2) Window Function
Commonly used window functions in voice processing are rectangular window and hamming window.

##### a) Rectangular Window:

$$w(n) = \begin{cases} 1 & (0 \le n \le N-1) \\ 0 & Other \end{cases}$$

b)

(1)

c) *Hamming Window:*

$$w(n) = \begin{cases} 0.54 - 0.46\cos\left[\dfrac{2\pi n}{N-1}\right] & (0 \le n \le N-1) \\ 0 & Other \end{cases} \quad (2)$$

### B. Text Preprocessing

Firstly we use the Chinese auto-segmentation system to do the process of word segmentation, and then move the stop words from target text. Finally we can get relatively clean data sets. The process of text pretreatment is showed in figure1.



Fig. 1. Flow chart of text preprocessing

### IV. EMOTIONAL FEATURE EXTRACTION

#### A. Acoustic Features

Speech signal is short-time stationary, calculating features based on short time frame are: short-time amplitude energy, short pitch and first-three formant. For the whole speech, every feature is calculated as one-dimensional sequence. However these sequences cannot be directly used as a feature vector for pattern recognition, commonly way to solve this problem is to calculate its statistical value, such as mean and slope.

Speech emotion recognition based on prosodic features has strong robustness and adaptability. Statistical characteristics can better reflect the rhythmic structure of speech. On the basis of previous experiment, we chose 37 identification features, which are shown in Table 1.

#### B. Textual Features

##### 1) Feature extension

Text orientation classification is different from general classification, words or phrases with semantic orientation or emotional tendencies play a crucial role for classification. In this paper we use three-step feature extension [5] to reconstruct the data sets, which can extent features of the data sets by using list of tendency words, negative words and degree adverbs. This method can enhance expression ability of the textual features by adding words or phrases with semantic orientation to feature sequence.

##### 2) Feature Selection

Commonly used feature selection methods are: document frequency, mutual information, information gain, expects cross-entropy, chi-square statistics etc.

TABLE I ACOUSTIC FEATURES

| Feature | Describe | Feature | Describe |
|---|---|---|---|
| 1 | Maximum energy | 20 | Mean duration of pitch frequency contour ascent |
| 2 | Mean value of energy | 21 | Maximum of pitch frequency contour decline |
| 3 | Median value of energy | 22 | Mean value of pitch frequency contour decline |
| 4 | Rate of change of energy | 23 | Maximum duration of pitch frequency contour decline |
| 5 | Maximum of energy contour ascent | 24 | Mean duration of pitch frequency contour decline |
| 6 | Mean value of energy contour ascent | 25 | Maximum of the first formant |
| 7 | Maximum duration of energy contour ascent | 26 | Mean value of the first formant |
| 8 | Mean duration of energy contour ascent | 27 | Median value of the first formant |
| 9 | Maximum of energy contour decline | 28 | Rate of change of the first formant |
| 10 | Mean value of energy contour decline | 29 | Maximum of the second formant |
| 11 | Maximum duration of energy contour decline | 30 | Mean value of the second formant |
| 12 | Mean duration of energy contour decline | 31 | Median value of the second formant |
| 13 | Maximum of pitch frequency | 32 | Rate of change of the second formant |
| 14 | Mean value of pitch frequency | 33 | Maximum of the third formant |
| 15 | Median value of pitch frequency | 34 | Mean value of the third formant |
| 16 | Rate of change of pitch frequency | 35 | Median value of the third formant |
| 17 | Maximum of pitch frequency contour ascent | 36 | Rate of change of the third formant |
| 18 | Mean value of pitch frequency contour ascent | 37 | Speed(voice frames / statement of words) |
| 19 | Maximum duration of pitch frequency contour aascent | | |

However these methods are not much suitable for text orientation classification. In this paper we chose the document frequency feature selection formula presented in literature [6] which considered the words tendentiousness.

$$DF\_Sen(t,c) = \frac{\lg(DF_t)}{\lg(N_c)} * \frac{\alpha_t(|\beta_t|+1)}{\alpha_t + \gamma} \quad (3)$$

Among it, $DF_t$ means the number of documents showed in class $c$ of feature $t$, $N_c$ means the whole numbers of documents in class $c$, $\beta_t$ means the intensity values of orientation, $\alpha_t$ means the number of words feature $t$ contains, $\gamma$ means the weighing coefficient which can be adjusted in experiment. When selecting parameter, we set a threshold $DF\_SEN_{min}$. If the threshold of a feature is less than a certain value, it will be deleted. In this paper, based on experiments we select 0.04 as the value of threshold. When a feature word appeared at multiple classes, we selected it according to its feature score in

each category. If the absolute value of difference value of feature score in two deferent categories is more than 0.12, the word will be selected as textual feature.

## V. BIMODAL FUSION RECOGNITION ALGORISM

Currently, there are two ways to combine different pieces of information: a) feature level fusion, b) decision level fusion. The problem with feature level fusion is the potential of having to face the nurse of dimensionality due to the increase in the input feature dimension. In our case, we have two parallel classifiers for acoustic information and two serial classifiers for textual information, and a final decision is made by combing the outputs results from these classifiers in decision level fusion.

### A. Classification of the Acoustic Set

In this paper, we have two parallel classifiers for acoustic information. They are support vector machines (SVM) and BP neural nets.

#### 1) Support Vector Machines

A great interest in Support Vector Machines (SVM) in classification can be observed recently. They tend to show a high generalization capability due to their structural risk minimization oriented training. Non-liner problems are solved by a transformation of the input feature vectors into a generally higher dimensional feature space by mapping function where linear separation is possible. Maximum discrimination is obtained by an optimal placement of the separation plane between the border of two classed.

#### 2) BP Neural Nets

The BP Neural Nets is proposed by a team of scientists led by Rumelhart and McCelland, and it is one of the most widely used neural network model. BP network can learn and store a large amount of mapping relationship of input-output model without pre-revealing the mathematical equations that describe the mapping relationship.

### B. Classification of the Textual Set

In this paper, we use the two-step classification proposed in literature [7]. We construct two serial classifiers CF1 and CF2, both of them use equation (4) to select features. Firstly use CF1 for classification. For unreliable part of the classification results, we use CF2 for secondary classification.

#### 1) To construct classifier CF1

CF1 is Naive Bayes classifier. Text $d$ is expressed as $d = (t_1, t_2, K, t_n)$, $t_k$ is feature item of the text. Then Naive Bayes classifier is expressed as follows:

$$P(C_i \mid d) = \max(P(d \mid C_i)P(C_i)), i = 1, 2 \qquad (4)$$

In formula (4), $P(d \mid C_i) = \prod_{j=1}^{n} P(t_j \mid c_i)$

#### 2) To construct classifier CF2

Classifier CF2 is expressed as follows:

$$f(d) = \sum_{i=1}^{n} \left( \left| \frac{\lg(P(t_k \mid C_1) / P(t_k \mid C_2))}{\lg(P(t_k \mid C_1)P(t_k \mid C_2))} \right| * \lg(\frac{P(t_k \mid C_1)}{P(t_k \mid C_2)}) + \lambda Q(t_k) \right) \qquad (5)$$

Where $c_1$ represents positive emotion, $c_2$ represents non-positive emotion, $Q(t_k)$ represents the intensity values of feature $t_k$, $Q(t_k) > 0$ means it's a positive word, $Q(t_k) < 0$ means it's a negative word. If feature $t_k$ is not in the tendency word list, then $Q(t_k) = 0$, $\lambda$ is adjustment coefficient, if $f(d) > 0$, then text $d$ is positive, otherwise text $d$ is non-positive.

### C. Fusion Algorism in Decision Level

In this paper we construct two parallel classifiers for acoustic information and two serial classifiers for textual information, and a final decision is made by combing these classifiers in decision level fusion. Flow chat is showed in figure2.



Fig. 2. Fusion algorism in decision level

Assuming that $w_1$, $w_2$, $w_3$ each represents the weight of SVM, BP neural nets and textual classifier. The right value of $w_1$, $w_2$, $w_3$ plays a very important role to the fusion result. As different classifier has a different recognition rate, a number of samples extracted from the training set as a check set. The value of $w_1$, $w_2$, $w_3$ is determined by the recognition rate of the check set. $P_1$, $P_2$, $P_3$ each represents the recognition rate of SVM, BP neural nets and textual classifier. Then the value of $w_i$ is calculated as follows:

$$w_i = \frac{P_i}{\sum_{i=1}^{3} P_i} \qquad (6)$$

The idea of weighted score voting strategy: if the three classifiers have the same recognition result, then the sample will be identify as such class; if two of the three classifiers have the same recognition result, then sum the two weights of the classifiers with the same result, and compare it with the weight of the classifier which has a different result, then the sample will be identify as class recognized by the classifier with bigger weight.

## VI. EXPERIMENT AND RESULTS

Both of the training data sets and the testing data sets have two kinds of emotion: positive and non-positive. Each training set contains every emotion 200 speech samples and 200 text samples, each testing sets contains every emotion 100 speech samples and 100 text samples. In experiments, we use the same training sets and testing sets to test every single model classifier.

In this paper, we use 2*2 confusion matrix to evaluate the emotion recognition algorism. The element in row i and column j means the proportion that the real emotion state i is recognized as j. That is to say, the greater values on the diagonal matrix are, the better effect of the emotion recognition algorism is.

Experiment 1: recognition rate of single-mode SVM classifier based on acoustic features is shown in Table 2.

TABLE II        ACCURACY OF SVM (%)

| Sample Sets | Positive | Non-positive |
|---|---|---|
| Positive | 82 | 18 |
| Non-positive | 16 | 84 |

Experiment2: recognition rate of single-mode BP Neural Nets classifier based on acoustic features is shown in Table 3.

TABLE III        ACCURACY OF NEURAL NETS (%)

| Sample Sets | Positive | Non-positive |
|---|---|---|
| Positive | 76 | 24 |
| Non-positive | 22 | 78 |

Experiment3: recognition rate of single-mode classifier based on textual features is shown in Table 4.

TABLE IV        ACCURACY OF TEXTURL (%)

| Sample Sets | Positive | Non-positive |
|---|---|---|
| Positive | 90 | 10 |
| Non-positive | 12 | 88 |

Experiment4: recognition rate of decision level fusion algorism is shown in Table 5.

TABLE V        ACCURACY OF FUSION METHOD (%)

| Sample Sets | Positive | Non-positive |
|---|---|---|
| Positive | 94 | 6 |
| Non-positive | 8 | 92 |

From table 2 we can see the recognition rate of single-mode SVM classifier based on speech signal is around 83%. Non-positive emotion has shown its importance using value in practical application. Average recognition rate of non-positive emotion is around 81%, which means the acoustic features extracted in this paper have a higher correlation with non-positive emotion, can be used to recognize non-positive emotions. From table 3 we can see the recognition rate of single-mode BP Neural Nets classifier based on speech signal is around 77%. From table 4 we can see the recognition rate of single-mode classifier based on textual features is around 89%. From table 5 we can see the recognition rate of decision level fusion algorism proposed by this paper is around 93%, it's better than all the single-mode classifiers.

From experiment result we can see that the bimodal fusion algorism presented by this paper obtained the expected effect. The advantage of decision level fusion algorithm is that each classifier is independent from each other, when emotional data not available of with low quality in one channel, decision lever can still recognize emotion state, it has good robustness.

## VII. CONCLUSION AND PROSPECTS

This paper presents an approach to bimodal emotion recognition from speech signals and textual content. We conduct two parallel classifiers for acoustic information and two serial classifiers for textual information, and a final decision is m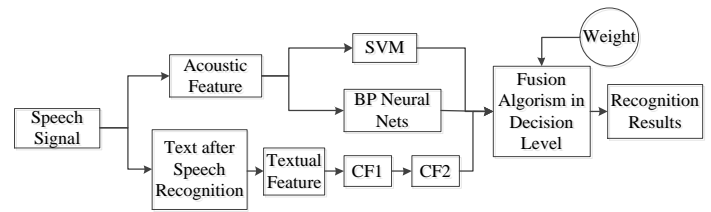ade by combing these classifiers in decision level fusion. Experimental results show that the emotion recognition accuracy of the integrated system is better than that of either of the two individual.

Emotion recognition cannot only combine speech and text, but also heart rate, blood pressure, skin current etc. physiological characteristics, which can be applied to polygraph, entertainment and many other areas. Affective computing will help artificial intelligence become more and more humanized.

REFERENCES

[1] Z. Zeng, M Pantic and Gl Roisman, "A survey of affect recognition methods: Audio, visual, and spontaneous expressions," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 33(1), pp. 39-58, 2009.

[2] S. Hoch, F Althoff, G McGlaun and G. Rigoll, "Bimodal fusion emotional data in an automotive environment," IEEE International Conference on Acounstics, Speech, and Signal Processing, USA, 2005, pp. 1085-1088.

[3] C. Busso, Z. Deng, S. Yildirim, M. Bulut, CM. Lee, A. Kazemzadeh et al. "Analysis of emotion recognition using facial expressions, speech and multimodal information," Prcoceeding of the Sixth International Conference on Multimodal Interfaces, USA, 2004, pp. 205-211.

[4] J. Wagner, J. Kim and E. Andre, "From physiological signals to emotions: implementing and comparing selected method for feature extracton and classification, " IEEE International Conference on Multimedia & Expo. Netherlands, 2005, pp. 940-943.

[5] XH. Fan, P. Wang and P. Zhou. "An two step text tendency analytical method based on extension," Computer engineering and Applications, China, vol. 48(1), pp. 162-165, 169, 2012.

[6] XH. Fan, and H. Wu, "Research of text orientation in the opinion leaders identification", Computer Application Study, China, vol. 30(9), pp. 2613-2615, 2636, 2013.

[7] XH. Fan, and MS. Sun, "A high performance two-class chinese text categorization method," Journal of Computers, China, vol. 29(1), pp. 124-131, 2006.

# Performance of window synchronisation in coherent optical ofdm system

Sofien Mhatli[1], Bechir Nsiri[2]
[1]SERCOM, EPT
Université de Carthage, 2078, La Marsa, Tunis , Tunisia
[2]Sys'com Lab, ENIT, BP.37 Le Belvédère 1002
Tunis, Tunisia

Mutasam Jarajreh[3], Basma Hammami[2], Rabah Attia[1]
[3]Faculty of Engineering & Environment, Northumbria
University, Ellison Place, Newcastle upon Tyne, NE1
8ST, UK

*Abstract*—**In this paper we investigate the performances of a robust and efficient technique for frame/symbol timing synchronization in coherent optical OFDM. It uses a preamble consisting of only two training symbol with two identical parts to achieve reliable synchronization schemes. The performances of the timing offset estimator at correct and incorrect timing in coherent optical OFDM are compared in term of mean and variance of the timing offset, and finally, we study the influence of number of subcarriers and chromatic dispersion.**

*Keywords*—*COFDM; timing offset; time synchronization; training symbol )*

## I. INTRODUCTION

Orthogonal frequency division multiplexing (OFDM) has been used in many telecommunication applications because of its high spectral efficiency and simple hardware implementation. OFDM has also been considered for optical systems as a candidate for future long range high data rate communication systems. The principle of frequency multiplexing is to group the digital data by N packets, which will be called OFDM symbol and modulating each given a different carrier simultaneously. One of the central features that set orthogonal frequency-division multiplexing (OFDM) apart from single-carrier modulation is its uniqueness of signal processing.

After transmission in fiber optic, the transmitted signal has distorted and rotated because the effect of white Gaussian noise (AWGN) and chromatic dispersion (CD) which are the reason why the system need channel estimation. For a high performances of the equalizer, we need have a precisely synchronization for the following module of channel estimator and channel equalizer.

T. Schmidl and D. Cox's algorithms make use of the correlation of training symbol to realize the synchronization. The synchronization algorithms depending on correlation need set a detect threshold of timing metric to detect whether the symbol synchronization succeed.

In this paper, we present the formulation of the Shmidl and Cox algorithm [4, 5] which is developed in radiofrequency, and we focus in depth on the performance of this algorithm in coherent optical OFDM. The rest of the article is organized as follows: in Section 2, the COFDM system is described. In section 3, the COFDM synchronization technique is described. The performance evaluation of synchronization system presented and discussed in Section 4. Finally, Section 5 concludes the work.

## II. COFDM SYSTEM DESCRIPTION

Figure 1 shows the conceptual diagram of a generic CO-OFDM system, including five basic functional blocks: RF OFDM transmitter, RF-to-optical (RTO) up-converter, optical link, optical-to-RF (OTR) down-converter, and RF OFDM receiver. In the RF OFDM transmitter, the input digital data are first converted from serial to parallel into a "block" of bits consisting of $Nsc$ information symbol, each of which may comprise multiple bits for m-ary coding. This information symbol is mapped into a two-dimensional complex signal Cki, for instance, using Gray coding, where Cki stands for the mapped complex information symbol. The subscripts of Cki correspond to the sequence of the subcarriers and OFDM blocks. The time domain OFDM signal is obtained through inverse discrete Fourier transform (IDFT) of Cki , and a guard interval is inserted to avoid channel dispersion [6,7]. The resultant baseband time domain signal can be described as :

$$s_B(t) = \sum_{i=-\infty}^{+\infty} \sum_{k=-\frac{N_{sc}}{2}+1}^{k=N_{sc}/2} c_{ki}\Pi(t - Ts)e^{j2\pi f_k(t-iT_s)} \qquad (1)$$

Where:

$$f_k = \frac{k-1}{t_s} \qquad (2)$$

$$\Pi(t) = \begin{cases} 1, (-\Delta G < t \le t_s) \\ 0, t \le -\Delta G, t > t_s \end{cases} \qquad (3)$$

where Cki is the i[th] information symbol at the k[th] subcarrier; $f_k$ is the frequency of the k[th] subcarrier; is the number of OFDM subcarriers Ts, ∆G, and ts are the OFDM symbol period, guard interval length, and observation period, respectively; and $\Pi(t)$ is the rectangular pulse waveform of the OFDM symbol. The extension of the waveform in the time frame of [-∆G, 0] in (1) represents the insertion of the cyclic prefix, or guard interval. The digital signal is then converted to an analog form through a DAC and filtered with a low-pass filter to remove the alias signal. The baseband OFDM signal can be further converted to an RF pass band through an RF IQ. The subsequent RTO up-converter transforms the base band signal to the optical domain using an optical IQ modulator comprising a pair of Mach–Zehnder modulators (MZMs) with

a 90 degree phase offset. The baseband OFDM signal is directly up-converted to the optical domain given by:

$$E(t) = e^{j(\omega_{LD1}t + \varphi_{LD1})} . s_B(t) \qquad (4)$$

Where $\omega_{LD1}$ and $\varphi_{LD1}$, respectively, are the angular frequency and phase of the transmitter laser. The up-converted signal E (t) traverses the optical medium with an impulse response of E (t), and the received optical signal becomes:

$$E'(t) = e^{j(\omega_{LD1}t + \varphi_{LD1})} . s_B(t) \otimes h(t) \qquad (5)$$



Fig. 1    Conceptual diagram for a generic CO-OFDM system with a direct up/down conversion architecture

Where $\otimes$ stands for convolution. The optical OFDM signal is then fed into the OTR downconverter, where the optical OFDM signal is converted to an RF OFDM signal. The directly down-converted signal can be expressed as:

$$r(t) = e^{j(\omega_{off}t + \Delta\varphi)} . r_0(t) \qquad (6)$$

$$r_0(t) = s_B(t) \otimes h(t) \qquad (7)$$

$$\omega_{off} = \omega_{LD1} - \omega_{LD2} \qquad (8)$$

$$\Delta\varphi = \varphi_{LD1} - \varphi_{LD2} \qquad (9)$$

$$c'_{ki} = e^{j\Phi_i} e^{j\Phi_D(f_k)} T_k c_{ki} + n_i \qquad (10)$$

$$\Phi_D(f_k) = \pi.c.D_t . \frac{f_j^2}{f_0^2} \qquad (11)$$

$$c'_{ki} = E'(t) \qquad (12)$$

Where $\omega_{off}$ and $\Delta\varphi$ are respectively the angular frequency offset and phase offset between the transmitted and receive lasers. $\Phi_D(f_k)$ Is the phase dispersion due to the fiber chromatic dispersion [8]. $\Phi_i$ Is the ofdm common phase error (CPE) [9] due to the phase noises from lasers and RF local oscillators at both the transmitter and receiver.

## III.    COFDM SYSTEM SYNCHRONISATION DESCRIPTION

In the RF OFDM receiver, the down-converted OFDM signal is first sampled with an ADC. Then the signal needs to go through the following three levels of sophisticated synchronizations before the symbol decision can be made:

- DFT window synchronization in which OFDM symbols are properly delineated to avoid intersymbol interference

- Frequency synchronization, namely frequency offset $\omega_{off}$ being estimated, compensated, and, preferably, adjusted to a small value at the start.

- The subcarrier recovery, where each subcarrier channel is estimated and compensated.



Fig. 2    Time domain structure of an OFDM signal



Fig. 3    the Schmidl synchronization format

Synchronization is one of the most critical functionalities for a CO-OFDM receiver. As discussed in the previous section, it can be divided into three levels of synchronization: DFT Window timing synchronization, carrier frequency offset synchronization, and subcarrier recovery. Figure 2 shows the time domain structure of an OFDM signal consisting of many OFDM symbols. Each OFDM symbol comprises a guard interval and an observation period. It is imperative that the start of the DFT window (i.e., the observation period) be determined properly because an improper DFT window will result in intersymbol interference (ISI) and intercarrier interference (ICI) [7].

One of the popular methods for window synchronization was proposed by Schmidl and Cox.[4]. In such a method, a pilot symbol or preamble is transmitted that consists of two identical segments, as shown in Figure 3, which can be expressed as:

$$s_m = s_{m - \frac{N_{sc}}{2}} , m \in \left[ \frac{N_{sc}}{2} + 1, N_{sc} \right] \qquad (13)$$

Where $s_m$ is the $m$ th sample with a random value when $m$ is from 1 to $N_{sc}/2$ . Assuming a time-invariant channel impulse response function h(t) the sampled received signal has the following form:

$$r_m = r\left(\frac{mt_s}{N_{sc}}\right) = r_m^0 e^{j\left(\omega_{off}\frac{mt_s}{N_{sc}}\right)} + n_m \qquad (14)$$

We have assumed that the constant phase across the entire OFDM symbol, or $\Delta\varphi$ equals zero in Eq. (9). The delineation can be identified by studying the following correlation function [4] defined as :

$$R(d) = \sum_{m=1}^{N_{sc}/2} r_{m+d}^* r_{m+d+N_{sc}/2} \qquad (15)$$

The principle is based on the fact that the second half of $r_m$ is identical to the first half except for a phase shift. Assuming the frequency offset $\omega_{off}$ is small to start with, we anticipate that when $d = 0$, the correlation function R (d) reaches its maximum value. The correlation function can be normalized to its maximum value given by:

$$M(d) = \left|\frac{R_d}{S_d}\right| \qquad (16)$$

And

$$S_d = \sqrt{\left(\sum_{m=1}^{N_{sc}/2} r_{m+d}^2\right)\left(\sum_{m=1}^{N_{sc}/2} r^2_{m+d+N_{sc}/2}\right)} \qquad (17)$$

Where $M(d)$ is defined as the DFT window synchronization timing metric. The optimal timing metric has its peak at the correct starting point of the OFDM symbol that is:

$$d_{opt} = \arg\{\max [M(d)]\} \qquad (18)$$

Where $\arg\{\max[M(d)]\}$ in general stands for searching the optimal argument of d that maximizes the objective function of $M(d)$, and $d_{opt}$ stands for the optimal timing point.

## IV. SIMULATION AND RESULT

We have conducted a matlab simulation for the coherent optical OFDM presented in figure 1 to confirm the DFT window synchronization using the Schmidl format for a CO-OFDM system at 10 Gb/s under the influence of chromatic dispersion, linewidth, and optical-to-signal noise ratio (OSNR).

TABLE I.      SIMULATION PARAMETERS

| Data rate (G/s) | 10 |
|---|---|
| Modulation | QAM |
| M | 4 |
| FFT Size | 128 |
| L(number of samples of the training symbol) | 256 |
| Number of samples per bit | 4 |
| Number of blocks | 1 |
| Number of symbols | 100 |
| Number of frames | 100 |
| CP length | 1/8 |
| Wavelength (nm) | 1.55 |
| Dispersion (ps/nm) | 17000 |
| Laser linewidth (KHz) | 100 |

The OFDM system parameters used for the simulation are described in the TABLE I mentioned above.

### A. Timing metric for CO-OFDM systems

This plot shows that the peak of the timing metric decreases from an ideal value of 1 to approximately 0.3 when the SNR is 0 dB. All curves show the flat platform corresponding to the guard interval inserted in the beginning of each symbol.

Fig. 4    Timing metric for CO-OFDM systems



Fig. 5    Expected value of timing metric. (Dashed lines indicate three standard deviations)

This figure shows a plot of the expected value of timing metric $M(d)$ versus SNR at both the best timing instant and a point outside the training symbol. The dashed lines indicate three standard deviations from each curve.

We concluded that outside the training symbol the timing metric is null.

### B. Mean and Variance at correct timing

The correct timing point was chosen at the start of the useful part of the first training symbol, the plots show that variance of timing metric is lower than the mean of timing metric.

Fig. 6    Mean and Variance at correct timing



Fig. 7    Mean and variance at incorrect timing

Figure 6 shows the result of simulation for the correct timing.

### C. Mean and Variance at incorrect timing

Figures 7 and 8 show the mean and variance of the incorrect timing. The incorrect timing point was chosen one symbol after the training symbol, the performance of the variance of timing metric is lower than the performance of the mean of the timing metric.

Some plot shows that there's a differentiation between the theory and simulation result and this lead on the two terms $r_{m+d}^{*}$ and $r_{m+d+N_{sc/2}}$ which are dependent and we consider in the simulations that are independent.

### D.  Influence of chromatic dispersion

This plot show that the increase of the chromatic dispersion leads on an increase of the timing offset and the minimum timing offset correspond on a null chromatic dispersion.

Fig. 8    Influence of chromatic dispersion



Fig. 9    Influence of number of subcarriers

### E.  Influence of number of subcarriers

This figure shows that the increase of the number of subcarrier lead on an increase of the timing offset; so to avoid this timing offset we must insert in the OFDM symbol in the beginning and in the end a sufficient number of zeros to avoid this increase of timing offset.

### V.    CONCLUSION

In this paper, we evaluated the performance of shmidl and cox algorithm of synchronization in coherent optical OFDM system. This method present a flat plateau in the timing offset versus SNR, so we aim in future works to eliminate this plateau and improve the performance of estimators in terms of mean and variance of timing metric and the complexity of synchronization algorithm.

#### REFERENCES

[1]   Shieh.W and Authadage C. coherent optical orthogonal frequency division multiplexing; electronic letters.42(10),587–589 (2006).

[2]   Shieh.W,Yang.Q, and MA.Y.107 Gb/s coherent optical ofdm transmission over 1000km SSMF fiber using orthogonal band multiplexing, optics express,16(9),6378–6386 (2008).

[3]   Ma CP, Kuo JW. Orthogonal frequency division multiplex with multi-level technology in optical storage application. Jpn J Appl Physics Part 1: Regular Papers Short Notes Rev Papers 2004;43:4876–8.

[4]   Schmidl TM, Cox DC. Robust frequency and timing synchronization for OFDM. IEEE Trans Commun 1997;45:1613–21.

[5]   Minn H, Bhargava VK, Letaief KB. A robust timing and frequency synchronization for OFDM systems. IEEE Trans Wireless Commun 2003;2:822–39.

[6]   Hara S, Prasad R. Multicarrier Techniques for 4G Mobile Communications. Boston: Artech House; 2003.

[7]   Hanzo L, Munster M, Choi BJ, Keller T. OFDM and MC-CDMA for Broadband Multi-User Communications, WLANs and Broadcasting. New York: Wiley; 2003.

[8]   Bauml RW, Fischer RFH, Huber JB. Reducing the peak-to-average power ratio of multicarrier modulation by selected mapping. IET Electron Lett 1996;32:2056–7.

[9]   Friese M. OFDM signals with low crest-factor. In: Proc. 1997 IEEE Global Telecommun. Conf; 1997. pp. 290–4.

# SentiTFIDF – Sentiment Classification using Relative Term Frequency Inverse Document Frequency

Kranti Ghag
Information Technology Department
MET's SAKEC, Mumbai University
Mumbai, India

Ketan Shah
Information Technology Department
SVKM's NMIMS MPSTME
Mumbai, India

*Abstract*—Sentiment Classification refers to the computational techniques for classifying whether the sentiments of text are positive or negative. Statistical Techniques based on Term Presence and Term Frequency, using Support Vector Machine are popularly used for Sentiment Classification. This paper presents an approach for classifying a term as positive or negative based on its proportional frequency count distribution and proportional presence count distribution across positively tagged documents in comparison with negatively tagged documents. Our approach is based on term weighting techniques that are used for information retrieval and sentiment classification. It differs significantly from these traditional methods due to our model of logarithmic differential term frequency and term presence distribution for sentiment classification. Terms with *nearly* equal distribution in positively tagged documents and negatively tagged documents were classified as a Senti-stop-word and discarded. The proportional distribution of a term to be classified as Senti-stop-word was determined experimentally. We evaluated the SentiTFIDF model by comparing it with state of art techniques for sentiment classification using the movie dataset.

*Keywords— Sentiment Classification; Term Weighting; Term Frequency; Term Presence; Document Vectors*

## I. INTRODUCTION

The web which is massively increasing resource of information has changed from read only to read write. Organizations now provide opportunity to the user to express their views on the products, decisions and news that are released [1]. Users can express their emotions as well can comment on the earlier user sentiments. Understanding consumer opinion for a product as well as for competitor's products is important for an organisation to take crucial decisions. Large amount of sentiment data is generated by various users for different features of products and services. Automatically processing this sentiment data needs to be handled systematically.

Sentiment Analysis involves extracting, understanding, classifying and presenting the emotions and opinions expressed by the users. Sentiment Classification generally involves classifying the polarity of a piece of text or classifying its subjectivity [2]. Polarity of a term, sentence, paragraph or document is classified as positive or negative [3]. At a deeper level of granularity it also involves identifying intensity such as moderately positive or strongly negative [4]. Some work has also focused on extracting the mood or emotion such as joy, surprise or disgust [5]. G. Paltoglou and

M. Thelwall systematically attempted to classify two emotions on scale of five leading to a combination of twenty-five classes [6]. The later part of subjectivity classification focused on identifying whether the term, sentence, paragraph or document is sentimental i.e. subjective or is informative i.e. objective. Lin, He and Everson have also worked on extracting only subjective part from a given document before classifying its polarity [7]. Most of the researcher used 6 to 9 emotions for sentiment classification where as the hour glass of emotions has many more emotions.

Sentiment analysis techniques can be broadly classified as supervised learning and unsupervised learning techniques [8].

Many unsupervised learning techniques use existing lexical resources (like WordNet) and language specific sentiment information (like sentiment seed words, their Synonyms and antonyms) to construct and update sentiment lexicons. [9], [10]. Very few sentiment lexicons are domain specific whereas most of these are generalized. Cross Domain lexicons were methodically extended to adapt for other related domains if the sentiment classes for one domain are available [11]. Unsupervised learning techniques assigned a generalized polarity and weight to a term failing to capture its domain specific context.

Supervised learning techniques constructed sentiment model trained with the help of tagged reviews. As these reviews are collection of domain-wise tagged set, the model constructed served well for specific domains [12]. It was also noted in our survey that most of the research in Sentiment Analysis has focused on supervised learning techniques such as Naïve-Bayes, Maximum-Entropy and Support Vector Machine (SVM) [13]. It was also marked that SVM was popularly used technique for Sentiment Classification. Supervised learning techniques entirely depend on the availability and the quality of tagged dataset.

A set of documents is used as training set to the classifier. These documents are represented as vectors. Every term in the document is an element in the vector in SVM approach for text mining. Term Presence and Term Frequency are two popular techniques for Information Retrieval when representing documents as vectors [8]. In Term Presence technique an element can take a binary value. This element is set to one if the term is present in document otherwise set to zero if the term is not present in document. In Term Frequency technique an element in the document vector is a non-negative integer that is set to count of the given term in a document.

For Sentiment Classification the training dataset consists of reviews tagged as positive and negative. All reviews tagged positive are called positively tagged documents whereas all reviews tagged negative are called negatively tagged documents. Every element in the vector represents a term that occurred in some document/s of training set. Each element of vector has two counts associated with it. One count is number of times of occurrence of that term (element) in positively tagged documents and other is number of times of occurrences in negatively tagged documents.

Our approach is based on traditional term weighting functions that are based on Term Frequency Inverse Document Frequency (TF-IDF) where the vectors are processed to identify and sequence index terms. Some of these are techniques are adapted for sentiment classification [3] [14]. These methods utilize combination of overall frequency count of term and proportional presence count distribution. Although our approach is based on traditional techniques of Information Retrieval, we examine whether addressing sentiment classification as special case of information retrieval can improve classification accuracy.

Accordingly we have attempted to adapt the model for sentiment classification, considering the similarities and differences with information retrieval techniques. In this paper a term was classified as positive if its TFIDF in positively tagged documents was more than negatively tagged documents and vice versa. This can be calculated using document vectors. The $i^{th}$ element of each vector that was constructed from positively tagged documents contributed to positivity of $i^{th}$ term and similarly $i^{th}$ element of each vector that was constructed from negatively tagged documents contributed to negativity of the same term.

Our approach differs significantly from traditional approaches on the basis of usage pattern of term presence and term count vectors. We focus on *proportional* frequency count distribution and proportional presence count distribution whereas traditional approaches such as delta TFIDF and other term weighting techniques rely on combination of overall frequency count of term and proportional presence count distribution.

The rest of the paper is organized as follows. Sentiment Classification techniques are surveyed in section 2. Section 3 focuses on the proposed model for Sentiment Classification. Experimental setup is discussed in section 4. Results are presented in section 5. Concluding remarks and future scope are put forth in section 6.

## II. PRIOR WORK

Pang, Lee and Vaithyanathan laid the foundation of harnessing supervised machine learning techniques for Sentiment Classification. They are also the pioneers for extracting, transforming and making available the popular movie review dataset. Naive Bayes, maximum entropy classification, and support vector machines algorithms were applied on unigrams and bigrams features and their weights, extracted from this movie dataset [15]. They concluded that sentiment analysis problem needs to be handled in a more sophisticated way as compared to traditional text categorization techniques. SVM classifier applied on unigrams produced best results unlike information retrieval where bigrams generate remarkable accuracy as compared to unigrams.

Mullen and Collier used SVMs and expanded the feature set for representing documents with favorability measures from a variety of diverse sources [16]. They introduced features based on Osgood's Theory of Semantic Differentiation, using Word-Net to derive the values of potency, activity and evaluative of adjectives [17] and Turney's semantic orientation [18]. Their results showed that using a hybrid SVM classifier that uses as features the distance of documents from the separating hyper plane, with all the above features produces the best results.

Zaidan, Eisner, and Piatko introduced "annotator rationales", i.e. words or phrases that explain the polarity of the document according to human annotators [19]. By deleting rationale text spans from the original documents they created several contrast documents and constrained the SVM classifier to classify them less confidently than the originals. Using the largest training set size, their approach significantly increased the accuracy on movie review data set.

Prabowo and Thelwall [20] proposed a hybrid classification process by combining in sequence several ruled-based classifiers with a SVM classifier. The former were based on the General Inquirer lexicon by lin, Wilson, Wiebe and Hauptmann. [21] and the MontyLingua part-of-speech tagger by Liu [22] and co-occurrence statistics of words with a set of predefined reference words. Their experiments showed that combining multiple classifiers can result in better effectiveness than any individual classifier, especially when sufficient training data isn't available.

Bruce and Wiebe made an effort to manually tag sentences as subjective or objective by different judges and the resultant confusion matrix was analyzed [23]. 14 articles were randomly chosen and every non-compound sentence was tagged. Also a tag was attached to conjunct of every compound sentence. Authors then attempted to identify if pattern exists in agreement or disagreement between human judges. Authors observed that manual tagging suffered due drawback of biased nature of human beings during tagging phase.

Dave, Lawrence and Pennock used a self tagged corpus of sentiments [24] available on major websites such as Amazon and Cnet as training set. Naïve Bayes classifier was trained and refined using the above corpus. The classifier was then tested on other portion of self-tagged corpus. The sentences were parsed to check semantic correctness and then tokenized. Techniques such as co- allocation substrings and stemming were applied for generalisation of tokens. When pre-processed, N-grams (bi-gram and tri-gram) improved the results as compared to unigram. They also applied smoothing so that non-zero frequencies were available. Score were then assigned to features.

Zhang constructed computational model that explored reviews linguistics properties to judge its usefulness [25]. Support Vector Regression (SVR) algorithm was used for

classification. In contrast to major studies which filter out subjective information in any review or are not considered important, Zhang claimed that the quality of review was reasonably good if it was a good combination of subjective and objective information.

Wang and Dong encoded semantic information and grammatical knowledge into a lower dimension vector to represent text for the purposes of sentiment classification [26]. Grammatical knowledge-embedding representation methods were used to provide extra information for the classification algorithm. This reduced the space complexity. Longer text contained more information about the semantic orientation features. Sometimes longer text contained subjective information that had contradictory sentiments.

Ghosh and Iperrotis found that reviews with combination of objective and subjective sentences had more impact on sales of a product as compared to reviews with purely subjective sentences [27]. Random forest based classifier was used to classify reviews. The impact of subjectivity, information, readability and linguistic correctness in reviews affected in influencing sales and perceived usefulness. Li and Liu obtained stable clustering for opinion clustering by applying a TF-IDF weighting method, voting mechanism and importing term scores [28].

Yu, Liu and Huang attempted to identify hidden sentiment factors in the reviews [29]. Bag of words approach was used for sentiment identification in the review. Along with sentiment identification, product sales prediction methods were also proposed.

Lin, Everson and Ruger preprocessed reviews to extract words and noise such as punctuation, numbers, and non-alphabet characters were removed [30]. Stemming was applied so that the related terms fall in same clusters, thus reducing the vocabulary classes. MPQA and appraisal lexicons were merged stemmed and cleaned to form a new lexicon which was used to classify the document irrespective of the domain.

Kumar and Ahmad proposed a preliminary prototype ComEx Miner System for mining experts in virtual communities [31]. They constructed a collaborative interest group known as the virtual community which grouped researchers with similar interests to facilitate collaborative work. The expertise from the virtual community was retrieved using sentiment analysis of each group member's blog & comments received on it. Authors with top ranks were identified based on the ranks that there blogs received.

Hamouda and El-taher developed a corpus using different machine learning algorithms such as decision tree, support vector machines and naive bayes for Arabic Facebook news pages [32]. They constructed corpora for supportive comments, attacking comments, and neutral comment for different posts. They claimed that best result was obtained by the support vector machine classifier with 73.4% of accuracy on their test set.

93.75% of the synonymous sets in SentiWordNet are ignored as they have a stronger objective tendency [33]. Hung and Lin re-evaluated these objective words in SentiWordNet based on their presence in positive and negative sentences.

Apart from traditional pre-processing and document classification, an additional step that may re-assign a sentiment polarity to objective word was incorporated. They also marked that concluding sections of a review were strongly sentimentally oriented, which could be used for dimensionality reduction.

TFIDF is a popular statistical technique to index the term as per their importance. TFIDF is based on documents and term vectors that represent term frequency as well as term presence [34] [35]. Term presence could be constructed if term frequency vector is available but vice-versa is not possible.

$$d^{(i)} = TF(w_i, d).IDF(w_i) \tag{1}$$

Where,
$w_i$ = $i^{th}$ term.
d = document.
$d^{(i)}$ = TFIDF of term $w_i$ in document d.
$TF(w_i,d)$ = Term Frequency of term $w_i$ in document d.
and $IDF(w_i)$ = Inverse Document Frequency.

TFIDF of term $w_i$ in document d can be computed using "(1)". Term frequency $TF(w_i,d)$ is count of a term $w_i$ in document d. Larger value of a Term Frequency indicates its prominence in a given document. Terms present in too many documents were suppressed as these tend to be stop words. This suppression was handled by the second component IDF.

$$IDF(w_i) = \log\left(\frac{|D|}{DF(w_i)}\right) \tag{2}$$

Where,
$IDF(w_i)$ = Inverse Document Frequency.
$w_i$ = $i^{th}$ term.
$|D|$ = the total count of documents.
$DF(w_i)$ = count of documents that contain term $w_i$.

If a term is present in all the documents then numerator equals denominator in "(2)". As a result of this $IDF(w_i)$= log 1 which is zero. But if term occurred in relatively less number of document then $DF(w_i) < |D|$. As a result $IDF(w_i)$ = log (>1) which is a positive integer. Term presence vector was used for calculation of IDF. TFIDF identified important terms in given set of documents but as per Martineau and Finin top ranked index terms were not the top ranked sentimentally polarized terms [3].

In connection with the occurrences of rare words, different variations of TFIDF scores of words, indicating the difference in occurrences of words in different classes (positive or negative reviews), have been suggested by Paltoglou and Thelwall [14]. They surveyed many term weighting techniques as well proposed "smart" and "BM25" term weighting techniques for sentiment classification.

TFIDF identified important terms in given set of documents but as per Martineau and Finin top ranked index terms were not the top ranked sentimentally polarized terms [3]. Martineau and Finin constructed vectors to classify a term

based on term frequency vector as well as term presence vectors. Unlike TFIDF which used single term presence vector, two vectors were separately constructed for presence in positively tagged documents and negatively tagged documents [3].

$$V_{td} = C_{td} \times log\left(\frac{|N_t|}{|P_t|}\right) \qquad (3)$$

Where,

$V_{td}$ = Polarity of term t in document d.

$C_{td}$ = count of a term in a given document.

$|N_t|$ = count of negatively tagged documents with term t.

$|P_t|$ = count of positively tagged documents with term t.

More importance was given in "(3)" to terms that occurred frequently. The later part of the model i.e. log $(|N_t|/|P_t|)$ contributed to polarity of a term. Term presence count of a term was number of documents that term was present. log $(|N_t|/|P_t|)$ component returned a negative value if a term occurred in more number of positively tagged documents as compared to negatively tagged documents and vice-versa.

If a term was present in equal number of positive and negative document then this component returned zero. Since this value was multiplied with $C_{td}$, resulting $V_{td}$ value was also grounded. These terms were classified as stop words.

Delta TFIDF returned a negative value if the term was classified as positive and vice-versa.

It considered overall count of terms in all documents ignoring the frequency distribution of terms across positively and negatively tagged documents. For example if a term was present in more number of negatively tagged documents as compared to positively tagged document, term was classified as negative. Although the term was present in less number of positively tagged documents, its frequency count in these positively tagged documents may be more which contributed to $C_{td}$ part. This incorrectly boosted the $V_{td}$ value.

$C_{td}$ being frequency count of terms over all the documents did not correctly relates to second part of the model that dealt with distribution of presence.

To calculate polarity of i[th] term summation of i[th] element of the vectors was taken in which log $(|N_t|/|P_t|)$ was common. Sum of $C_{td}$ which was always a positive number acted as a boosting factor.

### III. INTRODUCING SENTITFIDF

Our model SentiTFIDF works on the principle logarithmic proportion of TFIDF of a term across positively tagged documents and negatively tagged documents. If the TFIDF of a term in positively tagged documents is larger than TFIDF of same term in negatively tagged documents the term is assigned positive polarity and vice-versa.



Fig. 1. The proposed SentiTFIDF based on relative TFIDF

Figure 1 presents the system flow of SentiTFIDF. It can be divided into three parts. In first part the positivity of a term is calculated. Similarly negativity of a term is calculated in second part. Third part classifies the term as positive, negative or neutral based on its proportion of positivity and negativity calculated in previous steps.

In the first part term presence vector and term frequency vector are constructed in "(6)" and "(7)" for positively tagged documents using "(4)" and "(5)",.

$$TFP = \begin{bmatrix} c_{11} & \cdots & c_{1d} \\ \vdots & \ddots & \vdots \\ c_{t1} & \cdots & c_{td} \end{bmatrix} \qquad (4)$$

Where,

TFP=Term Frequency Matrix for positively tagged documents

t = term        d = document

TF[i][j] = $c_{ij}$ = count of term i in document j

$$TPP = \begin{bmatrix} p_{11} & \cdots & p_{1d} \\ \vdots & \ddots & \vdots \\ p_{t1} & \cdots & p_{td} \end{bmatrix} \qquad (5)$$

Where,

TPP = Term Presence Matrix for positively tagged documents.

t = term.        d = document.

TF[i][j] = $p_{ij}$ = presence of term i in document j

        = 1 if term i is present in document j otherwise 0.

Using the term frequency and presence matrix, frequency of term in positively tagged documents and number of positive documents that contain term t can be computed as follows.

$$P_{ctd} = \sum_{j=1}^{d} c_{tj} \qquad (6)$$

Where,
$P_{ctd}$ = Frequency of term t in positively tagged documents.
$c_{tj}$ = count of term t in $j^{th}$ document.

$$P_t = \sum_{j=1}^{d} p_{tj}$$

(7)

Where,
$P_t$ = Number for positively tagged documents with term t.
$p_{tj}$ = presence of term t document j

In "(8)" TFIDF of the terms is calculated by using the vectors in "(6)" and "(7)" of positively tagged documents. This value contributes of the positivity of the terms.

$$Pos_t = P_{ctd} \times log \frac{P}{P_t}$$

(8)

Where,
$Pos_t$ = Positivity of term t.
$P_{ctd}$ = Frequency of term t in positively tagged documents.
P = Total Number of positively tagged documents.
$P_t$ = Number for positively tagged documents with term t.

Similarly in second part the negatively tagged documents are represented as document vectors i.e. in form of term presence vector and term frequency vector. Negativity of the terms is calculated using TFIDF on negatively tagged documents, vectors, as in "(9)".

$$Neg_t = N_{ctd} \times log \frac{N}{N_t}$$

(9)

Where,
$Neg_t$ = Negativity of term t.
$N_{ctd}$ = Frequency of term t in negatively tagged documents.
N = Total Number of negatively tagged documents.
$N_t$ = Number for negatively tagged documents with term t.

If positivity of a term is larger than negativity of the same term than the term is classified as positive. Conversely, if negativity of a term is larger than positivity of the same term then the term is classified as negative in the third part using "(10)" and "(11)". A term is classified as neutral if its positivity equals negativity.

$$LDT_t = log \frac{Pos_t + 0.001}{Neg_t + 0.001}$$

(10)

Where,
$LDT_t$ = Logarithmic differential TFIDF.
$Pos_t$ = Positivity of term t.
$Neg_t$ = Negativity of term t.

$$1 \qquad\qquad > 0$$

$$Pol_t = \quad 0 \qquad if \qquad LDT_t \quad = 0 \qquad (11)$$
$$-1 \qquad\qquad\qquad\qquad < 0$$

Where,
$Pol_t$ = Polarity of term t.
$LDT_t$ = Logarithmic differential TFIDF.

If negativity of a term was zero the model would have been affected by divide by zero error. So we added a small value 0.001 to $Neg_t$ i.e. denominator part. As a result, the term was classified as positive if positivity of term was nonzero and neutral if positivity of the term was also zero.

Similarly if $Pos_t$ =0 and $Neg_t$ is not equal to zero then the term should be classified as negative. This is accomplished by adding a negligible value 0.001 to the numerator.

Stop-words are handled while computing $Pos_t$ and $Neg_t$, in "(8)" and "(9)". If a term occurred in all positive and negative documents then $Pos_t$ = 0 and also $Neg_t$ = 0. Thus $LDT_t$ and $Pol_t$ would be zero using "(10)" and "(11)". These are classified as neutral terms.

Our model handles specialized class of stop-words called as senti-stop-words. The terms whose positivity equals negativity are classified as neutral using "(10)" as log 1 equals zero.

There might be Senti stop-words whose positivity might are not exactly evenly distributed so we varied parameters in experiment 1 to further improve our model. Senti-stop word differs from stop word depending on its distribution across documents of both classes.

## IV. EXPERIMENTS CONDUCTED

Pang and Lee's movie dataset with 1000 positively tagged text documents and 1000 negatively tagged text document were used in all the experiment. These review text files size varied from 1KB to 15 KB. Number of words per document varied from 17 to 2678.

A list of terms that occurred in the documents was prepared. A term is entered only once in this term list although it may appear may times in documents. A document vector was constructed for every document. Every $i^{th}$ element in this vector was count of $i^{th}$ term in this document. If a term in term list was not present in the document the count associated with that term was set to zero. These vectors were used to calculate term polarity for the terms in the term list. Polarity was calculated using proposed SentiTFIDF model as well as Delta TFIDF model described in section 3 and 2 respectively. A term was classified either as positive or negative or neutral

A document was classified by our model as positive if total number of positive terms in the document were more than negative terms. Similarly a document was classified as negative if total number of negative terms in the document were more than positive terms.

If a document was originally tagged as positive and also classified as positive then it contributed to True Positive in confusion matrix. If a document was originally tagged as negative and also classified as negative then it contributed to

True Negative in confusion matrix. If a document was originally tagged as positive but classified as negative then it contributed to False Negative in confusion matrix. If a document was originally tagged as negative but classified as positive then it contributed to False Negative in confusion matrix.

### A. Experiment 1

Experiment 1 was conducted to determine that a term t should be classified as Senti-stop-word if $LDT_t$ exactly equals zero or it was within a specified range. For this accuracy was computed using, 10 Fold Cross Validation (10 fold CV), varying the range of $LDT_t$ between 0 to 5 at step of 0.5 and simultaneously 0 to -5 at step of -0.5.

To calculate accuracy dataset was divided in 10 parts. At every fold this 10% dataset was used for testing and remaining 90% dataset was used for training the classifier.

Confusion matrix was constructed as well as accuracy was calculated at every fold and then averaged to form the accuracy of the model.

### B. Experiment 2

10 Fold Cross Validation (10 fold CV) technique [4] was used to calculate accuracy of RTFSC and Delta TFIDF. Dataset was divided in 10 parts. At every fold this 10% dataset was used for testing and remaining 90% dataset was used for training the classifier. $LDT_t$ range was now set to -0.5 to 0.5 as determined in experiment 1 for a term to be classified as neutral. Confusion matrix was constructed as well as accuracy was calculated at every fold and then averaged to form the accuracy of the model at that value of $LDT_t$.

### C. Experiment 3

Accuracy was calculated using 10% of the dataset as training set and remaining 90% as the test set. Training dataset was incremented by 10% and remaining was used as test data in further iterations till 90% data was used for training and 10% for testing. Accuracy was calculated at every repetition.

### D. Experiment 4

Accuracy was calculated using the entire dataset for training set as well as for testing.

## V. RESULTS AND DISCUSSION



Fig. 2. 10 Fold Cross Validation and varying LDTt

Figure 2 represents accuracy at variations of the parameter $LDT_t$ defined in "(10)". Every series represents a fold and the order of fold is not important. In every fold 90% data of dataset is used for training and remaining 10% for testing. It can be marked that accuracy is maximum for all the folds when $LDT_t = 0.5$, which is even larger than when $LDT_t = 0$. This indicated that theoretical model needed to be adapted efficiently classify Senti-stop-words. So the term was classified as positive if $LDT_t > 0.5$ and negative if $LDT_t < -0.5$. A term was classified as Senti-stop-word if $-0.5 < LDT_t < -0.5$. All further experiments were performed using "(12)" instead of "(11)".

$$Pol_t = \begin{cases} 1 & > 0.5 \\ 0 & \text{if} \quad LDT_t = -0.5 \text{ to } 0.5 \\ -1 & < 0 \end{cases} \quad (12)$$

Where   $Pol_t$ = Polarity of term t.
         $LDT_t$ = Logarithmic differential TFIDF.

Fig. 3.   10 Fold CV of SentiTFIDF and Delta TFIDF

Figure 3 represents accuracy at each of 10 folds. The order of fold is not important. Accuracy of SentiTFIDF was more than Delta TFIDF in all folds. The accuracy at every fold was averaged for comparison. The average accuracy of all folds of SentiTFIDF was 73.8% and that of Delta TFIDF was 66.5%. This indicates that our method performs better than Delta TFIDF and the improvements are independent of the data used for training set.



Fig. 4.   Accuracy by incrementing training set

Figure 4 represents accuracy of delta TFIDF and SentiTFIDF when the training dataset is incremented by 10% for all iteration and remaining data is used for testing. For any good classifier the accuracy should increase when training data is increased. Accuracy of SentiTFIDF as well as delta TFIDF increases as training data is incremented. At every iteration accuracy of Senti-TFIDF was more delta TFIDF.

As the size of training dataset was incremented the accuracy of our algorithm increased. The accuracy was also always more than Delta TFIDF. This indicates that even if any percentage of data is used for training our SentiTFIDF outperformed Delta TFIDF. Irrespective of the size of dataset SentiTFIDF performed well.



Fig. 5.   Accuracy with complete dataset for training and testing

Figure 5 represents accuracy of Delta TFIDF and SentiTFIDF when entire dataset was used for training as well as for testing. Accuracy of SentiTFIDF was 92% and Delta TFIDF was 85%. Even if maximum dataset is used for training SentiTFIDF performs better than DeltaTFIDF

## VI.   CONCLUSION AND FUTURE WORK

From the results of the experiments conducted it can be observed that accuracy of SentiTFIDF is more than Delta TFIDF. Unlike Delta TFIDF, SentiTFIDF efficiently handles absence of a term in positively and / or negatively tagged documents, thus eliminating divide by zero error as well as term getting wrongly classified. Frequency and presence of a term in all the documents is an important aspect of Information Retrieval. SentiTFIDF considers frequency and presence *distribution* of a term across positively and negatively tagged documents as compared to delta TFIDF which considers frequency of a term in all documents and distribution of presence even for Sentiment Classification.

Accuracy of SentiTFIDF was 92%. The accuracies of surveyed techniques that were tested using movie review dataset were between 84.6% and 92.2%. Although these accuracies cannot be directly compared as the experimental parameters may vary, SentiTFIDF performs better than most existing techniques except for Zaidan, Eisner, and Piatko model of SVM with annotator rationales which achieve accuracy of 92.2%.

Our classifier is based on term frequency and presence distribution. In future we aim to experiment the effect of other distributional count associated with terms.

REFERENCES

[1] H. Chen, "Business and Market Intelligence 2.0, "*In IEEE Intelligent Systems,*vol. 25, issue no. 01, pp. 68–71, 2010.

[2] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and Trends in Information Retrieval*, vol. 2, nos. 1–2, pp. 1–135, 2008.

[3] J. Martineau and T. Finin, "Delta TFIDF: An Improved Feature Space for Sentiment Analysis," *In Proc. of 3rd Int'l AAAI Conf. on Weblogs and Social Media,* pp.258-261, 2009.

[4] L. Lee and B. Pang, "Seeing Stars: Exploiting Class Relationships for Sentiment Categorization with Respect to Rating Scales," *In Proc. of Ann. Meeting Assoc. Computational Linguistics*, pp. 115-124, 2005.

[5] E. Cambria, R. Speer, C. Havasi, and Amir Hussain, "SenticNet: A Publicly Available Semantic Resource for Opinion Mining," *In Proceedings of AAAI Conf. on Commonsense Knowledge*, pp. 14-18, 2010.

[6] G. Paltoglou and M. Thelwall, "Seeing Stars of Valence and Arousal in Blog Posts," *In IEEE Trans. on Affective Computing*, vol. 04, issue no. 01, pp. 116-123, 2013.

[7] C. Lin, Y. He, and R. Everson, "Sentence subjectivity detection with weakly-supervised learning," *In Proc. of 5th Int'l Joint Conf. on Natural Language Processing*, pp. 1153–1161, 2011.

[8] E. Cambria, B. Schuller, Y. Xia, and C. Havasi, *"*New Avenues in Opinion Mining and Sentiment Analysis,*" Intelligent Systems, IEEE*, vol.28, no.2, pp. 15-21, 2013.

[9] S. Baccianella, A. Esuli, and F. Sebastiani, " SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining," *In Proc. of 7th Int'l Conf. on Language Resources and Evaluation,* pp 2200-2204, 2010.

[10] A. Neviarouskaya, H. Prendinger, and M. Ishizuka, "SentiFul: A Lexicon for Sentiment Analysis," *In IEEE Trans. On Affective Computing*, vol. 02, issue no. 01, pp. 22-36, 2011.

[11] Bollegala, D. Weir, and J. Carroll, "Cross-Domain Sentiment Classification using a Sentiment Sensitive Thesaurus," *Knowledge and Data Engineering, IEEE Transactions, vol.25, issue no.08,* pp. 1, 0

[12] R. Xia and C. Zong, "A POS-based Ensemble Model for Cross-domain Sentiment Classification," *In Proc. of 5th Int'l Joint Conf. on Natural Language Processing*, pp. 614–622, 2011.

[13] K. Ghag and K. Shah, "Comparative analysis of the techniques for Sentiment Analysis," *In Proc. of Int'l Conf. on Advances in Technology and Engineering,* pp. 1-7, 2013.

[14] G. Paltoglou and M. Thelwall, "A study of Information Retrieval weighting schemes for sentiment analysis," *In Proc. of 48th Annual Meeting of the Association for Computational Linguistics,* pp. 1386-1395, 2010.

[15] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? sentiment classification using machine learning techniques," *In Proc. of Conf. on Empirical Methods in Natural Language Processing,* pp 79-86, 2002.

[16] T. Mullen and N. Collier, "Sentiment analysis using support vector machines with diverse information sources,' *In Proc. of Conf. on Empirical Methods in Natural Language Processing,*" pp 412–418, 2004.

[17] Charles E. Osgood, George J. Suci, and Percy H. Tannenbaum. The measurement of meaning, 2$^{nd}$ ed.. University of Illinois Press Urbana, 1967.

[18] P. D. Turney, "Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews," *In Proc. of the 40th Annual Meeting on Association for Computational Linguistics ACL,* pp 417–424, 2002.

[19] O.F. Zaidan, J. Eisner, and C.D. Piatko, "Using Annotator Rationales to Improve Machine Learning for Text Categorization," *In Proc. of Conf. of North American Chapter of the Association for Computational Linguistics,* pp 260–267, 2007.

[20] Rudy Prabowo and Mike Thelwall,"Sentiment analysis: A combined approach," *Journal of Informetrics,*3(2):143–157, 2009.

[21] Wei-Hao Lin, Theresa Wilson, Janyce Wiebe, and Alexander Hauptmann, "Which side are you on? identifying perspectives at the document and sentence levels," In *Proceedings of the Conferenceon Natural Language Learning* ,2006.

[22] Hugo Liu., "MontyLingua: An end-to-end natural language processor with common sense". *Technical report,MIT*, 2004.

[23] R. F. Bruce and J. M. Wiebe, "Recognizing Subjectivity: A Case Study in Manual Tagging," *In Natural Language Engineering, ACM, vol. 05, issue 02*, pp 187-205, 1999.

[24] K. Dave, S. Lawrence, and D. M. Pennock, "Mining the Peanut Gallery: Opinion Extraction and Semantic Classification of Product Reviews," *In Proc of the 12th Int'l Conf. on World Wide Web*, pp 519 - 528, 2003.

[25] Z. Zhang, "Weighing Stars: Aggregating Online Product Reviews for Intelligent E-commerce Applications," *In Intelligent Systems, IEEE, vol. 23, issue no.05*, pp.42-49, 2008.

[26] J. Wang and A. Dong, "A Comparison of Two Text Representations for Sentiment Analysis," *In Proc. of Int'l Conf. on Computer Application and System Modeling,* pp V11-35–V11-39, 2010.

[27] Ghose and P. G. Ipeirotis, "Estimating the Helpfulness and Economic Impact of Product Reviews: Mining Text and Reviewer Characteristics," *In Knowledge and Data Engineering, IEEE Transactions, vol.23, issue no.10*, pp. 1498-1512, 2011.

[28] G. Li and F. Liu, "A Clustering-Based Approach on Sentiment Analysis," *In Proc. of Int'l Conf. on Intelligent Systems and Knowledge Engineering*, pp 331-337, 2011.

[29] X. Yu, Y. Liu, X. Huang, and A. An, "Mining Online Reviews for Predicting Sales Performance: A Case Study in the Movie Domain," *In Knowledge and Data Engineering, IEEE Transactions, vol.24, issue no.04,* pp. 720-734, 2012.

[30] Lin, Y. He, R. Everson, and S. Ruger, "Weakly Supervised Joint Sentiment-Topic Detection from Text," *In Knowledge and Data Engineering, IEEE Transactions*, vol.24,issue no.06, pp. 1134-1145, 2012.

[31] A. Kumar and N. Ahmad, "ComEx Miner: Expert Mining in Virtual Communities", *In International Journal of Advanced Computer Science and Applications, vol.03, issue no 06,* pp 54-65, 2012.

[32] A. Hamouda and F. El-taher, "Sentiment Analyzer for Arabic Comments System," *International Journal of Advanced Computer Science and Applications, vol. 04, issue no.03*, pp 99-103, 2013

[33] C. Hung and H. Lin, "Using Objective Words in SentiWordNet to Improve Word-of-Mouth Sentiment Classification," *In IEEE Intelligent Systems,* vol.28, no.2, pp. 47-54, 2013.

[34] G. Salton and M. McGill, "*Introduction to modern Information Retrieval,*" McGraw-Hill, pp. 105-107 & 205,1983.

[35] J. Han and M. Kamber, *"Data Mining Concepts and Techniques,"* Morgan Kaufmann Publishers, 02$^{nd}$ edition, pp. 364-365,2006.

# New Technique to Insure Data Integrity for Archival Files Storage (DIFCS)

Mohannad Najjar

Computer and Information Technology
University of Tabuk
Tabuk, Saudi Arabia

*Abstract*—**In this paper we are developing an algorithm to increase the security of using HMAC function (Key-Hashed Message Authentication) to insure data integrity for exchanging archival files. Hash function is a very strong tool used in information security. The algorithm we are developing is safe, quick and will allow the University of Tabuk (UT) authorities to be sure that data of archival document will not be changed or modified by unauthorized personnel through transferring in the network; it will also increase the efficiency of network in which archived files are exchanged. The basic issues of hash functions and data integrity will be presented as well.**

**In this research: The developed algorithm is effective and easy to implement using HMAC algorithm to guarantee data integrity for archival scanned documents in the document management system.**

*Keywords—cryptography; hash functions; data integrity; authentication; HMAC; file archiving*

## I. INTRODUCTION

Information and data security in the different systems at UT are one of the most critical issues for the university authorities. Ensuring data in these systems are not modified in an unauthorized fashion is a fundamental goal. UT departments use different kinds of information systems: the academic system, the ERP system, document management system etc. All of these systems don't have any tool to guarantee the integrity of their data.

Data Integrity is one of the fundamental components of information security. Data integrity is a tool used to insure that data (documents, messages, emails, files, etc.) can't be changed, modified, deleted by unauthorized personnel, thereby insuring accuracy and consistency.

When a message is sent through the local network or Internet to a Receiver; data integrity tools are used to insure that the message was not altered and that it is identical to that sent from the Sender. There are many tools to insure data integrity, such as: parity bit, checksum, encryption and hash functions. Hash functions are one of the most used tools because of simplicity, speed and being free of charge.

Insuring Data Integrity is already an important tool used in data exchange in telecommunications and networking systems. For UT the use of DIFCS (Data Integrity File Checking System) algorithm will guarantee that data stored in all applications will be safe and reliable. This solution also will

increase the safety of the university information systems, in a convenient and effective method. Additionally the DIFCS algorithm will increase the effectiveness of the whole files archive system.

We depend in our improved DIFCS algorithm on using HMAC function to insure data integrity, authentication and we will add additional improved techniques to increase the effectiveness of the algorithm in the local network.

**List of important symbols used in the paper:**

| | |
|---|---|
| H | The hash function, MD5 or SHA-1 |
| B | The number of bits in the block in the hash function |
| IV | The initial value for the hash function |
| M | The data input to HMAC |
| $Yi$ | The $i^{th}$ block of m, $0 \leq i \leq (l-1)$ |
| L | The number of blocks in m after padding |
| N | The length of hash code |
| K | The secret key, if K length is greater than *b* then K=$h$(K) |
| K+ | The K padded with zeros on the left so the result has *b* bits |
| ipad | The inner pad; the byte 36 (in hexadecimal) repeated b/8 times |
| opad | The outer pad; the byte 5c (in hexadecimal) repeated b/8 times |
| $h(m)$ | The value of the HMAC; the length of the data is n bits, where the maximum value for n depends on the hash function used, MD5 or SHA-1 |
| Y | Set of all possible hash results |

## II. HASH FUNCTIONS

Hash function is a function h: M$\rightarrow$Y that has, as a minimum, two properties:

- it compresses a sequence $m \in M$ of bits of arbitrary length, including the empty sequence, into a sequence $h(m) \in Y$ of the constant (fixed) length,

- for any $m \in M$ it is easy to compute $h(m)$.

The hash function transformation of the message $m = m_1 \| m_2 \| \ldots \| m_t$ divided into fixed length blocks $m_1, m_2, \ldots, m_t$ can be described as follows (see Fig. 1):

$$H_0 = IV,$$

$$H_i = \varphi(m_i, H_{i-1}) \text{ for } i = 1, 2, \ldots, t;$$

Where; IV is an initial value, $H_i$ is a chaining variable, $\varphi$ is a compression function (also called a round function) and ψ is an output transformation. As a result we obtain h(m) of fixed length. In cryptographic literature [2,5] the resulting sequence h(m) has been given a wide variety of names: hash result, hash code, hash total, imprint, fingerprint, message digest, cryptographic checksum, authenticator, authentication tag, compression, compressed encoding, condensation, Message Integrity Code (MIC), etc. In the sequel h(m) will be called hash result.

The structural model of the hash function is presented in Figure 1. [2]. It works well if the length of $m_t$ is of the same length as each previous block $m_1,m_2,\ldots,m_{t-1}$. If it is not a case then extra bits must be appended to an input string before hashing to make $m_t$ as long as $m_1,m_2,\ldots,m_{t-1}$.



Fig. 1.  General model of the hash function $h$

### III.  DATA INTEGRITY

Any information system is deemed secure if it has at least three properties: Confidentiality, Data Integrity and Availability. So data integrity is one of the most important aspects of security according to data.

It insures the accuracy and consistency of data stored or transmitted from one point to another. There are many methods for insuring data integrity: physical and logical.

Physical tools like RAID (Redundant Array of Independent Disks). And logical like parity bit, CRC, Checksum, Encryption and Hash functions. In our paper we will improve a logical tool that will use hash function to insure data integrity of archived documents and files.

We will focus on insuring data integrity by using hash functions. And we will explain some algorithms that use hash functions (by using SHA-256 hash algorithm) to insure data integrity and (something more like) authentication and confidentiality.

*Algorithm1:*

Process file $m_j$ by using a hash function SHA-256 $h$ to calculate hash result $h(m_j)$. Save file $m_j$ in the archive folder and save $y_j=h(m_j)$ in the secure folder of hash results. When you want to read $m_j$ from its original folder then hash $m_j$ by the same hash function $h$ to calculate actual $x_j=h(m_j)$. If $y_j=x_j$ then the file was not changed, if not then the file was changed.



Fig. 2.  Algorithm1

In this algorithm it is required to download the original file and hash result each time from the files storage and the hash storage, which are usually located on server decreasing the effectiveness of the whole reading process. Also there is no confidentiality for the files, or authentication for the source of the file where Man in the Middle attack can be a big threat.

*Algorithm2:*

Process file $m_j$ by using a hash function $h$ to calculate hash result $h(m_j)$ and encrypt it by using private key $k_d$. Save file $m_j$ in the archive folder and save $k_d(y_j)= k_d(h(m_j))$ in the secure folder of hash results. When you want to read $m_j$ from its original folder then hash $m_j$ by the same hash function $h$ to calculate actual $x_j=h(m_j)$. and decrypt $k_d(y_j)$ by using system public key $k_e$ to recover $y_j$. If $y_j=x_j$ then the file was not changed if not then the file was changed.



Fig. 3.  Algorithm 2

*Algorithm 3:*

Process file $m_j$ by using a hash function $h$ to calculate hash result $h(m_j)$ and encrypt it by using secret symmetric key $k$. Save file $m_j$ in the archive folder and save $k(y_j)= k(h(m_j))$ in the secure folder of hash results. When you want to read $m_j$ from its original folder then hash $m_j$ by the same hash function $h$ to

calculate actual $x_j=h(m_j)$. Decrypt $k(y_j)$ by using same symmetric key $k$ to recover $y_j$. If $y_j=x_j$ then the file was not changed, if not then the file was changed.

In this algorithm it is required to download the original file and hash result each time from the files storage and the hash storage, which are usually located on server decreasing the effectiveness of the whole reading process. Also there is no confidentiality for the files. In the other hand, authentication of file source is insured.



Fig. 4.   Algorithm 3

### Algorithm 4:

Pad secret $p$ serial of bits to $m_j$ and then process file $m_j\|p$ by using a hash function $h$ to calculate hash result $h(m_j\|p)$. Save file $m_j$ in the archive folder and save $y_j= h(m_j\|p)$ in the secure folder of hash results. When you want to read $m_j$ from its original folder then pad secret $p$ serial of bits to $m_j$ and hash $m_j\|p$ by the same hash function $h$ to calculate actual $x_j=h(m_j\|p)$. If $y_j=x_j$ then the file was not changed, if not then the file was changed.



Fig. 5.   Algorithm 4

In this algorithm it is required to download the original file and hash result each time from the files storage and the hash storage, which are usually located on server decreasing the effectiveness of the whole reading process. Also there is no confidentiality for the files but the authentication of file source is insured. Additional powerful cryptographic characteristic is fulfilled, where for $m_1= m_2$ then $h(m_1)\neq h(m_2)$.

If we want to make the saved files secret we can apply an additional operation where we encrypt $m_j$ by using symmetric or asymmetric encryption algorithm.

In this paper we will use a special case of the fourth algorithm, where we will use HMAC (Key-Hashed Message Authentication code), which is used as an authentication cryptographic tool.

## IV.   HMAC

The main goals behind the HMAC construction [20] are:

- To use available hash functions without modifications; in particular, hash functions that perform well in software, and for which the code is freely and widely available.

- Preserve the original performance of the hash function without incurring a significant degradation.

- Use and handle keys in a simple way.

- Gain a well-understood cryptographic analysis of the strength of the authentication mechanism based on reasonable assumptions on 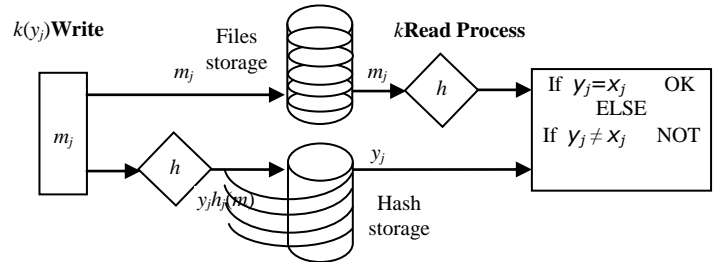the underlying hash function, and to allow easy replacement ability of the underlying hash function if it will be faster or more secure.

HMAC requires a cryptographic hash function, which we denote by $h$, and a secret key $K$. We assume $h$ to be a cryptographic hash function where data is hashed by iterating a basic compression function on $l$ blocks of data. We denote by $b$ the bit-length of such blocks (where $l*b$ equal to the length of $m$ in bits after padding), and by $n$ the bit-length of hash outputs ($n=128$ bits for MD5, $n=160$ bits for SHA-1). The authentication key $K$ can be of any length up to $b$, the block length of the hash function. Applications that use keys longer than $b$ bits will first hash the key using $h$ and then use the resultant $n$ bit string as the actual key to HMAC. In any case the minimal recommended length for $K$ is $n$ bits (as the hash output length).

HMAC can be calculated as follows (Fig. 6):

$$\text{HMAC}_K(m) = h(K^+ \text{ XOR } opad, \ h(K^+ \text{ XOR } ipad, m))$$

It can be done in 7 steps:

*1)   Append zeros to the left end of K to create a b-bit string K+.*

*2)   XOR (bitwise exclusive-OR) K+ with ipad to produce the b-bit block Si.*

*3)   Append m to Si.*

*4)   Apply h to the stream generated in step 3.*

*5)   XOR K+ with opad to produce the b-bit block S0.*

*6)   Append the hash result calculated in Step 4 to S0.*

*7)   Apply h to the stream calculated in step 6 and output the result.*

Because of using such different fixed values of *ipad* and *opad* and doing two times hashing function we avoid the situation where the XORing operation between $K^+$ and *ipad* or $K^+$ and *opad* to have zero's value.

**Keys**

The key for HMAC [21] can be of any length (keys longer than $b$ bits are first hashed using $h$). However, less than $n$ bits is strongly discouraged as it would decrease the security strength of the function. Keys longer than $n$ bits are acceptable but the extra length would not significantly increase the function's strength. A longer key may be advisable if the randomness of the key is considered weak.

Keys need to be chosen randomly (or using a cryptographically strong pseudo-random generator seeded with a random seed), and periodically refreshed. Current attacks do not indicate a specific recommended frequency for key changes as these attacks are practically infeasible. However, periodic key refreshment is fundamental security practice that helps against potential weaknesses of the function as well as the keys, and therefore limits the damage of an exposed key.



Fig. 6. HMAC function

## V. DATA INTEGRITY FILE CHECKING SYSTEM (DIFCS)

We will focus in our research on insuring data integrity by using HMAC [19]. As HMAC is open to use any hash function with it. So in our paper we recommend to use at least SHA-256, which still secure against brute-force attack. In the future we recommend using even hash results with 1024 bits length. In any document management system, each department in the organization has to archive its uploaded files in a central archival warehouse. In the implementation of such solution we will face two important issues: the insuring of data integrity for archived files through transmission and the performance of the network where the transfer of these files is done from the server to the local computers.

Usually each department has an access to its own archived files only and not to the files of the whole archival warehouse. The improved algorithm we developed depends on this factor, that most of the retrieved files requested by the department's user are usually uploaded by the same department.



(a)
Fig. 7. Sending and saving process of file F to SAV



(b)

Fig. 7. Sending and saving process of file F to SAV

In this paper we are implementing an efficient algorithm to insure data integrity and authentication for the archived files and at the same time to insuring better performance for the network. HMAC algorithm will be used to insure data integrity and authentication and a temporary local storage on local PC of most used archival files, which will increase the efficiency of the network.

In the proposed solution uploaded files will be saved in two storage devices: in the local PC of the uploaded user (LPC) and in the Central Archive Server (SAV). Additionally in SAV and LPC we will apply HMAC with a secret key.

**Uploading process:**

When the user uploads F on his LPC, this file is saved in the temporary matrix storage on LPC and it is also sent and saved in SAV server. This saving process is explained in fig. 7, where each file F will have unique identifier $f_{id}$ identifying F in a unique way on LPC and SAV. LPC will calculate hash result $h_{kid}(F)$ for F by using HMAC algorithm and random secret

unique key $k_{id}$ then it will encrypt $h_{kid}$(F) and $k_{id}$ by using User public key LPC$_{kd}$ to insure authentication and then result is encrypted by SAV public key SAV$_{ke}$ to insure confidentiality.

Encrypted result $c$ and F together are sent through network to SAV.

On SAV encrypted $c$ is decrypted by using SAV private key SAV$_{kd}$ and then again decrypted by using LPC public key LPC$_{ke}$ to recover $h_{kid}$(F) and $k_{id}$. By using $k_{id}$ recovered from $c$ SAV calculates hash result $h'_{kid}$(F) for F by using HMAC with key $k_{id}$. SAV compares recovered sent hash result $h_{kid}$(F) of F with the calculated one $h'_{kid}$(F), If they are equal then F and $f_{id}$ and $k_{id}$ are saved on SAV else SAV must sent a request to retransmit all again from LPC.

**Downloading process:**

We will have two situations, when file F with $f_{id}$ and $k_{id}$ exist on LPC, where only LPC will request for $h_{kid}$(F) from SAV, fig 8. (a). And second one when you have only $f_{id}$, Where we need file F with $f_{id}$ and $k_{id}$ and $h_{kid}$(F), fig. 8. (b).

When LPC requires a file F with $f_{id}$ identifier from SAV, the following steps will be done:

*1) Check if file F with $f_{id}$ and $k_{id}$ exist on LPC, If yes then go to 2 else go to 7,*

*2) Send a request to the SAV with $f_{id}$ to retrieve hash result $h_{kid}(F)$,*

*3) SAV Search for $h_{kid}(F)$ according to $f_{id}$,*

*4) SAV sends hash result $h_{kid}(F)$ and $f_{id}$ to LPC m= $h_{kid}(F)$ || $f_{id}$ ),*

*5) LPC retrieve $k_{id}$ and calculates hash result $h'_{kid}(F)$ for F,*

*6) If $h_{kid}(F)= h'_{kid}(F)$ then retrieve F from LPC and*

*end, else go to 7,*

*7) LPC sends request to the SAV with $f_{id}$ to retrieve*

*8) SAV Search for F according to $f_{id}$,*

*9) SAV sends F and hash result $h_{kid}(F)$ and $f_{id}$ and secret $k_{id}$ encrypted by LPC public key LPC$_{ke}$, (m= F || $h_{kid}(F)$ || $f_{id}$ || LPC$_{ke}(k_{id})$),*

*10) LPC receives m= F || $h_{kid}(F)$ || $f_{id}$ || LPC$_{ke}(k_{id})$ and recovers $k_{id}$ by using the private key of LPC$_{kd}$,*

*11) User application on LPC calculates the hash result $h'_{kid}(F)$ for F,*

*12) If $h_{kid}(F)= h'_{kid}(F)$ then retrieve F from LPC and save F and $h_{kid}(F)$ and fid and kid on LPC, else file is corrupted and resend again.*

If additional security is required like confidentiality then symmetric key algorithm is used to insure confidentiality to F. Public key algorithm will be used to exchange the secret key between SAV and user working on the LPC.

In our improved algorithm DIFCS we increased the cryptographic characteristics of the whole process of saving the file and its hash result on server and reading the files and their hash results from the same server. If we will compare the developed algorithm cryptographic characteristics with the other mentioned algorithms in this paper we can easily conclude the following:

*a. In DIFCS algorithm the original file is saved on the local machine so it is not required to download each time the original file from the files storage located usually on server, which increases the effectiveness of the whole archive file retrival process.*

*b. Authentication of file source is insured.*



(a)



(b)

Fig. 8. Downloading process from SAV

*c. Additional powerful cryptographic characteristic is fulfilled, where for if we have two messages m1 and m2, where m1= m2 then h(m1)≠ h(m2).*

*d. Confidentiality for the files or hash results can be implemented according to the user requirements*

## VI. CONCLUSIONS

In this research we developed a new algorithm called DIFCS, which uses HMAC function to insure data integrity and authentication for archival file systems. DIFCS also uses a new technique for retrieving and checking if the archive files are authentic. The main function of DIFCS is to increase the efficiency of the files archival system and the local network. Such an algorithm insures data integrity for archived files and makes them immune against unauthorized manipulation and Man in the Middle attack. It also insures authentication between LPC and SAV.

In future work, we will develop the algorithm to make it a distributed algorithm: where archival files will be distributed and saved in different places according to a known mechanism. Such a development will increase the efficiency of the system.

### REFERENCES

[1] R.C. Merkle, A Certified Digital Signature. In proceedings of Advances in Cryptology, Lecture Notes in Computer Science (435), Springer-Verlag, California, USA, 1989, pp. 218-238.

[2] Menezes A. J., van Oorschot P.C., Vanstone S. A., Handbook of Applied Cryptography. CRC Press, Boca Raton, FL, 1997.

[3] Pieprzyk J., Sadeghiyan B., Design of Hash Algorithms. LNCS 756, Springer, Berlin, 1993.

[4] Wayner P., Digital Cash. AP Professional, Bostan, 1996.

[5] Preneel B, The state of the cryptographic hash functions. Damgård I. (ed.), Lectures on Data Security. Modern Cryptology in Theory and Practice. LNCS 1561, Springer, Berlin, 1999, 158−182.

[6] Qu C., Sebbery J., Pieprzyk J., On the symmetric properties of homogeneous bent functions. Pieprzyk J., Safavi-Naini R., Seberry J. (eds.), Information Security and Privacy. LNCS 1587, Springer, Berlin, 1999, 26−35.

[7] R. Tamassia, N. Triandopoulos, On the Cost of Authenticated Data Structures. In Proc. European Symposium on Algorithms, LNCS (2832), Budapest, Hungary, 2003.

[8] Y.H. Chen, E.J. Lu, Design of a secure fine-grained official document exchange model for e-government, Information & Security 15(1), 2004, pp. 55-71.

[9] J. Woerner, H. Woern, A security architecture integrated co-operative engineering platform for organised model exchange in a Digital Factory environment, Computers in Industry 56(4), 2005, pp. 347-360.

[10] G. Yee, Y. Xu, L. Korba, K. El-Khatib, Privacy and Security in ELearning, Future Directions in Distance Learning and Communication Technologies. Idea Group, Inc. 2006.NRC Publication Number: NRC 48120.

[11] IBM, 2008. Data integrity. Available at:

[12] http://publib.boulder.ibm.com/infocenter/tpfhelp/current/index.jsp?topic=/com.ibm.ztpf-ztpfdf.doc_put.cur/gtps5/s5dint.html (Accessed 12 December 2008).

[13] H. Maruyama, K. Tamura, N. Uramoto, Digest Values for DOM (DOM-HASH), RFC2803. Available at: http://www.landfield.com/rfcs/rfc2803.html (Accessed 13 November 2008).

[14] Bret Mulvey, Evaluation of SHA-1 for Hash Tables, in Hash Functions. Accessed April 10, 2009.

[15] Boritz, J. Efrim. "IS Practitioners' Views on Core Concepts of Information Integrity". International Journal of Accounting Information Systems. Elsevier. http://www.fdewb.unimaas.nl/marc/ecais_new/files/boritz.doc. Retrieved 12 August 2011

[16] Trust and Privacy in Digital Business: Third International Conference, TrustBus 2006, Krakow, Poland, September 4-8, Springer 2006, Proceedings (Lecture Notes in Computer Science / Security and Cryptology, ISBN-13: 978-3540377504.

[17] RFC1321, The MD5 Message Digest Algorithm, R.Rivest, April, 1992.

[18] FIPS-180-1, SHA-1 Secure Hash standard algorithm, April, 1995.

[19] Mihir Bellare Ran Canettiy Hugo Krawczykz, "Keying Hash Functions for Message Authentication", Crypto 96 Proceedings,Lecture Notes in Computer Science Vol. 1109, N. Koblitz ed., Springer-Verlag, 1996.

[20] H. Krawczyk, M.Bellare, R. Canetti HMAC: Keyed-Hashing for Message Authentication, RFC 2104, 1997.

[21] H. Krawczyk, M.Bellare, R. Canetti: Message Authentication using Hash Functions – The HMAC Construction, CryptoBytes, Vol. 2, No. 1 Spring 1996.

[22] NIST FIPS PUB 198, The Keyed-Hash Message Authentication Code (HMAC), Federal Information Processing Standards Publication Issued March 6, 2002.

[23] RFC4868, Using HMAC-SHA-256, HMAC-SHA-384, and HMAC-SHA-512 with IPsec, S. Kelly, S. Frankel, May,2007.

# A Greedy Algorithm for Load Balancing Jobs with Deadlines in a Distributed Network

Ciprian. I. Paduraru

Department of Computer Science,
University of Bucharest
Bucharest, Romania

*Abstract*—**One of the most challenging issues when dealing with distributed networks is the efficiency of jobs load balancing. This paper presents a novel algorithm for load balancing jobs that have a given deadline in a distributed network assuming central coordination. The algorithm uses a greedy strategy for global and local decision making: schedule a job as late as possible. It has an increased overhead over other well-known methods, but the load balancing policy provides a better fit for jobs.**

*Keywords—scheduling; greedy; coordination; network*

## I. INTRODUCTION

Distributed architecture of computers can represent the underlying of a web or network-based service used for processing user requests. In this context, the performance of the processing system is closely related to user experience and service availability, and can therefore play an important role in the success or failure of the respective service on the market. As sufficient hardware resources for processing a large number of requests are generally expensive, a good algorithm for the distribution of load - between the processing units in the distributed system - is necessary to save costs in addition to increase clients' satisfaction.

This paper proposes an algorithm for load balancing of jobs in a distributed network assuming central coordination. Jobs are non-preemptive, received by a single machine in the distributed network (master) and sent to workers. Each job has a given deadline which is assigned by the owner of the request. The master must decide if the job can be executed by one of the workers considering its deadline and an error window, and if it does, then who the best worker to execute it is. Once received by a worker, it must decide where on its own waiting list of jobs the new job should be added. The algorithm can work both for homogeneous and heterogeneous workers. It all depends on the ability of a worker to determine the execution time of a job. If workers can estimate how much time it will take to execute a given job within the considered error window then we can have heterogeneous machines in the distributed system. Various methods for doing such estimation are presented in [1]. One method is to assign to each job a length-class and test each worker in the system how much time it will take to execute each kind of length-class. There is a linear search overhead determining the best fit for a new job but it makes the load balancing better and provides better results. Also, we assume that there is a communication link between the master machine and each worker, and we can

estimate the average communication time for each job. For simplicity, this communication time is included in the execution time of a job.

The rest of the paper is organized as follows: In Section 2 there is a discussion about research made on load balancing or scheduling algorithms with deadlines. Section 3 presents how the algorithm is designed and a pseudocode for its implementation. Results obtained from running a simulator over some test samples are given in Section 4. Conclusions are presented in the last section.

## II. RELATED WORK

At the time when this paper is written there is no paper dealing with load balancing tasks with deadlines in a distributed network with central coordination. However, there are various papers presenting techniques for load balancing / scheduling of tasks which are a point of inspiration and a possible comparison for the algorithm presented here.

Some theoretical aspects with high-importance for this paper are presented in [5]. A conclusion is that Earliest DeadLine First (EDF) policy is not optimal for non-preemptive tasks or when there are multiple processors in a system. [2] Presents a new algorithm for load balancing in grid architecture for fair scheduling. It addresses the fairness issues by using mean waiting time. It schedules the tasks by using fair completion time and reschedules them by using mean waiting time of each task to obtain load balance. In [3], there is a comparison between two important task schedulers such as EDF scheduler and Ant Colony Optimization Based (ACO) scheduler. Paper [4] presents a greedy algorithm for scheduling jobs with deadlines and profits with the main objective to maximize the profit, for a single processing unit.

## III. ALGORITHM DESIGN AND IMPLEMENTATION

Jobs are received and sent further by the master machine. The algorithm doesn't move any job from a worker to another because each time when we give a job to a worker, we know that it can satisfy its deadline constraint. Also, as Section 1 states, jobs are non-preemptive. There are two separated views of the algorithm:

- Master view: responsible for assigning a new job to a worker if there is one available to execute this job satisfying its deadline constraint.

- Worker view: responsible for managing a data structure that holds jobs and extracting / executing these jobs.

A data type that defines a job can be defined as: *JobType*={timeToExecute, deadline, timeToStart, timeToEnd, dataContext}. *timeToExecute* is the time needed by an worker to execute the job, while *dataContext* is the data associated with the job execution. *timeToStart* represents the time when a job can start on a worker. *timeToEnd* is the computed value of timeToStart + timeToExecute.

### A.  MASTER VIEW

At this level, the main idea is to send a new job to the worker which can start it as late as possible but still satisfying the deadline constraint. It is a greedy solution which can keep workers available for earlier deadlines. Considering that function *GetTimeToStart* returns the time when a worker can start a job given as parameter (-1 is considered to be the return result for not being able to execute it and satisfy its deadline goal) then the pseudocode that master runs when a new job is received is presented below.

```
OnNewTaskArrived( JobType task)
    bestWorkerId = -1
    bestWorkerTime = 0
    foreach worker W do
    {
        Wtime = Controller[W]->GetTimeToStart()
        if  Wtime != -1 AND  bestWorkerTime < Wtime
        {
            bestWorkerTime = Wtime
            bestWorkerId = W
        }

        if bestWorkerId != -1
        {
            SendTask(job, bestWorkerId)
        }
        else
        {
            // Code for refused job
        }
}
```

The *Controller* array is stored on master and represents the state of each worker − the data structure which stores informations about the currently assigned jobs for each client, excluding the "dataContext" field which is only needed by workers. This is actually logic part of worker's view.

### B.  WORKER VIEW

There are two issues at worker's view: the *GetTimeToStart* function implementation (called by the master and having its context data stored on the master in the *Controller* array) and how a worker manages its internal data structure to execute jobs.

Same greedy idea as in section 3.1 is used here: schedule a new job as late as possible (Figure 1). Workers are using a linked list to store the assigned jobs. In this linked list, jobs are sorted in ascending order by the value of field *timeToStart*. The value of this field for a new job should ideally be: deadline − timeToExecute, because the main objective is to promote free spaces for new jobs that have earlier deadlines. But if this is not possible due to existing jobs, then it should be set to the last gap found between assigned jobs that allows us to execute the new job and satisfy its deadline.
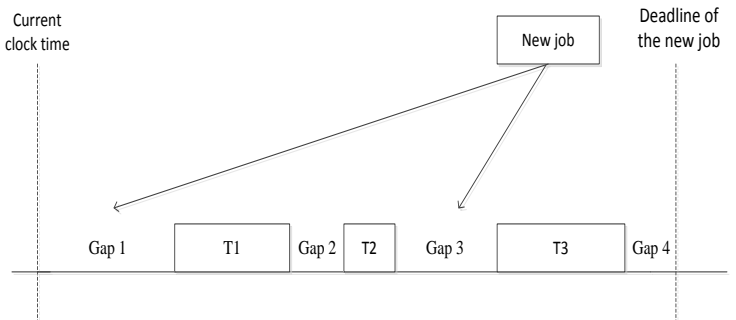


Fig. 1.  Adding a new task to an existing list of jobs (T1, T2 and T3). Considering that the width of the rectangle represents the execution time of a job, then gaps with index 1 and 3 can fit the new job. The end of gap 3 will be preferred for the new job to schedule it as late as possible.

Pseudo-code for this operation is presented below. The

```
GetTimeToStart(job)
  // Step 1:
  // Check the back of the list first
  If mList.isEmpty()  OR
   mList.last().timeToEnd<=(job.deadlinejob.timeToExec)
  {
    job.timeToStart=job.deadline – job.timeToExec;
    mLastCachedPos = mList.end;
    return job.timeToStart;
  }


  // Step 2:
  // Check for gaps between tasks, starting from last to first
  foreach job T in mList (reverse order)
  {
    prevT = T->prev
    if  prevT == NULL continue

    deadlineImp=min(job.deadline,T.timeToStart)
    if prevT.timeToEnd+job.timeToExec<= deadlineImp
    {
      job.timeToStart = deadlineImp – job.timeToExec
      mLastCachedPos = position of T in mList
      return newJob.timeToStart
    }
  }

  // Step 3:
  // Check for a gap between current clock time and first job
  // begin
  Tfirst = mList.first()
  deadlineImp = min(job.deadline, Tfirst.timeToStart)
  if (clock() <= (deadlineImp – newJob.timeToExec))
  {
    mLastCachedPos = mList.first()
    job.timeToStart = Tfirst.timeToStart – job.timeToExecute
    return  job.timeToStart
  }
```

*mList* variable represents the linked list where the jobs are stored. *mList.end/mList.start* represents the last/first element in the list. If there is no other assigned job in the list or we can schedule the new job after the last one, then the ideal value for *timeToStart* will be deadline – timeToExecute. Otherwise, the algorithm tries then to fill the first gap found starting from the end of the list and going to its beginning. Each time we compare two consecutive elements in the list (T and *prevT)* and check if the new job can be added between the *timeToEnd* of prevT and the minimum between its deadline and the *timeToStart* of T.  If we find such a position then we set the *timeToStart* as late as possible in this gap. Finally, if no gap was found yet, we try to add it in the gap starting from current clock time to the beginning of the first job in the list. The implementation of *GetTimeToStart* will also cache the *timeToStart, timeToEnd* and the position in the linked list where it should be added (mLastCachedPos). This information will be sent together with the job in the SendTask function to avoid doing the linear search twice.

On each worker there is a function which continuously polls for jobs in the jobs list and grabs them for execution. The trick is to allow grabbing and execution of the first job from the list even if the current clock time is less than its *timeToStart* field. Doing this will keep the machines busy. It is possible that some of the jobs with a deadline close to current clock time will be refused, but it has the same probability (assuming a normal distribution of jobs in time) that other new jobs will benefit from this.

The complexity of searching for the best place to add a new job inside a worker is linear in the number of existing jobs on that worker. The worst case happens where there are many jobs received in a short time interval while the jobs execution time is higher than the arrival time rate of new jobs.

## IV. SIMULATION RESULTS

To test the performance of the proposed algorithm, a simulation was made in order to see its behavior in comparison with other two load balancing algorithms. The first one sends the jobs in a round robin policy while the second one to the worker that can execute the job as late as possible. Both have just a simple policy at the worker's view: add the new job to the end of the queue if it can be executed before        its deadline expires. Ideally, the load balancing should use the resources correctly by keeping the hardware busy most of the time and minimizing the number of refused jobs.

Final results were obtained by averaging a number of test samples which creates 1000 of jobs with a normal distribution of execution times between 10 and 200 milliseconds. The arrival time rate of new jobs was between 1 and 10 milliseconds. Deadline time extension of each job (time since a new job was received to when it should finish) was also chosen by a normal distribution in interval [10, 2000] milliseconds (considering the job execution time too). Samples where run on 24, 16, 12, 8 and 4 machines in a local network (workers) each having single hardware process dedicated for our job execution. The process of receiving and assigning a new job to a worker was done by a separate machine called master.

Figure 1. shows how many jobs where refused depending on the number of machines and the algorithm used. The results graph shows that the proposed algorithm is better than the other two methods, despite its overhead. The difference between it and the other two algorithms increases with the number of machines used. When using 16 machines, the number of jobs refused by the other two algorithms is with 49% higher than the proposed algorithm.

With 24 machines, the proposed algorithm succeeded to obtain 0 refused jobs while the other two solutions had 18/21 refused jobs.  The samples used creates jobs in a short time interval (defined at the beginning of this section) to represent a worst case scenario for the proposed algorithm. When jobs have longer execution time and the arrival time has a different time distribution in the proposed algorithm can perform even better than this because there is less overhead spent on decision making.

Fig. 2. The average number of refused jobs (1000 was the total number of jobs) in test samples by each algorithm.

## V. CONCLUSION AND FUTURE WORK

This paper presented a novel algorithm for load balancing jobs having deadlines in a distributed network with central coordination. By using two different greedy strategies, one from master view and the other from worker's view, the proposed algorithm provide an increased performance than the classical methods for load balancing jobs. Keeping the same hardware and being able to increase the performance with over 49%, as the results sections shows, represents an important issues for most of the web services on the market.

One important topic to study in continuation is to consider that each job also has a profit assigned and find a load balancing policy that maximize the profit instead of considering equal profits for jobs like the proposed algorithm does.

REFERENCES

[1] Ciprian Paduraru, "A New Online Load Balancing Algorithm in Distributed Systems", 14th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing, Timisoara, Romania, 2012

[2] U.Karthick Kumar, "A Dynamic Load Balancing Algorithm in Computational Grid Using Fair Scheduling", IJCSI International Journal of Computer Science Issues, Vol. 8, Issue 5, No 1, September 2011

[3] M.Kaladevi and S.Sathiyabama, "A Comparative Study of Scheduling Algorithms for Real Time Task", International Journal of Advances in Science and Technology,

Vol. 1, No. 4, 2010

[4] Antonina Kolokolova, "A Greedy Algorithm for Scheduling Jobs with Deadlines and Profits", Scheduling case study, Lecture notes

[5] C. L. LIU and JAMES W. LAYLAND, "Scheduling Algorithms for Multiprogramming in a Hard Real-Time Environment", Journal of the Associatlon ior Cornputmg Machinery, Vol. 20, No. I, January 1973

[6] Peter Brucker, "Scheduling Algorithms Fifth Edition", Springer, October 2006.

[7] Kirk Schloegel, George Karypis and Vipin Kumar, "A unified algorithm for load-balancing adaptive scientific simulations", Proceeding Supercomputing '00 Proceedings of the 2000 ACM/IEEE conference on Supercomputing Article No. 59 IEEE Computer Society Washington, DC, USA

[8] Abbas Karimi, Faraneh Zarafshan, Adznan b. Jantan, A.R. Ramli, M. Iqbal b.Saripan, "A New Fuzzy Approach for Dynamic Load Balancing Algorithm", (IJCSIS) International Journal of Computer Science and Information Security, Vol. 6, No. 1, 2009.

[9] Reinhard Lüling, Burkhard Monien, "A Dynamic Distributed Load Balancing Algorithm with Provable Good Performance", Proceedings of the 5th Annual ACM Symposium on Parallel Algorithms and Architectures, 1993.

# Evaluation of Different Hypervisors Performance in the Private Cloud with SIGAR Framework

P. Vijaya Vardhan Reddy

Department of Computer Science & Engineering
University College of Engineering, Osmania University
Hyderabad, India

Dr. Lakshmi Rajamani

Department of Computer Science & Engineering
University College of Engineering, Osmania University
Hyderabad, India

*Abstract*— **To make cloud computing model Practical and to have essential characters like rapid elasticity, resource pooling, on demand access and measured service, two prominent technologies are required. One is internet and second important one is virtualization technology. Virtualization Technology plays major role in the success of cloud computing. A virtualization layer which provides an infrastructural support to multiple virtual machines above it by virtualizing hardware resources such as CPU, Memory, Disk and NIC is called a Hypervisor. It is interesting to study how different Hypervisors perform in the Private Cloud. Hypervisors do come in Paravirtualized, Full Virtualized and Hybrid flavors. It is novel idea to compare them in the private cloud environment. This paper conducts different performance tests on three hypervisors XenServer, ESXi and KVM and results are gathered using SIGAR API (System Information Gatherer and Reporter) along with Passmark benchmark suite. In the experiment, CloudStack 4.0.2 (open source cloud computing software) is used to create a private cloud, in which management server is installed on Ubuntu 12.04 – 64 bit operating system. Hypervisors XenServer 6.0, ESXi 4.1 and KVM (Ubuntu 12.04) are installed as hosts in the respective clusters and their performances have been evaluated in detail by using SIGAR Framework, Passmark and NetPerf.**

*Keywords—CloudStack; Hypervisor; Management Server; Private Cloud; Virtualization Technology; SIGAR; Passmark*

## I. INTRODUCTION

Cloud computing is a model for enabling convenient, on-demand network access to a shared pool of configurable computing resources such as networks, servers, storage, applications, and services that can be rapidly provisioned and released with minimal management effort or service provider interaction [1].

Virtualization, in computing, refers to the act of creating a virtual version of something, including but not limited to a virtual computer hardware platform, operating system, storage device, or computer network resources. Storage virtualization is amalgamation of multiple network storage devices into what appears to be a single storage unit. Server virtualization is partitioning of a physical server into smaller virtual servers. Operating system-level virtualization is a type of server virtualization technology which works at the operating system (kernel) layer. Network virtualization is using network resources through a logical segmentation of a single physical network. Virtualization is the technology which increases the

utilization of physical servers and enables portability of virtual servers between physical servers. Virtualization Technology gives the benefit of work load isolation, work load migration and work load consolidation.

For being able to reduce hardware cost, cloud computing uses virtualization. Virtualization technology has evolved really quickly during past few years. Also it is particularly due to hardware progresses made by AMD and Intel. Virtualization is a technology that combines or divides computing resources to present one or many operating environments using methodologies like hardware and software partitioning or aggregation, partial or complete machine simulation, emulation, timesharing, and many others [2]. A virtualization layer provides an infrastructural support using the lower-level resources to create multiple virtual machines that are independent and isolated from each other. Such a virtualization layer is also called Hypervisor. [2].

Cloud computing allows customers to reduce the cost of the hardware by allowing resources on demand. Also customers of the service need to have guaranty of the good functioning of the service provided by the cloud. The Service Level Agreement brokered between the providers of cloud and the customers is the guarantees from the provider that the service will be delivered properly [3].

This paper provides a quantitative comparison of three hypervisors Xen Server 6.0, VMware ESXi Server 4.1 and KVM (Ubuntu 12.04) in the private cloud environment. Microsoft Windows 2008 R2 server is installed on three hypervisors as a guest operating system and a series of performance experiments are conducted on the respective guest OS and results are gathered using SIGAR [36], Passmark [16] and NetPerf [35]. This technical paper presents and analyses the results of these experiments. The discussion in this paper should help both IT decision makers and end users to choose the right virtualization hypervisor for their respective private cloud environments. The experimental results indicate that both XenServer and VMware ESXi Server deliver almost equal and near native performance in all the tests except in CPU test ESXi is performing marginally better than XenServer and in Memory test XenServer performing slightly better than that of ESXi Server. Furthermore, KVM performance is noticeably lower than that of XenServer and ESXi Server, hence it needs to improve in all the performance aspects.

## II. Virtualization Techniques

This section describes the different virtualization techniques namely, Full virtualization and Paravirtualization used by different hypervisors.

X86 operating systems are designed to run directly on the bare-metal hardware, so they naturally assume they fully 'own' the computer hardware. The x86 architecture offers four levels of privilege known as Ring 0, 1, 2 and 3 to operating systems and applications to manage access to the computer hardware. While user level applications typically run in Ring 3, the operating system needs to have direct access to the memory and hardware and must execute its privileged instructions in Ring 0. Virtualizing the x86 architecture requires placing a virtualization layer under the operating system (which expects to be in the most privileged Ring 0) to create and manage the virtual machines that deliver shared resources. Three alternative techniques now exist for handling sensitive and privileged instructions to virtualize the x86 Architecture. Full virtualization [17] approach, translates kernel code to replace non-virtualizable instructions with new sequences of instructions that have the intended effect on the virtual hardware. This combination of binary translation and direct execution provides Full virtualization as the guest OS is fully abstracted (completely decoupled) from the underlying hardware by the virtualization layer. The full virtualization approach allows datacenters to run an unmodified guest operating system, thus maintaining the existing investments in operating systems and applications and providing a non-disruptive migration to virtualized environments. VMware ESXi server uses a combination of direct execution and binary translation techniques [4] to achieve full virtualization of an x86 system. Paravirtualization [17], involves modifying the OS kernel to replace non-virtualizable instructions with hyper-calls that communicate directly with the virtualization layer hypervisor. The hypervisor also provides hyper-call interfaces for other critical kernel operations such as memory management, interrupt handling and time keeping. The paravirtualization approach modifies the guest operating system to eliminate the need for binary translation. Therefore it offers potential performance advantages for certain workloads but requires using specially modified operating system kernels [4]. The Xen open source project was designed initially to support paravirtualized operating systems. While it is possible to modify open source operating systems, such as Linux and OpenBSD, it is not possible to modify "closed" source operating systems such as Microsoft Windows. Hardware vendors are rapidly embracing virtualization and developing new features to simplify virtualization techniques. First generation enhancements include Intel Virtualization Technology (VT-x) and AMD's AMD-V which both target privileged instructions with a new CPU execution mode feature that allows the VMM to run in a new root mode below ring 0. The hardware virtualization [17] support enabled by AMD-V and Intel VT technologies introduces virtualization in the x86 processor architecture itself.

## III. Hypervisor Models

All three hypervisors which used in the experiment are discussed from viewpoint of their virtualization technique.

### A. Paravirtualized Hypervisor

*XenServer* - Citrix XenServer is an open-source, complete, managed server virtualization platform built on the powerful Xen Hypervisor. Xen [21] uses para-virtualization. Para-virtualization modifies the guest operating system so that it is aware of being virtualized on a single physical machine with less performance loss. XenServer is a complete virtual infrastructure solution that includes a 64-bit Hypervisor with live migration, full management console, and the tools needed to move applications, desktops, and servers from a physical to a virtual environment [8]. Based on the open source design of Xen, XenServer is a highly reliable, available, and secure virtualization platform that provides near native application performance [8]. Xen usually runs in higher privilege level than the kernels of guest operating systems. It is guaranteed by running Xen in ring 0 and migrating guest operating systems to ring 1. When a guest operating system tries to execute a sensitive privilege instruction (e.g., installing a new page table), the processor will stop and trap it into Xen [9]. In Xen, guest operating systems are responsible for allocating the hardware page table, but they only have the privilege of direct read, and Xen [9] must validate updating the hardware page table. Additionally, guest operating systems can access hardware memory with only non-continuous way because Xen occupies the top 64MB section of every address space to avoid a TLB flush when entering and leaving the Hypervisor [9]. XenServer is a complete virtual infrastructure solution that includes a 64-bit Hypervisor [8].

### B. Full virtualized Hypervisor

*ESXi Server* - VMware ESXi is a Hypervisor aimed at server virtualization environments capable of live migration using VM motion and booting VMs from network attached devices. VMware ESXi supports full virtualization [7]. The Hypervisor handles all the I/O instructions, which necessitates the installation of all the hardware drivers and related software. It implements shadow versions of system structures such as page tables and maintains consistency with the virtual tables by trapping every instruction that attempts to update these structures. Hence, an extra level of mapping is in the page table. The virtual pages are mapped to physical pages throughout the guest operating system's page table [6]. The Hypervisor then translates the physical page (often-called frame) to the machine page, which eventually is the correct page in physical memory.

This helps the ESXi server better manage the overall memory and improve the overall system performance [19]. VMware's proprietary ESXi Hypervisor, in the vSphere cloud-computing platform, provides a host of capabilities not currently available with any other Hypervisors. These capabilities include High Availability (the ability to recover virtual machines quickly in the event of a physical server failure), Distributed Resource Scheduling (automated load balancing across a cluster of ESXi servers), Distributed Power Management (automated decommissioning of unneeded servers during non-peak periods), Fault Tolerance (zero downtime services even in the event of hardware failure), and Site Recovery Manager (the ability to automatically recover virtual environments in a different physical location if an entire datacenter outage occurs) [7].

## C. Hybrid methods

*KVM* - KVM (Kernel-based Virtual Machine) is another open-source Hypervisor using full virtualization apart from VMware. And also as a kernel driver added into Linux, KVM enjoys all advantages of the standard Linux kernel and hardware-assisted virtualization thus depicting hybrid model. KVM introduces virtualization capability by augmenting the traditional kernel and user modes of Linux with a new process mode named guest, which has its own kernel and user modes and answers for code execution of guest operating systems [9]. KVM comprises two components: one is the kernel module and another one is userspace. Kernel module (namely kvm.ko) is a device driver that presents the ability to manage virtual hardware and see the virtualization of memory through a character device /dev/kvm. With /dev/kvm, every virtual machine can have its own address space allocated by the Linux scheduler when being instantiated [9]. The memory mapped for a virtual machine is actually virtual memory mapped into the corresponding process. Translation of memory address from guest to host is supported by a set of page tables. KVM can easily manage guest Operating systems with kill command and /dev/kvm. User-space takes charge of I/O operation's virtualization. KVM also provides a mechanism for user-space to inject interrupts into guest operating systems. User-space is a lightly modified QEMU, which exposes a platform virtualization solution to an entire PC environment including disks, graphic adapters and network devices [9]. Any I/O requests of guest operating systems are intercepted and routed into user mode to be emulated by QEMU [9].

## IV. RELATED WORK

The following papers are studied to understand about the relevant work which had happened in the selected research area.

Benchmark Overview - vServCon a white paper by FUJITSU [10], scalability measurements of virtualized environments at Fujitsu Technology Solutions are currently accomplished by means of the internal benchmark "vServCon" (based on ideas from Intel's "vConsolidate"). The abbreviation "vServCon" stands for: "virtualization enables SERVer CONsolidation. A representative group of application scenarios is selected in the benchmark. It is started simultaneously as a group of VMs on a virtualization host when making a measurement. Each of these VMs is operated with a suitable load tool at a defined lower load level. All known virtualization benchmarks are thus based on a mixed approach of operating system and applications plus an "idle" or "standby" VM, which represents the inactive phases of a virtualization environment and simultaneously increases the number of VMs to be managed by the Hypervisor [10].

The virtualization overhead involves performances depreciation rather to native performances. Research have been made to measure the overhead of the virtualization for different hypervisor such as XEN, KVM and VMware ESX [11]; [12]; [13]; [14]; [15]. For their researches Menon used a toolkit called Xenoprof which is a system wide statistical tool implemented specially for Xen [13]. Due to this toolkit they have managed to analyse the performances of the overhead of network I/O devices. Their study has been performed within

uniprocessor as well as multiprocessor. A part of their research has been dedicated to performance debugging of Xen using Xenoprof. Those researches have permitted to correct bugs and improve by that the network performances significantly. After the debugging part it has been focused on the network performances. It has been observed that the performance seems to be almost the same between Xen Domain0 and native performances. However if the number of interfaces increase, the receive throughput of the domain0 is significantly smaller than the native performances. This degradation of network performances is cause by an increasing CPU utilisation. Because of the overhead caused by the virtualization there are more instructions that need to be managed by the CPU. This involves more information to treat and bufferization by the CPU which cause a degradation of receive throughput compared to native performances. More recent studies try to compare the differences between hypervisors and especially the performances of each one according to their overhead [12];[15]. They are using three different benchmark tools to measure the performances: LINPACK, LMbench and Iozone. Their experiment is divided in three parts according to the specific utilisation of each tool. With LINPACK Jianhua had tested the processing efficiency on floating point. Different pick value has been observed over the different systems tested which are native performance, Xen and KVM. The result of this show that the processing efficiency of Xen on floating point is better than KVM because Fedora 8 virtualized with Xen have performances which represent 97.28% of the native rather than Fedora 8 virtualized with KVM represent only 83.46% of the native performances. The virtualization of Windows XP comes up with better performances than with the virtualization of fedora 8 on Xen. This is explained by the authors by the fact that Xen own fewer enhancement packages for windows XP than for fedora 8because of that the performances of virtualized windows XP are slightly better than virtualized fedora 8.

After having testing the processing efficiency with LINPACK, Jianhua have analysed memory virtualization of Xen and KVM compared to native memory performances with LMbench. It has been observed that the memory bandwidth in reading and writing of Xen are really close to native performances. However the performances of KVM are slightly slower for reading but significantly slower concerning the writing performances. The last tool used by Jianhua is IOzone which is used to perform file system benchmark. Once again the native performances are compared to the virtualization performances of Xen and KVM. Without Intel-VT processor the performances of either Xen or KVM are around 6 or 7 times slower than the native performances. However within the Intel-VT processor the performances of Xen increase significantly because the performances are even better than native performances. However KVM does not exploit the functionalities of the Intel-VT processors and because of that does not improve its performances.

After analysing the relevant work on hypervisors performance we have chosen the below experimentation to compare the respective hypervisors in the private cloud environment with CloudStack using SIGAR framework which is a novel idea.

## V. TEST METHODOLOGY - PRIVATE CLOUD: CLOUDSTACK WITH HYPERVISORS

In our experiment, the proposed test environment contains following infrastructure using open source cloud computing software. CloudStack is an Infrastructure as a service (IaaS) cloud based software which is able to rapidly build and provide private cloud environments or public cloud services. Supporting KVM, XenServer and Vmware ESXi, CloudStack is able to build cloud environments with a mix of multiple different hypervisors. With rich web interface for users and administrators with operations of cloud use and operation being performed on a browser. Additionally, the architecture is made to be scalable for large-scale environments [22]. CloudStack is open source software written in java that is designed to deploy and manage large networks of virtual machines, as a highly available, scalable cloud computing platform. CloudStack offers three ways to manage cloud computing environments: an easy-to-use web interface, command line and a full-featured RESTful API [22]. Private clouds are deployed behind the firewall of a company where as public cloud is usually deployed over the internet. It is always ideal to use open source solutions to perform any experiment related to cloud computing.

In our test environment XenServer, ESXi and KVM are used as hypervisors (Hosts) in the CloudStack (private cloud). One machine is Management Server, runs on a dedicated server. It controls allocation of virtual machines to hosts and assigns storage and IP addresses to the virtual machine instances. The Management Server runs in a Tomcat container and requires a MySQL database for persistence. In the experiment, Management Server is installed on Ubuntu (12.04 64-bit). On the host servers XenServer 6.0, ESXi 4.1 and KVM (Ubuntu 12.04) [31] hypervisors are installed as depicted in Fig. 1. Front end will be any base machine to launch CloudStack UI using web interface (with any browser software IE, Firefox, Safari) to provision the cloud infrastructure by creating zone, pod, cluster and host in the sequential order. After respective hypervisors are in place, guest OS Windows 2008 R2 64-bit [33] installed on them to carry out all performance tests.



Fig. 1. Test Environment Architecture – Private Cloud (CloudStack with Multiple hypervisors)

A typical enterprise datacenter runs a mix of CPU, memory, and I/O-intensive applications. Hence the test workloads chosen for these experiments comprise several well-known standard benchmark tests. Passmark, a synthetic suite of benchmarks intended to isolate various aspects of workstation performance, was selected to represent desktop-oriented workloads. Disk I/O performance is measured using Passmark. CPU and Memory performance on the guest OS are measured using SIGAR Framework. SIGAR (System Information Gatherer and Reporter) is a cross-platform, cross-language library and command-line tool for accessing operating system and hardware level information in Java, Perl and .Net. In the experiment, Java program has written to gather system information using SIGAR API by deploying sigar-amd64-winnt.dll for Windows. And for network performance Netperf is used in the experiment. Netperf was used to simulate the network usage in a datacenter. The objective of these experiments was to test the performance of the three virtualization hypervisors. The tests were performed using a Windows 2008 R2 64-bit as guest operating system. The benchmark test suites are used in these experiments only to illustrate performance of the three hypervisors.

## VI. RESULTS

This section provides the detailed results for each of the benchmarks run. Disk I/O and Network Performance results have been normalized to native performance measures. Native performance is normalized at 1.0 and all other various benchmark results are shown relative to that number. Hence benchmark results of 90% of the native performance would be shown as 0.9 on the scale in the graph. Higher numbers indicate better performance of the particular virtualization platform, unless indicated otherwise. Near-native performance also indicates that more virtual machines can be deployed on a single physical server, resulting in higher consolidation ratios. This can help even if an enterprise plans to standardize on virtual infrastructure for server consolidation alone. CPU utilization tests indicate lower CPU utilization is better for a hypervisor, which is evaluated by using SIGAR API. In case of Memory tests, High available memory indicated better performance of a hypervisor which gathered using SIGAR.

### A. SIGAR

CPU utilization on the guest Operating System is captured when it is running on the respective Hypervisor. CPU utilization details are captured through java program using SIGAR API on the guest OS for each hypervisor. As shown in Fig. 2, ESXi for its guest OS shows less utilization of CPU as compared to other hypervisors. Lower utilization CPU indicates the better performance for a hypervisor. XenServer also shows low utilization of CPU for its guest OS but little higher than ESXi hypervisor. On the other hand KVM's CPU utilization is slightly high for its guest OS as compared to other two hypervisors.

Memory performance is evaluated by considering the available memory in the respective hypervisor when the single guest Operating Systems is given full available memory.

Fig. 2. CPU Utilization captured using SIGAR (Lower value is better)

Fig. 3 shows Available memory on the respective hypervisor when guest OS is running. Memory details are captured using Java program with SIGAR API on the guest OS. XenServer for its guest OS shows maximum available memory as compared to other hypervisors. Higher available memory indicates better performance for a hypervisor. ESXi also exhibits higher available memory only but slightly less compared to XenServer. KVM indicates marginally less available memory compare to other hypervisors.



Fig. 3. Available Memory captured using SIGAR (Higher Value is better)

## B. PASSMARK

The following Fig. 4 shows benchmark results for Passmark Disk I/O read write tests. Sequential Read and Sequential Write are the disk mark tests which were conducted on the three hypervisors in the private cloud environment. Both XenServer and ESXi perform almost equal to native performance.



Fig. 4. Passmark – Disk I/O Read Write results compared to native (Higher values are better)

In Sequential Read and Sequential Write XenServer slightly shows better performance than that of VMWare ESXi Server. In overall disk mark performance XenServer shows 2.7% overhead vs native whereas ESXi shows 3.4% overhead vs native. KVM significantly falls behind other two hypervisors and native as well.

## C. NETPERF

For experiment, in the private cloud for all the three hypervisors, Netperf test involved running single client communicating with single virtual machine through a dedicated physical Ethernet adapter and port. All tests are based on the Netperf TCP_STREAM test. Fig. 5 shows the Netperf results for send and receive tests. XenServer and ESXi demonstrated near native performance in Netperf test, while KVM lags behind other hypervisors and native.



Fig. 5. Netperf results compared to native (higher values are better)

## VII. DISCUSSION ON RESULTS

Performance results show convincingly that XenServer and ESXi Server both perform equally well in all experiments close to near native performance without showing the signs of any virtualization overhead except KVM falling behind other two hypervisors and native as well.

In CPU utilization tests ESXi CPU utilization is 0.06% less than that of XenServer and 0.24% less than that of KVM thus exhibiting better performance in CPU utilization. In memory tests XenServer available memory is 1% more than that of ESXi Server and 6% more than that of KVM hence showing better memory performance among two other hypervisors. In I/O tests XenServer scores over ESXi and KVM, where XenServer shows 4% overhead in sequential read and 6% overhead in sequential write as compared to native. ESXi shows 5% overhead in sequential read and 7% overhead in sequential write as compared to native. And KVM shows 35% overhead in sequential read and 36% overhead in sequential write as compared to native. In Network performance tests both XenServer and ESXi gives near native performance and KVM falls marginally behind other two hypervisors. In Client-Receive tests both XenServer and ESXi gives performance equal to native and in Client-Send tests XenServer gives equal to native performance but ESXi shows 3% overhead as compared to native. In Client-Send and Client-Receive tests KVM shows 22% overhead as compared to native.

On overall XenServer and ESXi two hypervisors are reliable, affordable and offer the windows or any other guest operating system IT professional a high performance platform for server consolidation for production workloads. KVM needs to improve up on almost all fronts if it has to become on par with other two hypervisors. ESXi and XenServer are matured hypervisors as compare to KVM and their Reliability, Availability and Serviceability (RAS) is significantly higher than that of KVM.

## VIII. CONCLUSION AND FUTURE WORK

The objective of this experiment is to evaluate the performance of VMWare ESXi Server, XenServer and KVM Hypervisors in the private cloud environment. After evaluation results indicate that XenServer and ESXi hypervisors exhibit impressive performance in comparison with KVM. Virtualization infrastructure should offer certain enterprise readiness capabilities such as maturity, ease of deployment, performance, and reliability. From the test results VMware ESXi Server and XenServer are better equipped to meet the demands of an enterprise datacenter than the KVM hypervisor. And KVM needs significant improvement to become an enterprise ready hypervisor. The series of tests conducted for this paper proves that VMware ESXi Server and XenServer delivers the production-ready performance needed to implement an efficient and responsive datacentre in the private cloud environment.

The performance tests are conducted in the private cloud with 64-bit Windows guest operating system. While evaluating network performance, one client send and receive tests are performed on three hypervisors which are supported by CloudStack private cloud platform. The future work can include multiple client send and receive network tests for hypervisors. Experiments can also be carried out with paravirtualized Linux guest operating system as well. With more workloads scalability tests can be performed with other hypervisors which are not covered in the present experiment. And future work can also consider public cloud environment for experimentation.

### REFERENCES

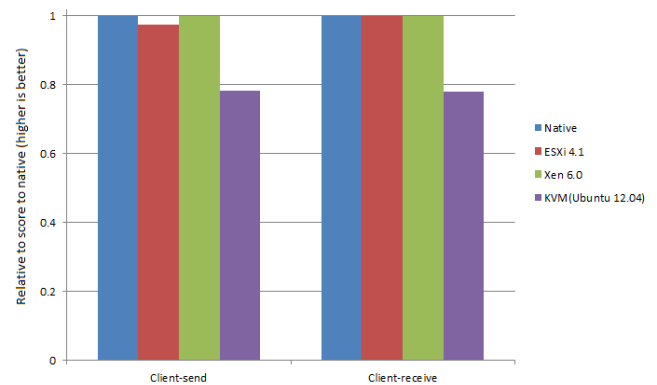[1] Mell, P. & Grance, T. (2009) The NIST Definition of Cloud Computing. Version 15, 10-7-09. National Institute of Standards and Technology, Information Technology Laboratory.

[2] Nanda, S., T. Chiueh, ―A Survey on Virtualization Technologies, Technical report, Department of Computer Science, SUNY at Stony Brook, New York, 11794-4400, 2005.

[3] Buyya, R., Yeo, C.S., Venugopal, S., Broberg, J., Brandic, I. (2009) "Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility". In: Future Generation Computer Systems, Elsevier B. V.

[4] Adams K. and Agesen O. A Comparison of Software and Hardware Techniques for x86 Virtualization. *ASPLOS* October 2006.

[5] AMD. (2005) Amd secure virtual machine architecture reference manual.

[6] Barham, P., B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield, Xen and the art of virtualization, Proceedings of the Nineteenth ACM Symposium on Operating systems Principles. ACM Press, New York, 2003, pp. 164–177.

[7] Hostway UK VMware ESXi Cloud Simplified, Comprehensive explanation of the features and benefits of VMware ESXi Hypervisor.

[8] Fujitsu Technology Solutions, DataSheet Citrix XenServer,

[9] Che, J., Q. He, Q. Gao, D. Huang, ―Performance Measuring and Comparing of Virtual Machine Monitors,College of Computer Science, Zhejiang University, Hangzhou 310027, China, IEEE/IFIP International Conference on Embedded and Ubiquitous Computing, 2008.

[10] FUJITSU, Benchmark Overview-vServCon, white paper, March 2010.

[11] Apparao, P. & Makineni, S. & Newell, D.Virtualization (2006) Characterization of network processing overheads in Xen. Technology in Distributed Computing, 2006. VTDC 2006.

[12] Jianhua, C. & Qinming, H. & Qinghua, G. & Dawei, H. (2008) Performance Measuring and Comparing of Virtual Machine Monitors. Embedded and Ubiquitous Computing, 2008. EUC '08.

[13] Menon, A. et Al. (2005) Diagnosing Performance Overheads in the Xen Virtual Machine Environment. Conference on Virtual Execution Environments (VEE'05).

[14] Shan, Z. & Qinfen, H. (2009) Network I/O Path Analysis in the Kernel-based Virtual Machine Environment through Tracing. Information Science and Engineering (ICISE).

[15] VMware (2007) A Performance Comparison of Hypervisors VMware. White paper feb 1, 2007.

[16] Passmark. Performance Test – PC Benchmarking

[17] VMware (2007) Understanding Full Virtualization, Paravirtualization, and Hardware Assist. VMware, white paper nov 10, 2007.

[18] Vallee, G. & Naughton, T. & Engelmann, C. & Ong, H. &Scott, S.L. (2008) System- Level Virtualization for High Performance Computing. Parallel, Distributed and Network-Based Processing, 2008. PDP 2008. 16th Euromicro Conference on. Varia, J. (2008) Cloud Architectures. White Paper of Amazon.

[19] VMware, ―The Architecture of VMware ESXi, white paper, 2007.

[20] Xu, X., F. Zhou, J. W. Y. Jiang, ―Quantifying Performance Properties of Virtual Machines, School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou 310018, China, International Symposium on Information Science and Engineering, 2008.

[21] Xen,―How does Xen work, Xen Orgnaization

[22] CloudStack – OpenSource Cloud Computing

[23] Greenberg, A & Hamilton, J & Maltz, D. A. & Patel, P. (2009) The Cost of a Cloud: Research Problems in Data Center Networks.

[24] He, Q. & Zhou, S. & Kobler, B. & Duffy, D. & McGlynn, T. (2010) Case study for running HPC applications in public clouds. HPDC '10 Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing.

[25] Karger, P.A. & Safford, D.R. (2008) I/O for Virtual Machine Monitors: Security and Performance Issues. Security & Privacy, IEEE.

[26] King, R. (2008) How Cloud Computing Is Changing the World. CEO guide to technology.

[27] Kloster, J. F. & Kristensen, J. & Mejlholm, A. (2007) A Comparison of Hardware Virtual Machines Versus Native Performance in Xen.

[28] Lui, J. & Huang, W. & Abali, B. & Panda, K. D. (2006) High performance VMM Bypass I/O in virtual machines. USENIX 2006 Annual technical conference refereed paper.

[29] Mollick, E. (2006) Establishing Moore's Law. Annals of the History of Computing, IEEE.

[30] XenSource (2007) A Performance Comparison of Commercial Hypervisors. XenEnterprise vs. ESX Benchmark Results. 2007 XenSource.

[31] Ubuntu 12.04 – Free Operating System

[32] Padala, P. & Zhu, X. & Wang, Z. & Singhal, S. & Shin, K. G. (2007). Performance Evaluation of Virtualization Technologies for Server Consolidation. Enterprise Systems and Software Laboratory, HP Laboratories Palo Alto.

[33] Windows 2008 R2 Server Operating sytesm

[34] Zhao, T. & Ding, Y. & March, V. & Dong, S & See, S. (2009) Research on the Performance of xVM Virtual Machine Based on HPCC. ChinaGrid Annual Conference, 2009. ChinaGrid '09. Fourth.

[35] Network Performance Benchmark Netperf.

[36] Hyperic's System Information Gatherer (SIGAR)

## AUTHOR'S PROFILE

Vijaya Vardhan Reddy (First Author): He received his M.Tech degree in Computer Science and Engineering from Osmania University in 2000. For the past 14 years he has been working in the IT industry. He had worked in Tokyo (2001) for CSFB Project and in London (2005-06) for ADP Freedom Payroll Project in Java / J2EE Technologies. Currently he is working as an Assistant Vice President in GE Capital. His research areas include distributed computing, grid computing and cloud computing.

Dr. Lakshmi Rajamani (Second Author): She is a retired professor from Osmania University. She has many papers published in journals across the world. She had served as a HOD for computer science department from 2010 to 2012.

# TCP- Costco Reno: New Variant by Improving Bandwidth Estimation to adapt over MANETs

Prakash B. Khelage

Asst. Professor, Information Technology,
UMIT, SNDT Women's University, Mumbai-400049,
India.

Dr. Uttam D. Kolekar

Principal Smt. Indira Gandhi College of Engineering,
Navi Mumbai - 400709,
India.

*Abstract*—The Transmission Control Protocol (TCP) is traditional, dominant and has been de facto standard protocol, used as transport agent at transport layer of TCP/IP protocol suite. Basically it is designed to provide reliability and assure guaranty to end-to-end delivery of data over unreliable networks. In practice, most TCP deployments have been carefully designed in the context of wired networks. Ignoring the properties of wireless Ad Hoc Networks, therefore it can lead to TCP implementations with poor performance. The problem of TCP and all its existing variations within MANETs resides in its inability to distinguish between different data packet loss causes, whenever the data loss occur traditional TCP congestion control algorithm assumes loss is due to congestion episode and reduces sending parameters value unnecessary. Thus, TCP has not always the optimum behavior in front of packet losses which might cause network performance degradation and resources waste. In order to adapt TCP over mobile Ad hoc environment, improvements have been proposed based on RTT and BW estimation technique in the literature to help TCP to differentiate accurate causes between the different types of losses. But still does not handle all the problems accurately and effectively. In this paper, a proposed TCP-Costco Reno a New Variant, accurately estimates the available bandwidth over Mobile Ad Hoc networks and sets sending rate accordingly to maximize utilization of available resources and hence improves performance of TCP over mobile Ad hoc networks. The results of the simulation indicate an improvement in throughput over interference, link failure and signal loss validation scenarios. Further, it shows highest average of average throughput then those variants which are most successful over MANETs.

*Keywords—mobile ad hoc network (MANET); Ccongestionl; Link failure;signal loss; interference; Retransmission timeout; RTT and BW estimation*

## I. INTRODUCTION

The phenomenal growth experienced by the Internet over the last decade has been supported by a wide variety of evolving mechanisms to meet the requirements of emerging, demanding applications. The basic TCP/IP protocol suite has been instrumental in developing today's Internet. In particular, TCP has been successful due to its robustness in reacting dynamically to changing network traffic conditions and providing reliability on an end-to-end basis. This Wide acceptance has driven the development of many TCP applications, motivating the extension of this protocol to wireless networks. These networks pose some critical challenges to TCP since it was not originally designed to work in such complex environments, where the level of bit error rate (BER) is not negligible due to the physical medium[3][9].

High mobility may further degrade the end-to-end performance because TCP reduces its transmission rate whenever it perceives a dropped packet.

Mobile ad hoc network is a collection of mobile nodes that offers different opportunities to TCP. Reduction in deployment cost due to absence of fixed infrastructure and elimination of administration cost since it is self-configurable.

However, MANET consists of unstable wireless communication links in compare to the wired network [1]. This instability is mainly due to mobility of nodes. Because TCP is originally invented for wired network [3], it ignores non-congestion loss which occurs rarely in this environment. Thus, TCP in present form cannot address frequent link breakage in MANET and suffers from performance degradation [4]. TCP is responsible for providing reliability of connection by retransmitting lost packet. Congestion control is the most controversial parts of TCP which degrades performance in front of packet loss [9]. Congestion control as its name appears, assumes all packet loss induced by congestion.

When link failure lasts greater than RTO (Retransmission timeout), Retransmission timer expires and TCP interprets packet loss as a congestion loss. Then congestion control executes back-off algorithm to grow RTO exponentially and retransmit packet. After a few successive back-off executions, RTO becomes too long. Hence when route recovered, sender resumes data transmission with long RTO which forces sender remains idle unnecessary in case of probable next losses[2]. Thus, traditional TCP fails in wireless network. As per the literature study, the proposals based on RTT (retransmission time) and BW (Bandwidth) estimation technique are successful till somewhat extent and not utilizing available resources optimized. Therefore further improvement or modification need to be done. This paper focused on the modification of TCP-Westwood which is sender side modification and based on BW (bandwidth) estimation Technique.

The remainder of the paper is organized as follows. In Section II, literature survey of TCP-Westwood and TCP-WELCOME with its algorithms, problems and comparative analysis is presented. Section III defines the problems with existing Solution. Section IV explains about the proposal for solution. Tools and Techniques, Basic validation Scenario with discussion and experimental results presented in section V. Section VI Acknowledges to those who motivated and

helped me lot for doing this research work. Section VII concludes the paper.

## II. Literature Survey

### A. TCP Westwood

TCP Westwood proposes an end-to-end bandwidth estimation algorithm based on TCP Reno. TCP Westwood implements slow start and congestion avoidance phases as TCP Reno, but instead of halving the congestion window size as in TCP Reno when congestion happens, TCP Westwood adaptively estimates the available bandwidth and sets the congestion window size and slow start threshold accordingly to improve the link utilization. In TCP Westwood, packet loss is indicated by the reception of 3 duplicated acknowledgements (DUPACKs) or timeout expiration. When 3 DUPACKs are received, TCP Westwood sets SSThreshHold and CWND as follows:

*if (3 DUPACKs are received)*

*SSThreshHold = (BE \* RTTmin)/seg_size;*

*if (CWND > SSThreshHold) /\* in congestion avoidance phase\*/*

*CWND = SSThreshHold;*

*endif*

*endif*

Where the seg_size is the length of the TCP segments and RTTmin is the minimum RTT experienced. BE is the estimated available bandwidth. It is assumed in TCP Westwood that when 3 DUPACKs are received in the congestion avoidance phase, the available bandwidth is fully utilized. So the values SSThreshHold and CWND should reflect the estimated bandwidth (BE). If a packet loss is indicated by timeout expiration, TCP Westwood sets SSThreshHold and CWND as follows:

*If (timeout expires)*

*CWND = 1;*

*SSThreshHold = (BE \* RTTmin)/seg_size;*

*if (SSThreshHold < 2)*

*SSThreshHold = 2;*

*endif*

*endif*

This sets the CWND to 1 and SSThreshHold to BE after the timeout event and then the TCP Reno behavior continues. In TCP Westwood, the setting of SSThreshHold and CWND is based on the bandwidth estimation, which is obtained by measuring the rate of the acknowledgments and collecting the information of the amount of packets delivered to the receiver in the ACK. Samples of bandwidth are computed as the amount of packet delivered divided by the inter-arrival time between two ACKs. Those sample bandwidth estimates are then filtered to achieve an accurate and fair estimation.

TCP Westwood modifies the Additive Increase and Multiplicative Decrease (AIMD) in TCP Reno and adaptively sets the transmission rates to remove the oscillatory behavior of TCP Reno and to maximize the link utilizations. But this also causes TCP Westwood to degrade the performance of TCP Reno connections when they coexist in the network [11].

**Problems:**

The behavior of the bandwidth estimation scheme is unpredictable, therefore TCP Westwood Perform poorly if it estimates incorrect bandwidth. Changes in the inter-arrival times of the acknowledgements cause improvement or worsening of the throughput in rather unpredictable ways. Additionally, the sensitivity of TCP Westwood Acknowledged Interval is variable.

### B. TCP WELCOME

It is a sender-based solution, known as TCP-WELCOME, to improve the TCP performance for route failure, wireless error and congestion losses in MANET based on RTT. TCP-WELCOME distinguishes between causes of packet loss and then triggers the most appropriate packet loss recovery according to the identified loss cause. It realizes its loss differentiation by observing the history of RTT samples evolution over the connection and the data packet loss triggers (3DuplACK and RTO). If loss is detected by 3DuplACK and RTT values are stable then it is wireless related packet loss or else it is congestion. On the other hand, if loss is detected by RTO and RTT values are stable then it is route failure related loss otherwise, it is congestion [4].



Fig. 1. TCP WELCOME Loss different algorithm [4]

Fig. 4 summarizes the main idea of loss differentiation Algorithm. After identifying the cause of a data packet loss using the proposed LDA, TCP-WELCOME react will concern on RTO calculation and data transmission rate. So in case of congestion loss no change is required to the standard TCP New Reno, the same goes for wireless error no change only retransmit the lost packet, however, once route failure is detected the RTO value and CWND would be updated based on the new route characteristics( length, load and quality) as follows:

*RTO new = (RTTnew / RTTold) × RTOold*

*CWND new = (RTTold/RTTnew) × CWNDold*

**Problems:**

Does not take network disconnection and frequent route change affecting into account during the evaluation, WELCOME uses RTT which include both delays of forward and reverse path while only delay of forward path must be

considered. In addition, it offers recovery method based on RTT Comparison TCP-Welcome claimed that RTO adjustment should be done based on the capabilities of discovered route such as length; load and link quality [6] after link breakage, total delay for new route varies from broken route.

TABLE I.        COMPARATIVE ANALYSIS OF TCP- VARIANTS

| TCP-Variant ⇨ Comparative Parameters ⬇ | TCP- Westwood | TCP- Welcome |
|---|---|---|
| Enhancement Proposed | End-to-End bandwidth estimation by Monitoring rate of returning ACKs. | End-to-End, implicit, loss differentiation and loss recovery algorithm solution. |
| Advantages | 1. Utilizes available bandwidth efficiently. 2. Handles delayed and cumulative ACKs and random loss problems solved. | 1. WELCOME outperforms then New Reno, SACK, vegas, Westwood in terms of average throughput and energy consumption. |
| Limitations | 1. Cannot distinguish between buffer overflow and random losses. 2. Performs poorly if it estimates incorrect Bandwidth. | 1. Does not take network disconnection and frequent route change affecting into account during the evaluation. 2. WELCOME uses RTT which include both delays. |
| Throughput | 550% over TCP Reno | Improved than New Reno, SACK, Vegas, and Westwood in interference and link failure scenario. |
| Handling BER | Westwood does not handle Bit Error Rate. | WELCOME does handle Bit Error Rate. |
| Energy Efficiency | Low | Good |
| Timeout  Measurements | No timeout measurements used. | Course-grained timeout used. |

III.    PROBLEM DEFINATION

In ad hoc networks, the principal problem of TCP lies in performing congestion control in case of losses that are not induced by network congestion. Since bit error rates are very low in wired networks, nearly all TCP Variants assume that packets losses are due to congestion. Consequently, when a packet is detected to be lost, either by timeout or by multiple duplicated ACKs, TCP slows down the sending rate by adjusting its congestion window. Unfortunately, wireless

networks suffer from several types of losses that are not related to congestion, making TCP not adapted to this environment. With the help of Literature survey of TCP over Ad Hoc Networks, identified following problems. In Wireless ad-hoc networks nodes may be mobile therefore no predefined topology.

As nodes can join and leaves network, so accordingly topology may change. When the topology of the network changes every time, then routing mechanism needs to trigger

to find alternative roots to do the reliable end to end communication between sender and receiver.

Thus due to frequently changes in topology communication links may failures and there will be the loss of data segment. To recover segment loss TCP sender reduces sending rate by triggering congestion control mechanism which sets size of congestion window of its lowest value. Assuming that the loss of packet is due to congestion in network which is totally misjudged and there will be the underutilization of available bandwidth. Hence the performance of TCP degrades. These problems are due to lossy channels, Hidden and exposed stations, path asymmetry which may appear in several forms like BW asymmetry, loss rate asymmetry and route asymmetry.

A large number of approaches for RTT and BW estimation proposed but none of them work well in all scenarios without any drawback or side effect. So it is essential that an improved BW estimation technique is evolved. The focus of our research is to study the existing techniques, propose an improved BW estimation technique based on TCP-Westwood which will have to reduce under and over estimation of available bandwidth before transmission.

## IV. PROPOSAL FOR SOLUTION

### A. Introduction

TCP variants are well adapted to deal with all data packet loss situations that can be encountered within wireless ad hoc networks. The performance of TCP degrades significantly within wireless ad hoc networks [8]. Moreover, some of the studied TCP variants perform well in certain cases while they perform badly in other cases. The ability of TCP to distinguish among congestion-induced and wireless-related data losses (as in TCP Westwood) leads to an improved performance in some cases. However, TCP variants that incorporate a loss differentiation algorithm do not consider all types of data packet loss that can be encountered within wireless ad hoc network environments. In fact, they consider congestion-induced and wireless-channel related losses only. It also finds that the TCP variant which is able to adjust its performance parameters (CWND and RTO) after data losses (as in TCP Vegas) [5] [6], in certain cases, can improve the performance within the network.

In this paper, we propose TCP-Costco Reno a new TCP variant that is based on the modification of Westwood, which will be able to deal effectively with the under as well as over estimation  of available bandwidth.

### B. Motivation

The main motivation behind this research work is TCP Westwood and WELCOME Variants. Which are most successful in wireless Ad hoc network but TCP Westwood does wrong estimations, like over or under estimation of the available bandwidth due to variable delay in returning ACK Whereas, TCP WELCOME does not handle network disconnection and frequent route change. The research work has been concentrated on bandwidth estimation based technique.

### C. Problems with existing solution

TCP Westwood is designed to perform well in wired, wireless and mixed networks. In TCP Westwood the TCP Reno congestion control is used and is modified only on the sender side [10]. TCP Westwood defines that if the congestion window size divided by the minimum RTT is larger than the currently achieved rate the channel is congested and if it is equal to the current rate, the loss is of random nature. The key innovative idea behind TCP Westwood is that it takes advantage of an end-to-end bandwidth estimation mechanism called Bandwidth Share Estimates (BSE) to set the values of slow start threshold (SSThreshHold) and congestion window (CWND) after a random (due to the lossy nature of a wireless link) loss (indicated by 3 duplicated acknowledgements that have been received or by a specific timeout).

Instead of setting the slow start threshold after receiving 3 duplicated ACKs or after the timeout expires to half of the size of the congestion window (as in TCP Reno) TCP Westwood sets the SSThreshHold to the product of the BSE and the minimum RTT. This estimation is not based on measurements performed by lower layers and therefore retains to layer principles like separation and modularity of layers. The bandwidth is estimated by the source that monitors continuously the received TCP acknowledgements. It then estimates the data rate currently achieved by the connection. By doing this, TCP Westwood ensures both faster recovery and more effective congestion avoidance. However TCP-Westwood cannot distinguish between buffer overflow and random losses, Performs poorly if it estimates incorrect Bandwidth and it is not sufficiently evaluated. Therefore the mechanism of TCP Westwood modified and named as TCP Costco Reno, as this variant also retains the principle of New Reno.

### D. TCP-Costco Reno (New Variant for MANET)

As per the pseudo code of TCP Westwood, it has calculated the value of SSThreshHold and compared it with CWND and then according to it sets the value of CWND. It is noticed that, it should not forcefully increase CWND if it is smaller than SSThreshHold, as well as this mechanism can't set the value properly, Hence Bandwidth estimation is inaccurate in case of TCP-Westwood, which faces problem like over and under estimation of available bandwidth over wireless networks scenarios.

**Over bandwidth estimation:** to set bandwidth more than available bandwidth.

**Under bandwidth estimation:** to set bandwidth less than available bandwidth.

So, to overcome above problems the proposal estimates bandwidth properly by using following code:

*Set SSThreshHold = (min 4 or maximum calculated value)*

*Set slow start = 2*

*thresh_= (int)((current_bwe_/size_/8)\* in_rtt_estimate);*

*if (thresh_ > 4 )*

*SSThreshHold_ = thresh_*

*else*

*if (thresh_ < 4)*

*SSThreshHold_ = 4*

*if (CWND_ > SSThreshHold_)*

*CWND_ = SSThreshHold_*

Thus, TCP Costco Reno sets optimum value of bandwidth and hence it recovers problems of under and over estimation.

In wireless environment whenever packet lost happens traditional TCP assumes that, it is due to congestion, but packet lost in wireless environment may happens due to congestion, link failure or wireless channel errors. All the times packet loss is may not due to congestion, so TCP Variant has to consider link failure and wireless channel errors also. Every time when it considers packet lost is due to congestion then TCP Westwood sets slowstart as "1" but in case of other Scenario like link failure and wireless channels errors no need to set slowstart as 1 all the time. So this is the reason that, in TCP Costco Reno sets slowstart equal to 2.

### E. Pseudo code of TCP-Costco Reno

*A. Algorithm after 3 duplicate ACKS:*

*If(3 dupack received)*

*thresh_ = (int)((current_bwe_/size_/8)* in_rtt_estimate);*

*if (thresh_ > 4*

*SSThreshHold_ = thresh_*

*else*

*if (thresh_ < 4)*

*SSThreshHold_ = 4*

*if (CWND_ > SSThreshHold_)*

*CWND_ = SSThreshHold_*

*B. Algorithm after slowstart*

*if (CWND_ < SSThreshHold_)*

*slowstart = 2*

### V. Tools, Validation Model And Simulation Environment

NS-2 is a discrete event simulator written in C++, with an OTcl interpreter shell as the user interface that allows the input model files (Tcl scripts) to be executed. Most network elements in NS-2 are developed as classes, in object-oriented fashion [7]. NS2 provides substantial support for simulation of TCP, routing algorithms, queuing algorithms, and multicast protocols over wired and wireless (local and satellite) networks, etc.

It is freely distributed. So, in order to investigate and understand the behavior of the congestion scheme and to observe the improved performance of our approaches in our research work we are using NS-2 as network simulator tool. We have planned to compare it with different TCP variants using different simulation scenarios with different proactive and reactive routing protocols that will describe multiple data packet loss causes which are related to wireless ad hoc networks.

We have used NS2 as a network simulation tool for hypothesis testing and Study the effect of the different loss scenarios (link failure, congestion, signal loss and interference) Evaluation of TCP-Costco Reno and comparisons with other TCP variants such as Tahoe, Reno, New-Reno, Vegas, Sack and Westwood.

In our simulation experiments we used proactive and reactive routing protocol and investigate which protocol is suitable for wireless ad hoc network to adapting TCP over ad hoc network. Nodes communicate through identical wireless radio settings using the standard MAC 802.11.

TABLE II. Simulation Parameters Values

| Parameter | Values |
|---|---|
| Channel Type | Wireless channel |
| Radio Propagation Model | Two Ray ground |
| Queue type | Droptail/PriQue |
| Max. packet(buffer size) | 50 |
| Network interface | Wirelessphy |
| MAC Protocol | 802.11 |
| Data Rate | 1 Mbps |
| Transmission Radius | 250 |
| Interference Radius | 550 |
| Packet size | 1000 bytes |
| Routing protocol | AODV, DSDV |
| Simulation Time | 150 s |
| Value x | 700 |
| Value y | 500 |
| Agent trace | ON |
| Mac trace | OFF |
| Router trace | ON |
| Movement trace | ON |

## A. *Basic Validation scenarios*



Fig. 2.   basic Validation Scenario

The basic validation scenarios, using NS-2, are implemented as follows:

**1) Congestion validation scenario:** In this scenario, congestion created at the middle of a five node topology by generating three TCP data traffic flows that must pass by intermediate node to reach the other communicating end point. Different levels of data congestion generated by controlling the number of TCP data flows crossing intermediate node at a certain time. Congestion scenario created sending multiple TCP sources through a bottleneck link. In congestion also Costco Reno provides best average throughput than new Reno, Vegas, sack and nearly equal to Westwood.



Fig. 3.   Throughput of Costco Reno in congestion



Fig. 4.   average throughput of all TCP variants in     congestion s scenario

**2) Interference Validation Scenario:** In this case, two TCP connections are established in parallel. The main TCP connection (TCP data flow 1 in Figure 6) is disturbed by the interferences generated by the second TCP connection. Indeed, the node acting as forwarder for the main TCP connection is placed within the interference range of the second TCP connection sender. So, this situation creates interference and thus data packet losses. In this scenario also TCP Costco Reno gives better average throughput than Westwood and all other TCP variants except sack by 0.15.



Fig. 5.   Average throughput of all TCP variants in Interference scenario

**3) Link Failure Validation Scenario:** In link failure validation model it has been forced to TCP traffic to change its communication path by shutting down the intermediate node between the communicating ends nodes. In addition, it employs routes with different number of hops. Thus, each time TCP changes the communication route, the characteristics of the path between the communicating nodes changes. It is obvious that the choice and the establishment delay of the new route will be dependent on the implemented ad hoc routing protocol. Packet losses and delay changes will also be generated by the link loss and the new chosen route.

Link failure scenario is created by five nodes, source and destination provided with mobility. Due to mobility both sender and receiver will be out of transmission and reception range link goes down for few seconds and few packets will drop. Other node join the mobile ad hoc network, provides transmission and reception range to sender and receiver node again transmission starts through new link. Referring with fig. 6 in this scenario also TCP-Costco Reno gives better performance than all other TCP-variants.



Fig. 6.    Average throughput of all TCP variants in link failure scenario

**4) Signal loss Validation scenario:** This scenario illustrates the situation where the wireless signal is not stable. The communicating nodes loose the connection due to signal loss then they resume the communication when the signal comes back. Signal losses are generated by moving one of the intermediate nodes out of the radio range of its connection neighbours. This scenario created using three nodes. End nodes acts as sender and receiver and intermediate node as router Transmission of ftp traffic source flow through intermediate node. Intermediate node moves away for few second so signal loss occurs between source and destination, after few second intermediate node moves at original place and again retransmission starts. In signal loss scenario TCP-Costco Reno performs better than all other TCP variant as shown in figure: 7



Fig. 7.    Average throughput of all TCP variants in signal loss scenario

*B. Main Simulation scenarios Chain multi-hop network:-*

The network consists of variable length chain of static nodes, placed at a distance of 200m from one another. FTP traffic is transferred between the first and last node of the chain shown network.



Fig. 8.    Chain Topology

During the simulation one FTP connection kept active at a time. Sequential TCP connection are initiated and terminated. The TCP-Costco Reno used at the transport layer which is an improvement of TCP Westwood. For the simulation of chain topology three hop networks created consisting of four nodes whereas, end nodes acts as sender and receiver. In this case TCP Costco Reno performs same as Tahoe as well as Tahoe based other variant but very much better then Vegas.



Fig. 9.  Average throughput of all TCP variants in      Chain Topology

### Grid network:-

Fig. 10 shows a static grid network as experiment topology with *3X4* nodes. The distance between two adjacent nodes is set to be 150 m, and the transmission and interference radii are set to 250 and 550 m, respectively.



Fig. 10.  Grid Topology

In each row, a TCP connection is assumed to set up from the left end node to the right end, and similarly, in each column, a TCP connection is assumed to set up from the bottom end node to the top end node. The histogram of average throughput of grid network shows that, TCP Costco Reno gives good average throughput then all TCP -variant except Reno.



Fig. 11. Throughput plots of all TCP variant in grid networks



Fig. 12.  Average Throughput of TCP Variants in Grid Networks

TABLE III.　　AVERAGE VALUE OF AVERAGE THROUGHPUT IN ALL SCENARIO OF ALL TCP-VARIANT

| Average of average Throughput ⬇ | Network Scenarios ⇨  TCP Variants ⬇ | Congestion Scenario | Interference Scenario | Link Failure Scenario | Signal Loss | Chain Network Scenario | Grid Network Scenario |
|---|---|---|---|---|---|---|---|
| C.　429.94 | Tahoe | D.　524.50 | E.　678.7 | F.　497.07 | G.　156.93 | H.　227.68 | I.　494.80 |
| J.　440.855 | Reno | K.　525.63 | L.　679.4 | M.　494.57 | N.　156.24 | O.　227.68 | P.　564.61 |
| Q.　434.271 | New Reno | R.　508.07 | S.　679.62 | T.　497.03 | U.　156.96 | *V.*　227.68 | W.　536.27 |
| X.　433.276 | Sack | Y.　512.86 | Z.　679.77 | AA.　496.89 | BB.　146.90 | *CC.*　227.68 | DD.　535.56 |
| EE.　217.138 | Vegas | FF.　271.96 | GG.　348.5 | HH.　251.84 | II.　84.65 | *JJ.* 116.44 | KK.　229.44 |
| LL. 427.348 | Westwood | MM.　540.66 | NN.　678.88 | OO.　497.07 | PP.　156.19 | *QQ.*　227.68 | RR.　463.61 |
| SS.　437.501 | Costco Reno | TT. 515.48 | UU.　679.62 | VV.　497.18 | WW.　161.04 | *XX.*　227.68 | YY. 544.01 |



Fig. 13. Average of average Throughput of TCP Variants in all Network Scenario's

## VI. Acknowledgment

## VII. Conclusion

This paper presents a thorough literature survey of TCP-Westwood and TCP-WELCOME along with its comparative analysis study and problems encountered in MANETs. As in case of wireless networks, performance of TCP degrades because of its inability to handle wireless channel errors. The proposed article placed special emphasis on those TCP-Variants that preserve end-to-end semantic and most successful over mobile ad Hoc Network. This article Proposed Bandwidth estimation based TCP-Costco Reno, a new variant which is the modification of TCP-Westwood. With reference to the experimental results and data analysis, it shows that, TCP-Costco Reno handles much efficiently the wireless channel errors (signal loss and interference). Hence it improves overall throughput without degrading TCP performance in other scenarios and suitable for mobile ad hoc networks to avoid over or under estimation of BWE value. It outperforms in different types of wireless channel errors such as signal loss, link failure and interference and gives improved performance in grid network. Simulated result also indicates that, Average of average throughput of TCP Costco Reno in all scenarios is better than all other variants.

### References

[1] J. Li, C. Blake, D.S.J. De Couto, H.I. Lee, and R. Morris, "Capacity of Ad Hoc Wireless Networks," Proc. ACM MobiCom '01, July 2001.

[2] Z. Fu, P. Zerfos, H. Luo, S. Lu, L. Zhang, and M. Gerla, "The Impact of Multi hop Wireless Channel on TCP Throughput and Loss," Proc. INFOCOM '03, Apr. 2003.

[3] K. Chen, Y. Xue, and K. Nahrstedt, "On Setting TCP's Congestion Window Limit in Mobile Ad Hoc Networks," Proc. IEEE Int'l Conf. Comm. (ICC '03), May 2003.

[4] A. Seddik-Ghaleb, Y. Ghamri-Doudane, and S. M. Senouci, "TCP WELCOME TCP Variant for Wireless Environment, Link losses, and COngestion packet loss ModEls," in First International Communication Systems and Networks and Workshops, COMSNETS 2009.

[5] Bhaskar Sardar, Debashis Saha, "A Survey of TCP Enhancements for Last-Hop Wireless Networks" in IEEE Communications Surveys & Tutorials, 2006 3rd Quarter 2006,Volume 8, No 3.

[6] Haifa Touati, Ilhem Lengliz, Farouk Kamoun, "TCP Adaptive RTO to Improve TCP performance in mobile ad hoc networks," in The Sixth Annual Mediterranean Ad Hoc Networking WorkShop, Corfu, Greece, June 12-15, 2007.

[7] Network Simulator Ns2, http://www.isi.edu/nsnam/ns.

[8] H. M. El-Sayed, "Performance Evaluation of Tcp in Mobile Ad-Hoc Networks," in the Second International Conference on Innovations in Information Technology (IIT) 2005.

[9] Adib M.Monzer Habbal, "Loss Detection and Recovery Techniques for TCP in Mobile Ad Hoc Network" in Second International Conference on Network Applications, Protocols and Services, 2010.

[10] Mario Gerla, M. Y. Sanadidi, Ren Wang, and Andrea Zanella, "TCP Westwood: Congestion Window Control Using Bandwidth Estimation," UCLA Computer Science Department, 0-7803-7206-9/ 2001

[11] Ahmad Al Hanbali, Eitan Altman, And Philippe Nain, Inria Sophia Antipolis France," A SURVEY OF TCP OVER AD HOC NETWORKS", IEEE Communications Surveys & Tutorials ,Third Quarter 2005.

### Authors Profile

**Prakash B. Khelage** received his B.E. in Electronics and Telecommunication Engineering from Dr. Babasaheb Ambedkar Marathwada University Aurangabad, M.Tech in information Technology from NMIMS University Mumbai, Maharastra, India. He is currently working as Assistant Professor with UMIT, SNDT Women's University. He has 13 years of experience in industrial as well as educational field; His research interest includes Ad Hoc Networks, Mobile Computing, Wireless Networks, Co-operative Communication Networks and Network Security. He has also interest in Computer Architecture design, Cloud Computing and Data Mining.

**Uttam D. Kolekar** received his B.E. in Electronics and Telecommunication Engineering, M.E. in Electronics from Shivaji University, Kolhapur and Awarded Ph. D. in electronics from Bharati Vidhyapith Pune, Maharastra, India. He is currently working as Principal with Smt. Indira Gandhi College of Engineering, Mumbai University. He has more than 20 years of experience in educational institution; His research interest includes Ad Hoc Networks, Mobile Computing, Wireless Networks, Neural Network and Co-operative Communication Networks. He has published over 30 National and International Jurnals & conferences various papers accros India and other countries.

# Early Development of UVM based Verification Environment of Image Signal Processing Designs using TLM Reference Model of RTL

Abhishek Jain
[1]Imaging Group and [2]Jaypee Business School
[1]STMicroelectronics and [2]Jaypee Institute of Information
Technology Greater Noida, India

Dr. Hima Gupta
Jaypee Business School
Jaypee Institute of Information Technology
Noida, India

Sandeep Jana
SDS Group
STMicroelectronics Greater Noida,
India

Krishna Kumar
SDS Group
STMicroelectronics Greater Noida,
India

*Abstract*—With semiconductor industry trend of "smaller the better", from an idea to a final product, more innovation on product portfolio and yet remaining competitive and profitable are few criteria which are culminating into pressure and need for more and more innovation for CAD flow, process management and project execution cycle. Project schedules are very tight and to achieve first silicon success is key for projects. This necessitates quicker verification with better coverage matrix. Quicker Verification requires early development of the verification environment with wider test vectors without waiting for RTL to be available.

In this paper, we are presenting a novel approach of early development of reusable multi-language verification flow, by addressing four major activities of verification –

1. **Early creation of Executable Specification**

2. **Early creation of Verification Environment**

3. **Early development of test vectors and**

4. **Better and increased Re-use of blocks**

Although this paper focuses on early development of UVM based Verification Environment of Image Signal Processing designs using TLM Reference Model of RTL, same concept can be extended for non-image signal processing designs.

*Keywords—SystemVerilog; SystemC; Transaction Level Modeling; Universal Verification Methodology (UVM); Processor model; Universal Verification Component (UVC); Reference Model*

## I. INTRODUCTION

Image signal processors (ISP) address different markets, including high-end smartphones, security/surveillance, gaming, automotive and medical applications. The use of industry standard interfaces and rich set of APIs makes the integration



Fig. 1. Verification Environment of Image Signal Processing Design

of image processors a straightforward process and helps to reduce end-product time to market.

Image signal processing algorithms are developed and evaluated using C/Python models before RTL implementation. Once the algorithm is finalized, C/Python models are used as a golden reference model for the IP development. To maximize re-use of design effort, the common bus protocols are defined for internal register and data transfers.

A combination of such configurable image signal processing IP modules are integrated together to satisfy a wide range of complex image signal processing SoCs [1].

In Verification Environment of Image Signal Processing design as shown in figure 1, Host interface path is used to do programming of configurable blocks using SystemVerilog UVM based test cases. UVM_REG register and memory model [20] is used to model registers and memories of DUT. DUT registers are written/read via control bus (AXI3 Bus here) UVC. RTL control bus interface acts as target and control bus UVC acts as initiator. The target control interface of the ISP RTL is driven by control bus UVC (configured as initiator).After register programming is done, image data(random/user-defined) is driven to the data bus interface by the data bus UVC and the same data is also driven to the reference model. Output of the ISP RTL is received by the receiver/monitor of the data bus UVC. Scoreboard compares the output of RTL and reference model and gives the status saying whether the both output matches or not.

'C' test cases are used for programming of RTL registers/memories via CPU interface. C test cases control the SystemVerilog Data Bus UVC using Virtual Register Interface (VRI) [15], [18]. VRI layer is a virtual layer over verification components to make it controllable from embedded software. It gives flexibility to Verification Environment users to use the Verification IPs without knowing SystemVerilog.

Generally, development of Verification Environment for verification of designs is started after availability of the RTL. Thus, significant time is spent for setup and debugging of verification environment after release of RTL which results in delay in start and completion of verification of the designs. It is required to find ways to start developing the Verification Environment much before the arrival of the RTL so that when RTL is available, Verification Environment can be easily plug and play and verification of the designs can be started quickly. Use of TLM reference model of RTL for development of Verification Environment much before arrival of RTL proves to be good solution for the above mentioned problem.

This paper is focusing on early development of UVM based Verification Environment of Image Signal Processing designs using TLM Reference Model of RTL before availability of the RTL. Early development of Verification Environment of Image Signal Processing designs is described in detail in Section II.

## II. EARLY DEVELOPMENT OF UVM BASED VERIFICATION ENVIRONMENT

### A. Modeling of ISP designs

A loosely timed high level model of the ISP block is generated at algorithmic functional level using C/C++/SystemC and with TLM-2 interface.

The SCML – SystemC Modeling Library, an open source SystemC library from Synopsys Inc. [26] is being used here.

The purpose of this model generation is to use this as a reference model. We may say it as a "Golden Reference Model" or "Executable Functional Specification" of the ISP

designs. From functional and structural perspective this model can be divided in two major spaces.

**First space -** the algorithmic computational part, is mainly responsible for image processing using various algorithms involved for image manipulation from the incoming image stream data.

**The second space** – a TLM interface, is responsible for all kinds of communication to external IPs and other system blocks.

Register interface of this model is generated using IP-XACT tools. And algorithmic part is manually implemented.

### B. Testing of Executable Spec only

To test the TLM ISP model, an environment is developed using Python (an open source scripting language) and Synopsys Pa-Virtualizer Tool Chain.

The test environment has following major components:

- Test bench in Python
- Configuration file reader in Python
- Raw Data Reader
- ISP model
- Input data injector in Python
- Output data receiver in Python
- Output data checker in Python
- Synopsys Pa-Virtualizer Tools Chain for GUI, debugging, and simulation

XML file format is used for test bench configuration and passing other parameter to the testing environment.



Fig. 2. ISP Model Testing Environment

## C. Use of TLM ISP Model for early development of RTL Verification Environment

After the ISP model is proved to be functionally correct, the same model is used for early development of RTL functional verification environment.

A suitable TLM sub-system is designed. This TLM sub-system consists of various models namely; ISP functional model, AXI BFM, configurable clock generators model, configurable reset generator model, memory model, configurable interconnect etc. All these are pure SystemC models. AXI BFM is provided to interact with other part of the world.

ISP RTL block needs exhaustive verification, which is possible only when the RTL is ready. But, development of RTL design takes time, which means verification of RTL design can't be possible before it becomes available. To shorten this sequential activity, functional model of ISP is used to prepare the early verification environment.

A SystemVerilog test bench wrapper is created over SystemC/TLM ISP sub-system. This SystemVerilog test bench interface with the RTL verification environment.

## D. Virtual Platform Sub-system

When all components of platform are in TLM/C, means C/C++ are used as modeling language; we call it a Pure Virtual platform. In typical verification environment, generally all verification components are not only TLM based but also of different verification languages thus making it a Multi-language heterogeneous simulation environment. For developing early verification environment, TLM based sub-system is developed which consists of every block in TLM/C. This TLM based Sub-system is model of RTL.

In the above mentioned RTL verification environment, a processor model is used which enables us to early develop 'C' test cases for programming of RTL registers/memories via CPU interface. The challenge is to keep the verification environment independent of "C" test cases. We don't wish to compile every time whenever there is change in application code. To be able to achieve this, a sub-system is designed which consists of models of bus interfaces, like AXI BFM, a "generic" processor model, model of memory, etc. an independent "C" program/test case is written to do all the programming and configuration, which in turn runs on processor model of this sub-system. This sub-system is active element in programming phase, but becomes passive once the programming is complete.

Virtual platform sub-system can be represented in following block diagram.



Fig. 3. Virtual Platform Sub-system

## E. Virtual Register Interface (VRI)

Today, most of the embedded test infrastructure uses some adhoc mechanism like "shared memory" or synchronization mechanism for controlling simple Bus functional models (BFMs) from embedded software.

In order to provide full controllability to the "C" test developer over these verification components, a virtual register interface layer is created over these verification environments which provides the access to the sequences of these verification environment to the embedded software enabling configuration and control of these verification environments to provide the same exhaustive verification at SoC Level.

This approach addresses the following aspects of verification at SoC Level:

- Configuration and control of verification components from embedded software.
- Reusability of verification environments from IP to SoC.
- Enables reusability of testcases from IP to SoC.
- Providing integration testcases to SoC team which is developed by IP verification teams.

It has been achieved by using Virtual Register Interface (VRI) layer over Verification components [18]. VRI layer over verification components is –

➢ A virtual layer over verification environment to make it controllable from embedded software

➢ Provides high level C APIs hiding low level implementation



Fig. 4.   Virtual Register Interface (VRI)

An example of C test case using VRI interface is as follows –

```
vr_enet_packet pkt;
vr_enet_packet rx_pkt;
rx_pkt.data = new vri_uint8_t[2000]; //create buffer for
receiving data

pkt.packet_kind = ETHERNET_802_3;
pkt.data_length = 0; //RANDOM DATA
pkt.dest_addr_high = 0x11ff;
pkt.src_addr_high = 0x2288;
pkt.tag_kind = UNTAGGED;
pkt.tag_prefix = 0x1234;
pkt.s_vlan_tag_prefix = 0x5678;
pkt.err_code = 0;
for (int i=0;i<100;i++) {
  pkt.dest_addr_low = i;
  pkt.src_addr_low = i+1;
  enet_send_pkt(0,&pkt);        //send packet to ENET UVC
instance0 (MAC)
  enet_recv_pkt(1,&rx_pkt); //receive packet from ENET UVC
instance1 (PHY)
  compare_pkt(pkt,rx_pkt);
};
```

*F.  Flow used for Design Verification*

Much before arrival of RTL, C/Python model of image signal processor designs is developed for algorithm evaluation. Then, TLM/SystemC model of the design is created from C/Python model. After proper exhaustive validation of the model with required test vectors, the model qualifies as an Executable Golden Model or Executable Specification means a 'living' benchmark for design specification. Enabling the use of TLM Model as DUT expedites development and better proofing of the verification environment with wider test vectors without waiting for RTL to be available.

Standard 'interfaces' are used to enable the reuse of verification components. In addition to standard method of bus-interface or signals level connectivity, UVM Multi-Language Open Architecture is used to connect System Verilog TLM port directly to SystemC TLM port which gives advantage of better simulation speed and better development/debug cycle in addition of clean, better and easy connectivity/integration of blocks. Presence of TLM components gives us flexibility to make backdoor direct access to the DUT registers and memories.



Fig. 5.   Early development of Verification Environment using TLM Model

A processor model is used which enables us to early develop 'C' test cases for programming of RTL registers/memories via CPU interface. Same 'C' test cases are used for controlling the SystemVerilog UVC's using Virtual Register interface (VRI) layer. In our verification environment, alternative Host interface path is used to do programming of configurable blocks using SystemVerilog UVM based test cases.

In both above cases, control/data flows across both TLM and bus interface boundaries. This method enhances the chances of re-using different already existing blocks in flow. IP-XACT based tools are also used for automatically configuring the environment for various designs.

By the time RTL arrives, complete verification environment and test-vectors are ready with sufficient sanctity, thus eliminating the number of verification environment issues which may arise when actual RTL verification is started. When RTL arrives, the TLM/SystemC model is simply replaced with RTL block with reuse of maximum of other verification components. This enhances the rapid/regress testing of design immediately. Also same C test cases can be run on actual core.



Fig. 6.    : Reuse of early developed Verification Environment

### III.    RESULTS

Using TLM reference model of the RTL, UVM based Environment for verification of design is developed without waiting for RTL to be available. Significant reduction in overall verification time of the design is achieved.

### IV.    CONCLUSIONS

TLM/SystemC reference model of the design is the key component to enable the early development of Verification Environment without waiting for RTL to be available. UVM based early verification Environment is developed using TLM/SystemC reference model of the design. Verification Environment is developed both with Host interface and Core using Virtual Register Interface (VRI) approach. IP-XACT based tools are used for automatically configuring the Verification Environment. Testing of features of Verification Environment at TLM abstraction level runs faster and thus, it overall speeds up functional verification. Same environment can be reused from IP level to SOC level or from one SOC to another SOC with no/minimal change. Verification Environment is reusable both vertically and across projects thus saving further time across projects.

### REFERENCES

[1] Abhishek Jain, Giuseppe Bonanno, Dr. Hima Gupta and Ajay Goyal, "Generic System Verilog Universal Verification Methodology Based Reusable Verification Environment for Efficient Verification of Image Signal Processing IPs/SOCs", International Journal of VLSI Design & Communication Systems, 2012.

[2] Abhishek Jain, Piyush Kumar Gupta, Dr. Hima Gupta and Sachish Dhar, "Accelerating System Verilog UVM Based VIP to Improve Methodology for Verification of Image Signal Processing Designs Using HW Emulator", International Journal of VLSI Design & Communication Systems, 2013.

[3] Abhishek Jain, Mahesh Chandra, Arnaud Deleule and Saurin Patel, "Generic and Automatic Specman-based Verification Environment for Image Signal Processing IPs", Design & Reuse, 2009.

[4] Mark Glasser, "Open Verification Methodology Cookbook", Springer, 2009.

[5] Iman, S., "Step-by-Step Functional Verification with SystemVerilog and OVM", Hansen Brown Publishing, ISBN: 978-0-9816562-1-2, 2008.

[6] Rosenberg, S. and Meade, K., "A Practical Guide to Adopting the Universal Verification Methodology (UVM)", Cadence Design Systems, ISBN 978-0-578-05995-6, 2010.

[7] Stuart Swan, "An Introduction to System Level Modeling in SystemC 2.0",  Cadence Design Systems, Inc. May 2001.

[8] Adam Rose, Stuart Swan, John Pierce, Jean-Michel Fernandez, "Transaction Level Modeling in SystemC", Cadence Design Systems, Inc., 2005.

[9] Frank Ghenassia, "Transaction Level Modeling with SystemC - TLM Concepts and Applications for Embedded Systems", ISBN: 978-0-387-26232-1, 2010.

[10] Daniel D. Gajski,"System-Level Design Methodology", www. cecs.uci.edu /~gajski, 2003.

[11] Lukai Cai and Daniel Gajski, "Transaction Level Modeling: An Overview",   {lcai, gajski}@cecs.uci.edu, 2003.

[12] Farooq Khalid Chughtai, "Accurate Performance Exploration of System-  on-Chip using TLM", 2012.

[13] Thorsten, "System Design with SystemC", Kulwar Academic Publishers Group, 2002.

[14] Sandeep Jana, Geetika Agarwal and Kishore Sur, "Unique Approach for System Level Verification Using Scalable and Reusable Verification IP with TLM Infrastructure". CDNLive 2010.

[15] Sandeep Jana, Krishna Kumar, Sonik Sachdeva, Swami Venkatasen and Debjoyoti Mukherjee, "TLM Based software control of UVCs for Vertical Verification Reuse:, CDNLive 2012

[16] Accellera Organization, Inc. Universal Verification Methodology (UVM) May 2012.

[17]     IEEE Computer Society. IEEE Standard for System Verilog-Unified Hardware Design, Specification, and Verification Language - IEEE 1800-2009. 2009.

[18] Virtual Register Interface Layer over VIPs from Cadence Design System.

[19] Spirit information, http://www.spiritconsortium.org.

[20] Accellera VIP TSC, UVM Register Modelling Requirements, www.accellera.org /activities/vip/

[21] www.ovmworld.org

[22] www.SystemVerilog.org

[23] www.uvmworld.org

[24] http://www.accellera.org/community/uvm/

[25] www.systemc.org

[26] www.synopsys.com

### AUTHOR'S PROFILE

**Abhishek Jain, Technical Manager, STMicro-electronics Pvt. Ltd.**
**Research Scholar, JBS, Jaypee Institute of Information Technology, Noida, India.**
**Email: ajain_design@yahoo.co.in;**
**abhishek-mmc.jain@st.com**
Abhishek Jain has more than 11 years of experience in Industry. He is driving key activities on Functional Verification Flow in Imaging Division of STMicroelectronics. He has done PGDBA in Operations Management from Symbiosis, M.Tech in Computer Science from IETE and M.Sc. (Electronics) from University of Delhi. His main area of Interest is Project Management, Advanced Functional Verification Technologies and System Design and Verification especially UVM based Verification, Emulation/Acceleration and Virtual System Platform. Currently, he is doing Research in Advanced Verification Methods for Efficient Verification Management in Semiconductor Sector. Abhishek Jain is a member of IETE (MIETE).

**Dr. Hima Gupta, Associate Professor, Jaypee Business School (A constituent of Jaypee Institute of Information Technology University), A – 10, Sector-62, Noida, 201 307 India.**
**Email: hima_gupta2001@yahoo.com**
Dr. Hima has worked with LNJ Bhilwara Group & Bakshi Group of Companies for 5 yrs. and has been teaching for last 11 years as Faculty in reputed Business Schools. She also worked as Project Officer with NITRA and ATIRA at Ahmedabad for 5 years.

She has published several research papers in National & International journals.

**Sandeep Jana, Staff Engineer, STMicroelectronics Pvt. Ltd.**
**Email: sandeep.jana@st.com**
Sandeep Jana is Staff Engineer at STMicroelectronics managing the TLM based Verification activities at Greater Noida. He has an expertise of over seven years in various aspects of ESL domain such as TLM modeling, Architectural exploration, Platform Integration, Mixed language Platforms, Advanced Verification Methodologies etc. He has been with ST since last 6 years and was previously working in VLSI group of HCL Technologies in their ESL domain. He has a B.Tech degree in Electronics Engineering from MDU Rohtak.

**Krishna Kumar, STMicroelectronics Pvt. Ltd.**
**Email: krishna.kumar@st.com**
Krishna Kumar with almost 12 years at STMicro-electronics has experience in ESL and Placement and Routing of FPGA software tool chain. He holds B.Tech degree in Computer Engineering from Aligarh Muslim University, India.

# Development of Rest Facility Information Exchange System by Utilizing Delay Tolerant Network

Masahiro Ono

Department of Information and
Communication
Engineering, Tokyo Denki
University Tokyo,
Japan

Kei Sawai

Department of Information and
Communication
Engineering, Tokyo Denki
University Tokyo,
Japan

Tsuyoshi Suzuki

Department of Information and
Communication
Engineering, Tokyo Denki
University Tokyo,
Japan

*Abstract*—**In this paper, we propose temporary rest facilities information exchange system among many people unable to get home by utilizing Delay Tolerant Network (DTN) after a disaster. When public transportation services are interrupted by the disaster, those people try to get home on foot while taking a rest at the facility. However, it is difficult for those people to obtain information of temporary rest facilities provided hurriedly, because communication infrastructures in the disaster area are disconnected by the disaster damage. Therefore, we propose a method to exchange the information among those people mutually by using mobile device via DTN for diffusion of the information. By using DTN, those people can communicate with each other by using mobile device and use the rest facility on the basis of the information even if the communication infrastructures are disconnected. Then, we develop mobile device application software to exchange the rest facility information among the people via DTN. In order to evaluate the application, we verified the communication performance in practical experiments. The experimental results showed the developed application had sufficient performance to exchange the information of the rest facility via DTN. Then, we verify the diffusivity of the rest facility information by a network simulation. The simulation results showed that the rest facility information was diffused widely and effectively to those people.**

*Keywords*—*Delay Tolerant Network; rest facility; disaster; communication infrastructure; simulation*

## I. INTRODUCTION

The problem of people (e.g. commuters, students, etc.) unable to get home after a disaster has attracted attention since the Great East Japan Earthquake. Public transportation services are suspended by disaster damages when a major disaster occurs. For example, in the Great East Japan Earthquake of 2011, a lot of railways in the metropolitan area suspended a passenger transport service for a long time [1]. In such case, it has been pointed out that a large number of people who have usually commuted by public transportation are unable to get home by public transportation unavailable.

Those people try to return to their home on foot, and spend a lot of time in walk. However, there are risks of secondary disasters including accidents caused by fatigue due to long-time walk; e.g. myocardial infarction, depression from mental stress and so on. Therefore, they have to need a rest on their way home. In fact, about 40,000 people in the metropolitan area took a rest on their way home after 5 hours in the Great East

Japan Earthquake [2]. Hence, as a facility to assist those people, rest facilities are prepared.

People unable to get home need information of those rest facilities; e.g. location, capacity, etc. Information of rest facilities is provided by prefectural and city governments in advance and it is shown on the disaster prevention map which was made in preparation for a disaster. However, in the large scale disaster, there are cases where new temporary rest facilities will be provided hurriedly since a rest facility is overcrowded more than a capacity of the facility [3] [4]. To spread information of these rest facilities not listed in the disaster prevention map to a large number of people, the use of mobile devices such as cell phones or smart phones that many people use are effective. On basis of this information, they are able to move to a temporary rest facility without being at a loss by using a map or GPS function provided in each device.

On the other hand, high degree of risk about disconnection of various communications infrastructure in urban area in a major disaster is known. In the Great East Japan Earthquake, about 1.94 million landlines function was suspended, and wave transmission was stopped in about 28,650 base station [5]. Furthermore, there are cases that the communication function is down by congestion even if the communication infrastructure is available. Therefore, under these situations, it is difficult to spread the information of temporary rest facilities to people unable to get home by the communication infrastructure. Hence, a diffusion method that can spread such information to those people even if the communication infrastructure is unavailable is needed.

As a communication method to transmission information under those situations, the Delay Tolerant Network (DTN) is useful [6]. DTN is a technique for performing communication in an environment in which End-to-End communication path is not always connected. DTN performs direct communication between devices by utilizing mobile devices. Figure 1 shows how to transfer the data by DTN. Store-and-forward scheme is mainly used for data transfer method in DTN. In store-and-forward method, in order to send the data from a source to a destination, the mobile device stores the data if it cannot communicate with other devices, then the mobile device transfers the data if it can communicate with other devices. In this way, DTN can transmit information from a source to a destination by communicating with near devices while passing

each other even if it is difficult to maintain communication connectivity continuously. DTN is suitable for rest facilities information transmission method in post-disaster situation because it can be used in portable smart phones.

In this research, we propose information exchange method to diffuse information of temporary rest facilities to people unable to get home by utilizing DTN in communication infrastructure unavailable after the large scale disaster.



Fig. 1.  Delay Tolerant Network

## II.  RELATED WORKS

A portable satellite communications system is developed by NTT Laboratories as an information transmission method of disaster [7]. This system provides the means of communication at rest facilities or shelters in a disaster. By using the vehicle-mounted system with an antenna and telephone device in the post-disaster environment without the communication infrastructure, telephone call is possible via the satellite communications connection. However, a long distance communication is impossible when the base station is interrupted by the disaster damages because each base station has to be put within the communication range mutually. Further, it is difficult to transmit the rest facilities information for people unable to get home since a large number of people cannot use the system at a time. In addition, the movement of the vehicle-mounted system is difficult because it is considered that road of inner city is congested in a disaster.

As one of information transmission method of disaster, Heli-TV system is used by the disaster prevention center managed by Ministry of Land, Infrastructure, Transport and Tourism [8]. This system sends video data taken from a helicopter in the disaster area to disaster prevention center via satellite communications. A car equipped with satellite communication antenna moves in the post-disaster environment without communication infrastructure, and can relay the communication. However, it is difficult to utilize this system as a general communication system which sends the information to people unable to get home because this system is assumed to provide dedicated connection between the helicopter and the disaster prevention center. In addition, the movement of base station for this system by car is difficult because it is considered that road of inner city is congested in a disaster.

As a study utilizing DTN, MONAC of Android application was developed by Teranishi et al [9]. MONAC is a system that aims to transmit the disaster area information in an environment where communication infrastructure is disconnected by the disaster damage. MONAC exchanges the disaster area information between devices (e.g. mobile phone and smart phone) which have passed each other in the environment by DTN, then the information sends to Twitter when a device can connect to the communication infrastructure. However, MONAC assumes that the user is able to connect to the communication infrastructure in the mobile range; it might not work in the large environment where the communication infrastructure is interrupted.

A study of Data Gathering and Sharing based on DTN was proposed by Sun et al [10]. In this study, emergency personnel exchange the disaster area information with each other by utilizing DTN in the environment where the communication infrastructure is interrupted, and then the collected information is transmitted to the disaster response headquarters when the emergency personnel can connect the communication infrastructure. However, the proposed method assumes that the emergency personnel can connect to communication infrastructure in their moving range. Therefore, it might not work in the large environment where the communication infrastructure is interrupted.

In these studies, it is difficult to transmit information in the environment where the communication infrastructure is interrupted by disaster damages caused by the large scale earthquake. Therefore, we propose a system that can transmit the rest facility information to people unable to get home in such environment.

## III.  PROPOSED SYSTEM

### A.  Assumed environment

An assumed environment of the proposed system is large environment in which the communication infrastructure is disconnected by the disaster damage, and the communication function is stopped because of the communication amount increase. In such environment, information collection is difficult because communication services such as telephone, television broadcasting, e-mail and the Internet service cannot be used. In addition, many people in urban area cannot go back home for a long time because public transportation services are interrupted by disaster damages. Therefore, population density increases by the advent of people unable to get home on foot in urban area. These conditions show the increase of opportunity of communication by utilizing DTN since a large number of people pass each other on the sidewalk and those people move slowly by the traffic jam.

### B.  Request of function

This system aims to transfer the rest facility information to people unable to get home in the environment that has no communication infrastructure. Therefore, the communication method capable of transmitting information in such environment is needed. Further, this system is required to be used on the move because those people use it on their way home. It is also necessary to decide the rest facility information for the transmitting information source.

### C.  Outline of system

From the functional requirement in Section B, we propose the system to transmit the rest facility information which

received in the rest facility to people unable to get home by utilizing DTN. DTN can transmit the rest facility information even in the environment that has no communication infrastructure because DTN is available on mobile devices and transmit the information by the device communication. In addition, DTN is useful as a method for long distance communication by exchanging the information among those people who go back home because it is possible to transmit the information to the distant with a device movement.

On basis of the survey report by Japanese cabinet office about necessary information for the people who went back home when the Great East Japan Earthquake had occurred, we decided the rest facility information to be transmitted is the information about position of the rest facility, capacity, presence or absence of relief supplies and toilet [11]. To find the rest facilities, the location information of the rest facilities is required. And, as information to determine the availability of rest facilities, it is necessary to transmit information about the capacity of the rest facility. In addition, in the case that the people have to walk for a long time to get home, information of the location of toilets is important. It is considered demand of toilets increases in the large scale disaster, because a lot of people walk for a long time to get home [12]. Therefore, the information about the number and the presence of toilets is important in supporting those people. There is a case that the relief supplies such as foods and beverages in the victim support facilities such as rest facilities are distributed. In a disaster, many stores are closed and vending machines don't work by the disaster damage. Further, it is difficult to ensure offer of foods and beverages due to the increase of the demand for foods and beverages by those people [13]. Therefore, the information about the presence of foods and beverages distribution is important in the use of rest facilities.

Figure 2 shows the overview of the system. If the communication infrastructure cannot connect by the disaster damage, the rest facility information transmission server is installed in the rest facility. If people unable to get home come to the rest facility, the server directly sends the rest facility information to their devices. The device of those people saves the received rest facility information. If those people exit the rest facility after taking a rest, their devices send the rest facility information to other persons' devices within communication range by utilizing DTN. The information is widely propagated by being transmitted to surrounding device in communication range from a device repeatedly via DTN.

Therefore, it is possible to transmit the rest facility information in the environment where the communication infrastructure is interrupted, and it is possible to search the rest facility by using the received rest facility information. The proposed system can transmit the rest facility information to people unable to get home, because a lot of people try to get home on foot in urban area in disaster situation. In this study, we develop functions to exchange the rest facility information by utilizing DTN between those people, and then we verify the diffusivity of the rest facility information by a simulation.



Fig. 2. System outline

### D. Development of DTN application

In this system, the device needs to have functions of exchanging the rest facility information by utilizing DTN among people unable to get home. Therefore, the device to be used needs to have portability to be able to carry in a disaster, and the device-to-device communication function to communicate by utilizing DTN. Further, in this study, it is necessary to use the devices in which hardware and functions are relatively uniform because multiple devices are used. Therefore, in this study, we use the iPhone of Apple, Inc. in which these conditions are satisfied, and develop the system as the function of the iOS application.

Figure 3 shows the application of the algorithm. First, if the device detects the information transmission server in the rest facility by utilizing DTN, the device saves the rest facility information transmitted from the server. The device connects to the other devices when they were detected within communication range, and then the device transmits the rest facility information stored to other devices. The device, then, gets the rest facility information which had stored by other devices. If the device detects again other devices, it repeats the same operation from the detection of other devices of Figure 3.

Figure 4 shows the screen of the developed application. In the top part of the screen, there are the "add information screen" to enter the rest facility information and the "add button" to save the rest facility information. In the center part of the screen, there are the "start DTN" button to start the communication by utilizing DTN, the "finish DTN" button to terminate the DTN communication and the "log screen" that outputs an operation history of the application and communication status. In the lower part of the screen, there is the "information preservation screen" to output the rest facility information stored in the device. Application has the following three functions.

*1)* *Registration of the rest facility information in the add information screen.*
*2)* *Transmission of the rest facility information stored to other devices by utilizing DTN.*
*3)* *Save of the rest facility information received from other devices.*

The application saves the rest facility information input in the add information screen to the database that was implemented in the application. And, the application transmits the rest facility information stored to other devices in the communication range by utilizing DTN. Then, the application stores the rest facility information received from other devices to the database. The stored rest facility information can be viewed in the information preservation screen. The transmission of the rest facility information by utilizing DTN and preservation of the rest facility information obtained is carried out automatically. The registration function of the rest facility information is necessary for the transmission information system. We equipped the application with the function (1) for the creation of the rest facility information. The function of (2) is capable of transmitting the rest facility information obtained from other devices by the function of (3). We employed the Bluetooth system as the communication system, because it is mounted on a lot of mobile devices and is battery friendly [14].



Fig. 3. Algorithm of DTN application



Fig. 4. iOS application

## IV. PERFORMANCE EVALUATION OF COMMUNICATION CONNECTIVITY OF DEVELOPED DTN APPLICATION

### A. Experimental outline

In this experiment, we verified the communication connectivity of our developed DTN application for sending and receiving the rest facility information. Figure 5 shows the experimental overview. The performance of the communication connectivity was evaluated by utilizing two iPhone (Type: 3GS, 4S), it measured the processing time between two communication devices in the exchanging the rest facility information. Thus the processing time was measured by changing the distance between two communication devices. The processing time was defined the time between finish to receive the information and start to connect other device. We measured the processing time to the distance to disconnect area.



Fig. 5. Evaluation of communication quality

### B. Experimental conditions

The experiment was executed at the corridor of kita-senju campus in Tokyo Denki University. There were not some obstacles in the movement line. The processing time was measured at 10 [m] interval. The data size of sent information was configured to 1512 [Byte] by discussing the size of the rest facility information. The processing time was measured to 3 times in each distance, and then it was calculated the average value.

### C. Experimental results and Discussion

Figure 6 shows the processing time of each measurement distance. The maximum distance of the DTN by utilizing our developed application was 90 [m]. In the average value of the processing time, iPhone 4S spent to 5.65 [sec] and iPhone 3GS spent to 6.11 [sec].

The average of walking speed is defined to 4.0 [km/h] [15], and then the distance to walk in 6.11 [sec] is 6.79 [m]. Figure 7 shows the situation of the DTN in walking. Then the distance to walk in the processing time fit in the distance of the DTN. Therefore we confirmed that the DTN was able to exchange the rest facility information in walking situation.



Fig. 6. Experimental result of communication quality



Fig. 7. Verification of communication quality

## V. Verification Of DTN Connection At The Time Of The Walk

### A. Experimental overview

We verified exchange of the rest facility information between two pedestrians passing each other by utilizing our developed application in walking. Figure 8 shows the schematic diagram of the experiment. In this experiment, we used two iPhone (Type: 3GS, 4S) as communication device. Device holder started walking from out of communication range measured in previous chapter, and then we evaluated the capability of exchange of the rest facility information by utilizing DTN.

Fig. 8. Evaluation of proposed DTN application

### B. Experimental conditions

The experiment was executed at the corridor of kita-senju campus in Tokyo Denki University. In this experiment, the device holders activates the function of DTN by pressing the start DTN button of application in the out of communication range, and then the terminal holders start walking. We experimented in the environment that does not have an obstacle on a straight line between the devices. We measured the distance between two devices from the start to the finish of communication to other device and the processing time. These parameters were measured to 5 times.

### C. Experimental results and Discussion

The experimental results showed that developed application could exchange the information by utilizing DTN in the between two pedestrian. Figure 9 shows the measurement value of the processing time and the end-to-end distance from the start to the finish of communication. Exchanging information was finished before device holders passing each other. Processing time of this experiment shows a value close to the results of previous chapter. In addition, the detectable distance between two devices was shorter than the experimental results in previous chapter. We considered that the detection of the other device is difficult due to move two devices. However, the performance of DTN is sufficient values, because exchanging information was finished before device holders passing each other.

Fig. 9. Experimental result of DTN application

## VI. Performance Evaluation By Using Environmental Simulation

### A. Experimental overview

To verifying the diffusivity of the rest facility information by utilizing our proposed method, it requires to evaluate the conveyance of information by connection DTN in the environment that has no communication infrastructure. However the evaluation of the information diffusivity for many people in actual environment is distant experiment. Therefore, we evaluated the conveyance of information in simulation test. Then we adopt the "Opportunistic Network Environment Simulator (The One simulator)" as DTN evaluation software [16]. The One simulator is able to configure the parameter of the wireless communication protocols, the mode of locomotion and the number of people, etc., and then it simulate the migration of various node and the conveyance of information in the setup map. In this simulation test, we expanded capacity that setting the number of communication node in the One simulator.

Figure 10 shows the simulation overview by utilizing the One simulator. We simulated the information diffusivity from a rest facility by utilizing our proposed method in disaster area. People unable to get home is deployed on the setup map, and then randomly placed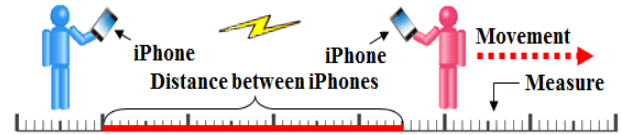 people moves on the setup pathway after running simulation. By the report of Japanese cabinet office, almost of those people started to move in after about 1 hour. Therefore we set those people with the rest facility information start to move in 1 hour after the start of simulation. The simulation time was set to 2 hours, we measured the number of those people conveyed the rest facility information and the maximum distance that information has been received at 10 minute intervals.

Fig. 10. Simulation of proposed system

### B. Experimental conditions

The parameters configuration in simulation test was based on the data of people unable to get home about the Great East Japan Earthquake to reproduce the environment difficult get home. We are shown in the parameters details below. Simulation area was defined the Chiyoda city, Tokyo as the area of a lot of the population of daytime. The area with many people is considered to be vulnerable to damage [17].

Determination of the migration pass of those people is based on the "person trip survey" conducted by Tokyo metropolitan government [18]. "Person trip survey" is discussed the movement of people in Tokyo, it is configured

the data that is the start and the finish point of traffic, the purpose of the movement, the means of transportation and the number of the people who are on the move. Then this survey involving the data that is the number of the person for the purpose of returning home is well suited to reproduce the activity of those people. In this study, we decided the start and the finish point of those people by the area data of "person trip survey", and then this data are reflected in simulation map.

The number of those people is also decided by "person trip survey". "Person trip survey" is the data at the time of no disaster. Therefore, we determined the number of people moving in per unit distance to get home by "person trip survey" and the data of the ratio of people that get home on foot in per unit distance in the Great East Japan Earthquake [19]. The number of the people who start returning home in per unit time of the simulation was defined by "person trip survey" and the data of the ratio of people that start returning home in per unit time in the Great East Japan Earthquake. Thus total number of people moving was defined as 7903 people in this simulation. Walking speed of the those people set to 4 [km / h] on the basis of the public data of Japanese cabinet office.

Then we set that every those people has an intelligent communication device, the entire device is installed our developed DTN application. The communication distance of the device set to 90 [m], communication time was configured to 5.65 [sec]. Communication protocol is adopted the Bluetooth, the number of device connecting at the same time set to 8.

### C. Experimental results and Discussion

Figure 11 shows the rate of the number that the person who received the rest facility information in the amount of people unable to get home. The number of the person who received the rest facility information based on the experimental result was 7871 people after 60 [min] from starting the information diffusion, and then the rate of received person was 99.6 [%]. Moreover we confirmed that the diffused information was received to those people that is the rate of 86.9 [%] in total number of people at 10 [min] after starting the information diffusion. The maximum propagation distance that was spread the information was shown to 2332.96 [m] in this simulation test. Then we confirmed that the our proposed DTN system was enable to spread the rest facility information to many those people, it was not needed enough time to spread the data in target area. This result was considered that the effect of staying a lot of those people in around environment of those people from starting simulation.

We considered the information diffusion capability of our proposed system in Tokyo metropolitan area by the experimental results. Area of the Tokyo metropolitan government without islands is 1782.89 [km$^2$], and there is the rest facility of 1030 places [20] [21]. Therefore the area that a rest facility place superintends environment defined 1730961.17 [m$^2$]. Figure 12 shows the environment area. Diffusion of the rest facility information is significant to conveying the information to the distance more than the radius shown in Figure 12. In this simulation test, the maximum distance of the conveyed distance was measured to 2332.96 [km], thus the distance of the conveyed information was longer than the 742.47 [m] as the radius.

Thus, it was considered that our proposed system diffused effectively to wide range by these results. Moreover the rate of the number of people unable to get home received the information was measured to 99.6 [%] at 60 [min] after starting diffusion of the information. Therefore our proposed system was suggested to diffuse effectively to wide range without the environment existing communication infrastructure such as disaster area.



Fig. 11. Ratio of number of people who have information



Fig. 12. Area to be responsible of rest facility

### VII. CONCLUSION

In this study, we proposed the exchange system of temporary rest facility information to people unable to get home by utilizing Delay Tolerant Network in urban area in post-disaster environment. The iOS application based on proposed method was developed, and then the communication performance was verified in the field test. In the test, we confirmed proposed method can exchange the temporary rest facility information between two pedestrian passing each other. Moreover, information diffusivity by the proposed method was verified by using the One Simulator. The result showed the proposed method is available for diffusion of the temporary rest facility information to a large number of people in urban area. In the simulation, we used data of Chiyoda-ku, a heavily populated city in Japan, for evaluation of proposed method by the One Simulator. Therefore, it is estimated that the information was frequently exchanged by many people passing each other.

In future work, to evaluate availability of the proposed method, the experiment should be carried out under several conditions; e.g. using data of sparsely populated area, changing walking speed of people on basis of a crowded situation, etc. Then we consider the simulation test estimate the precision

result by referring the increase and decrease rate of density of population in chronological order.

### REFERENCES

[1] Ministry of Land, Infrastructure, Transport and Tourism, http://www.mlit.go.jp/common/000208774.pdf, access: 10 January 2014.

[2] Mitsubishi Research Institute, http://www.mri.co.jp/NEWS/press/2011/2029010_1401.html, access: 10 January 2014.

[3] Setagaya Ward Office, http://www.city.setagaya.lg.jp/kurashi/104/141/563/571/573/d00033198_d/fil/33198_1.pdf, access: 10 January 2014.

[4] Itabashi Ward Office, http://www.city.itabashi.tokyo.jp/c_kurashi/041/attached/attach_41697_1.pdf, access: 10 January 2014.

[5] Ministry of Internal Affairs and Communications, http://www.soumu.go.jp/main_content/000113017.pdf, access: 10 January 2014.

[6] K. Fall, "A delay-tolerant network architecture for challenged internets", Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications, pp.27-34, 2003.

[7] H. Takashi, O. Kazuhiko, I. Yutaka, Y. Hidekuni, "Satellite communications systems used in disaster recovery operations after the Great East Japan Earthquake and tsunami", The Institute of Electronics, Information andCommunication Engineers Technical Report, Satellite Telecommunications 111(336), pp.109-113, 2011.

[8] Ministry of Land, Infrastructure, Transport and Tourism, http://www.mlit.go.jp/bosai/disaster/bousaicenter/bousaicenter.htm, access: 10 January 2014.

[9] Y. Teranishi, S. Shimojo, "MONAC: SNS message dissemination over smartphone-based DTN and cloud", The IEEE International Conference on Peer-to-Peer Computing, pp158-159, 2011.

[10] W. Sun, T. Kitani, N. Shibata, K. Yasumoto, "A Data Gathering and Sharing Proposal for Disaster Relief based on DTN", The Special Interest Group Technical Reports of IPSJ, 2009(20), pp.61-66, 2009.

[11] Tokyo Metropolitan Gaverment, http://www.bousai.metro.tokyo.jp/japanese/kitaku_portal/tmg/pdf/9-11122kaigi.pdf, access: 10 January 2014.

[12] Japan Toilet Labo, http://www.toilet.or.jp/dtinet/gaiyo.pdf, access: 10 January 2014.

[13] H. Nagato, N. Nishimura, K. Yamamoto, "Supporting Measures for the Difficult to Return Home at The Earthquake Disaster", Infrastructure Planning Review , No.22, pp.265-270, 2005.

[14] Agilent Technologies, http://cp.literature.agilent.com/litweb/pdf/5989-4204JAJP.pdf, access: 10 January 2014.

[15] Central Disaster Prevention Council, http://www.bousai.go.jp/kaigirep/chuobou/senmon/shutohinan/pdf/sanko03.pdf, access: 10 January 2014.

[16] A. Keranen, J. Ott, and T. Karkkainen, "The ONE simulator for DTN protocol evaluation", International Conference on Simulation Tools and Techniques, No.55, pp.1-10, 2009.

[17] Tokyo Metropolitan Gaverment, http://www.toukei.metro.tokyo.jp/tyukanj/tj-index.htm, access: 10 January 2014.

[18] Tokyo Metropolitan Gaverment, http://www.tokyo-pt.jp/person/index.html, access: 10 January 2014.

[19] U. Hiroi, N. Sekiya, R. Nakajima, S.Waragai, H. Hanahara, "Questionnaire Survey concerning Stranded Commuters in Metropolitan Area in the East Japan Great Earthquake", Institute of Social Safety Science, No.15, pp.343-353, 2011.

[20] Tokyo Metropolitan Gaverment, http://www.metro.tokyo.jp/PROFILE/map_to.htm, access: 10 January 2014.

[21] Tokyo Metropolitan Gaverment, http://www.bousai.metro.tokyo.jp/japanese/tmg/pdf/240113kitakukihon.pdf, access: 10 January 2014.

# New Method Based on Multi-Threshold of Edges Detection in Digital Images

Amira S. Ashour[1,2] , Mohamed A. El-Sayed[1,3]
[1]Dept of CS, Computer &IT College, Taif Univ., KSA;
[2]Dept of Electronics & Electrical Communications Eng.
Faculty of Engineering, Tanta Univ., Egypt
[3]Dept of Math, Faculty of Science, Fayoum Univ., Egypt

Shimaa E. Waheed[4,5] , S. Abdel-Khalek[4,6]
[4]Dept of Math, Faculty of Science, Taif University, KSA;
[5]Dept of Mathematics, Faculty of Science, Benha
University, Egypt
[6]Dept of Math, Faculty of Science, Azhar Univ., Egypt

*Abstract*—**Edges characterize object boundaries in image and are therefore useful for segmentation, registration, feature extraction, and identification of objects in a scene. Edges detection is used to classify, interpret and analyze the digital images in a various fields of applications such as robots, the sensitive applications in military, optical character recognition, infrared gait recognition, automatic target recognition, detection of video changes, real-time video surveillance, medical images, and scientific research images. There are different methods of edges detection in digital image. Each one of these methods is suited to a particular type of images. But most of these methods have some defects in the resulting quality. Decreasing of computation time is needed in most applications related to life time, especially with large size of images, which require more time for processing. Threshold is one of the powerful methods used for edge detection of image. In this paper, We propose a new method based on different Multi-Threshold values using Shannon entropy to solve the problem of the traditional methods. It is minimize the computation time. In addition to the high quality of output of edge image. Another benefit comes from easy implementation of this method.**

*Keywords—image processing; multi-threshold; edges detection; clustering*

## I. INTRODUCTION

In many applications of image processing, the gray levels of pixels belonging to the object are quite different from the gray levels of the pixels belonging to the background. Thresholding becomes then a simple but effective tool in edge detection to separate objects from the background. Edge detection using thresholding is significant importance in many research areas[1,2]. Since, the edge is a prominent feature of an image; it is the front-end processing stage in object recognition and image understanding system. The detection results benefit applications such as automatic target recognition [3], medical image applications [4], and detection of video changes [5].

Edge detection can be defined as the boundary between two regions separated by two relatively distinct gray level properties[6]. The causes of the region dissimilarity may be due to some factors such as the geometry of the scene, the radio metric characteristics of the surface, the illumination and so on [7]. An effective edge detector reduces a large amount of data but still keeps most of the important feature of the image. Edge detection refers to the process of locating sharp discontinuities in an image. These discontinuities originate from different scene features such as discontinuities in depth, discontinuities

in surface orientation, and changes in material properties and variations in scene illumination [8,9].

Most of the classical methods for edge detection based on the derivative of the pixels of the original image are Gradient operators, Laplacien and Laplacien of Gaussian (LOG) operators [7]. Many operators have been introduced in the literature, for example Roberts, Sobel and Prewitt [10-14]. Edges are mostly detected using either the first derivatives, called gradient, or the second derivatives, called Laplacien. Laplacien is more sensitive to noise since it uses more information because of the nature of the second derivatives.

Gradient based edge detection methods, such as Roberts, Sobel and Prewitts, have used two linear filters to process vertical edges and horizontal edges separately to approximate first-order derivative of pixel values of the image. Marr and Hildreth achieved this by using the Laplacien of a Gaussian (LOG) function as a filter [15]. The paper [9] used 2-D gamma distribution, the experiment showed that the proposed method obtained very good results but with a big time complexity due to the big number of constructed masks. To solve these problems, the study proposed a novel approach based on information theory, which is entropy-based thresholding. The proposed method is decrease the computation time. The results were very good compared with the well-known Sobel gradient [16] and Canny [17] gradient results.

The outline of the paper is as follows. In section 2, we have presented the classical edge detection methods that related to the paper. Image thresholding based on Shannon entropy is presented in section 3. Section 4, describes the proposed algorithm of edge detection. In section 5,we have presented the effectiveness of proposed algorithm in the case of real-world and synthetic images, is also, we compare the results of the algorithm against several leading edge detection methods. Conclusion and feature work are presented in Section 6.

## II. CLASSICAL EDGE DETECTION METHODS

Five most frequently used edge detection methods are used for comparison. These are: Gradient operators (Roberts, Prewitt, Sobel), Laplacian of Gaussian (LoG or Marr-Hildreth) and Gradient of Gaussian (Canny) edge detections [17, 18]. People which would like to read about this subject are referred to [19,20,21] evaluation studies of edge detection algorithms according to different criteria. The details of methods as follows:

## A. Roberts edge detector:

It was one of the first edge detectors and was initially proposed by Lawrence Roberts in 1963. It performs a simple, quick to compute, 2-D spatial gradient measurement on an image. It thus highlights regions of high spatial frequency which often correspond to edges [18]. The input to the operator is a grayscale image the same as to the output is the most common usage for this technique. Pixel values in every point in the output represent the estimated complete magnitude of the spatial gradient of the input image at that point, as shown in Figure 1.



Fig. 1. Roberts gradient estimation operator.

## B. Prewitt edge detector:

It based on the idea of central difference. It measures two components. The Prewitt edge detector is an appropriate way to estimate the magnitude and orientation of an edge. Although differential gradient edge detection needs a rather time consuming calculation to estimate the orientation from the magnitudes in the *x* and *y*-directions, the compass edge detection obtains the orientation directly from the kernel with the maximum response. The operator is limited to 8 possible orientations, however experience shows that most direct orientation estimates are not much more accurate. This gradient based edge detector is estimated in the 3×3 neighbourhood for eight directions as shown in Figure 2. All the eight convolution masks are calculated. One convolution mask is then selected, namely that with the largest module [18].



Fig. 2. Prewitt gradient estimation operator.

## C. Sobel edge detector:

The Sobel operators are named after Erwin Sobel. The Sobel operator relies on central difference, but gives greater weight to the central pixels when averaging. The Sobel operator can be thought of as 3×3 approximations to first derivative of Gaussian kernels. Sobel operators which are shown in the masks below (rotated by 90°):[18].



Fig. 3. Sobel gradient estimation operator.

## D. Laplacian of Gaussian Edge detection (LOG)

This LOG operator smoothes the image through convolution with Gaussian-shaped kernel followed by applying the Laplacian operator. Laplacian of Gaussian edge detection mask is:



Fig. 4. LOG gradient estimation operator.

## E. Canny edge detector:

The Canny edge detector is an edge detection operator that uses a multi-stage algorithm to detect a wide range of edges in images. It was developed by John F. Canny in 1986. Canny's aim was to discover the optimal edge detection algorithm. In this situation, an "optimal" edge detector means:

- Good detection – the algorithm should mark as many real edges in the image as possible.

- Good localization – edges marked should be as close as possible to the edge in the real image.

- Minimal response – a given edge in the image should only be marked once, and where possible, image noise should not create false edges.

The method can be summarized below:[22]

*1) The image is smoothed using a Gaussian filter with a specified standard deviation, to reduce noise.*

*2) The local gradient and edge direction are computed at each point using different operator.*

*3) Apply non-maximal or critical suppression to the gradient magnitude.*

*4) Apply threshold to the non-maximal suppression image.*

## III. SHANNON ENTROPY AND IMAGE THRESHOLDING

Entropy is a concept in information theory. It is used to measure the amount of information [23]. It is defined in terms of the probabilistic behavior of a source of information. In accordance with this definition, a random event $E$ that occurs with probability $P(E)$ :

$$I(E) = \log(1/P(E)) = -\log(P(E)) \tag{1}$$

The amount $I(E)$ is called information content of $E$. The amount of self information of the event is inversely related to its probability. If $P(E)=1$, then $I(E) = 0$ and no information is attributed to it. In this case, uncertainty associated with the event is zero. Thus, if the event always occurs, then no information would be transferred by communicating that the event has occurred. If $P(E) = 0.8$, then some information would be transferred by communicating that the event has occurred. The base of the logarithm determines the unit which is used to measure the information.

If the base of the logarithm is 2, then unit of information is bit. If $P(E)= ½$, then $I(E)= -log_2(½) = 1$ bit. That is, 1 bit is the amount of information conveyed when one of two possible equally likely events occurs. An example of such a situation is flipping a coin and communicating the result (Head or Tail) [24, 25]

The basic concept of entropy in information theory has to do with how much randomness is in a signal or in a random event. An alternative way to look at this is to talk about how much information is carried by the signal. Entropy is a measure of randomness. Consider a probabilistic experiment in which the output of a discrete source is observed during every unit of time (signaling interval). The source output is modeled as a discrete random variable $Z$. $Z$ is referred as a set of source symbols [26]. The set $Z$ of source symbols is referred to as the source alphabet, $Z= \{z_1, z_2, z_3, ..., z_k\}$.

The source symbol probabilities is $P= \{p_1, p_2, p_3, ..., p_k\}$. This set of probabilities must satisfy the condition $sum(p_i)=1$, $0 \le p_i \le 1$. The average information per source output, denoted $S(Z)$ [26], Shannon entropy may be described as:

$$S(Z) = -\sum_{i=1}^{k} p_i \, log(p_i) \qquad (2)$$

$k$ is the total number of symbols. If we consider that a system can be decomposed in two statistical independent subsystems $A$ and $B$, the Shannon entropy has the extensive property (additivity):

$$S(A+B) = S(A) + S(B) \qquad (3)$$

this formalism has been shown to be restricted to the Boltzmann-Gibbs-Shannon (BGS) statistics.

Let $f(x, y)$ be the gray value of the pixel located at the point $(x, y)$. In a digital image $\{f(x, y) \mid x \in \{1,2,...,M\}, y \in \{1,2,...,N\}\}$ of size $M{\times}N$, let the histogram be $h(a)$ for $a \in \{0,1,2,..., 255\}$ with $f$ as the amplitude (brightness) of the image at the real coordinate position $(x, y)$. For the sake of convenience, we denote the set of all gray levels $\{0,1,2,..., 255\}$ as $G$. Global threshold selection methods usually use the gray level histogram of the image. The optimal threshold $t^*$ is determined by optimizing a suitable criterion function obtained from the gray level distribution of the image and some other features of the image.

Let $t$ be a threshold value and $B = \{b_0, b_1\}$ be a pair of binary gray levels with $\{b_0, b_1\} \in G$. Typically $b_0$ and $b_1$ are taken to be 0 and 1, respectively. The result of thresholding an image function $f(x, y)$ at gray level $t$ is a binary function $f_t(x, y)$ such that $f_t(x, y) = b_0$ if $f_t(x, y) \le t$ otherwise, $f_t(x, y) = b_1$. In general, a thresholding method determines the value $t^*$ of $t$ based on a certain criterion function. If $t^*$ is determined solely from the gray level of each pixel, the thresholding method is point dependent [24, 25].

Let $p_i = p_1, p_2. . . p_k$ be the probability distribution for an image with $k$ gray-levels. From this distribution, we derive two probability distributions, one for the object (class $A$) and the other for the background (class $B$), given by:

$$p_A : \frac{p_1}{P_A}, \frac{p_2}{P_A}, \ldots, \frac{p_t}{P_A},$$
$$p_B : \frac{p_{t+1}}{P_B}, \frac{p_{t+2}}{P_B}, \ldots, \frac{p_k}{P_B} \qquad (4)$$

and where

$$P_A = \sum_{i=1}^{t} p_i, \qquad P_B = \sum_{i=t+1}^{k} p_i \qquad (5)$$

The Shannon entropy for each distribution is defined as:

$$S^A(t) = -\sum_{i=1}^{t} p_i \, log(p_i), \text{ and}$$
$$S^B(t) = -\sum_{i=t+1}^{k} p_i \, log(p_i) \qquad (6)$$

We try to maximize the information measure between the two classes (object and background). When $S(t)$ is maximized, the luminance level $t$ that maximizes the function is considered to be the optimum threshold value .

$$t^* = Arg \max_{t \in G}[S^A(t) + S^B(t)]. \qquad (7)$$

In the proposed scheme, first create a binary image by choosing a suitable threshold value using Shannon entropy. The *Threshold* procedure find the suitable threshold value $t^*$ for grayscale image $f$. It can now be described as follows:

---

**Procedure *Threshold*,**

**Input:** A grayscale image $f$ of size $m \times n$ with histogram $H$.

**Output:** $t^*$ of $f$.

Begin

Step 1: Let $f(x, y)$ be the original gray value of the pixel at the point $(x, y)$, $x=1..m$, $y=1..n$ .

Step 2: Calculate the probability distribution $0 \le p_i \le 255$ .

Step 3: For all $t \in \{0,1,…,255\}$,

    i. Calculate $p_A$, $p_B$, $P_A$, and $P_B$, using Eq.s (4 and 5).

    ii. Find optimum threshold value $t^*$, where $t^* = Arg \max_{t \in G}[S^A(t) + S^B(t)]$.

End.

---

## IV. THE PROPOSED MULTI-THRESHOLD ALGORITHM

This section presents the concept of object connectivity. It introduces a technique of edge detection based on entropy and geometric properties of the object. Geometric properties such as connectivity, projection, area, and perimeter are important components in binary image processing. An object in a binary image is a connected set of pixels. In what follows, we present some definitions related to connectivity of pixels in a binary image [25].

*Connected Pixels*: A pixel $f_0$ at $(i_0,j_0)$ is *connected* to another pixel $f_n$, at $(i_n,j_n)$ if and only if there exists a path from $f_0$ to $f_n$, which is a sequence of points $(i_0,j_0)$, $(i_1,j_1)$,…, $(i_n,j_n)$, such that the pixel at $(i_k,j_k)$ is a neighbor of the pixel at $(i_{k+1},j_{k+1})$ and $f_k = f_{k+1}$ for all, $0< k < n-1$.

*4-connected*: When a pixel at location $(i, j)$ has four immediate neighbors at $(i+1, j)$, $(i-1, j)$, $(i, j+1)$, and $(i, j-1)$, or four immediate neighbors at $(i+1, j+1)$, $(i-1, j+1)$, $(i+1, j-1)$, and $(i-1, j-1)$ they are known as, *4-connected*. Two four connected pixels share a common boundary as shown in Figure (5-a,5-b).

*8-connected*: When the pixel a t location $(i, j)$ has. in addition to above two types of four immediate neighbors, together, they are known as *8-connected*. Thus two pixels are eight neighbors if they share a common corner. This is shown in Figure (5-c).

*Connected component*: A set of connected pixels (4 or 8 connected) forms a *connected component*. Such a connected component represents an object in a scene as shown in Figure (5-d).



Fig. 5. (a) 4-connected, (b) Diagonal 4-connected, (c) 8-connected, and (d) Connected component.

In order to obtain edge detection, we find classification of all pixels that satisfy the criterion of homogeneousness, and detection of all pixels on the borders between different homogeneous areas. In the proposed scheme, first create a binary image by create a threshold value using Shannon entropy, using of the *Threshold* procedure. Region labeling in this system is done using 4-neighbor or 8-neighbor connectivity. A common alternative would be to use 4-neighbor connectivity instead (Figure 5).

The *Edge Detection* Procedure can be described as follows (using the 4-connected or diagonal 4-connected):

---

**Procedure *Edge Detection*;**

**Input:** A grayscale image $f$ of size $m{\times}n$ and $t^*$.

**Output:** The edge detection image $g$ of $f$.

Begin

Step 1: Create a binary image: For all $x$, $y$, If $f(x, y) \le t^*$ then $A(x, y) = 0$ Else $A(x, y) = 1$.

Step 2: Initialization of the output edge image of size $m{\times}n$, $g(x, y) = 0$ and for all $x$ and $y$.

Step 3: Checking for edge pixels:

For all $1< j< m$, and $1< i< n$ do

$$\lambda_1 = \left| A_{j,i} - A_{j,i-1} \right| + \left| A_{j,i} - A_{j,i+1} \right|, \lambda_2 = \left| A_{j,i} - A_{j-1,i} \right| + \left| A_{j,i} - A_{j+1,i} \right|,$$

$$\phi_1 = \left| A_{j,i} - A_{j-1,i-1} \right| + \left| A_{j,i} - A_{j+1,i+1} \right|, \quad \phi_2 = \left| A_{j,i} - A_{j-1,i+1} \right| + \left| A_{j,i} - A_{j+1,i-1} \right|,$$

If $\lambda_1 + \lambda_2 = 0$ or $\phi_1 + \phi_2 = 0$ then $g_{j,i} = 1$.

End For

End Procedure.



(a) single threshold value .



(b) Multi-threshold values .

Fig. 6. Histogram of test image and its multi-thresholds ($t_1$, $t_2$ and $t_3$).

The proposed Multi-Threshold Algorithm consists of the following steps:

**Algorithm Multi-Threshold;**

*1) Find the threshold value ($t_1$) using Threshold procedure based on Shannon entropy.*

*2) The histogram H of image with pixel values (0,1,2,...,255) is split by $t_1$ into two parts, $H_1$ pixel values (0,1,2,...,$t_1$) and $H_2$ with ($t_1$+1,...,255). See Figure 6-a.*

*3) Apply Threshold procedure with $H_1$ to find the threshold values ($t_2$) . then apply it with $H_2$ to find the threshold values ($t_3$). See Figure 6-b.*

*4) Create binary matrix A, using the three threshold values $t_1$, $t_2$ and $t_3$ according to the condition, For all $1 < j < m$, and $1 < i < n$ do: IF (($f(i,j) >= t_2$) and ($f(i,j) < t_1$)) or $f(i,j) >= t_3$) Then $A(i,j)=1$ else $A(i,j) = 0$.*

*5) Applying EdgeDetection procedure with A matrix to obtain the edge detection image g.*

End Algorithm.

## V. EXPERIMENTAL RESULTS AND DISCUSSION

In order to test the method proposed in this paper and compare with the other edge detectors, common gray level test images with different resolutions and sizes are detected by the proposed method, Gradient of Gaussian (Canny), Laplacian of Gaussian (LoG or Marr-Hildreth), Prewitt, Roberts and Sobel methods respectively.

The performance of the proposed scheme is evaluated through the simulation results using MATLAB. Prior to the application of this algorithm, no pre-processing was done on the tested images.

As the algorithm has two main phases – global and local enhancement phase of the threshold values and detection phase, we present the results of implementation on these images separately. Here, we have used in addition to the original gray level function $f(x, y)$, a function $g(x, y)$ that is the average gray level value in a 3×3 neighborhood around the pixel $(x, y)$.

Though the performance of the proposed entropic edge detector excels as a shape and detail detector, it is fraught with some drawbacks. It fails to provide all thinned edges. The weak edges are not eliminated but for some applications, these may be required.

This detector has another distinctive feature, i.e. it retains the texture of the original image. This feature can be utilized for the identification of fingerprints, where the ridges may have different intensities. We are experimenting on several images to come up with a useful selection guideline.



Fig. 7.   CPU time with 256×256 pixel test images or less

Fig. 8.   CPU time with 512×512 pixel test images



Fig. 9.   CPU time with 1024×1024 pixel test images

We run the previous methods and the proposed algorithm 10 times for each image with different sizes. As shown in Figures 7-10, the charts of the test images and the average of run time for the classical methods and proposed scheme. It has been observed that the proposed edge detector works effectively for different gray scale digital images as compare to the run time of Canny and LOG methods.

Image quality is a characteristic of an image that measures the perceived image degradation (typically, compared to an ideal or perfect image). Two parameters are there:

First, *MSE*, it is defined as the squared difference between the original image and estimated image.

$$MSE = \frac{1}{N} \sum_{i=1}^{N} \left( X - \hat{X} \right)^2$$

where $X$ = original value, $\hat{X}$ = stego value and $N$ = number of samples.

Second, *PSNR,* Peak Signal-to-Noise Ratio, often abbreviated PSNR, is an engineering term for the ratio between the maximum possible power of a signal and the power of corrupting noise that affects the fidelity of its representation [26]. Because many signals have a very wide dynamic range, PSNR is usually expressed in terms of the logarithmic decibel scale.



Fig. 10. The average of run time of proposed method ,Canny, log, Prewitt, Roberts and Sobel with different size test images

TABLE I.    AVERAGE MSE AND PSNR VALUES OF DIFFERENT EDGE DETECTION  METHODS ON TESTED IMAGES

| Method | Proposed Alg. | Canny | LOG | Prewitt | Roberts | Sobel |
|--------|---------------|-------|-----|---------|---------|-------|
| *MSE* | 0.0238 | 0.0200 | 0.1975 | 0.3086 | 0.3244 | 0.0278 |
| *PSNR* | 64.3448 | 68.1308 | 55.1745 | 53.2363 | 58.1951 | 63.6938 |

PSNR is most easily defined via the mean squared error (*MSE*):

$$PSNR = 10 \; \log_{10} \left( \frac{L^2}{MSE} \right) = 20 \; \log_{10} \frac{L}{\sqrt{MSE}}$$

where *L*= maximum value, *MSE* = Mean Square Error. See the Table 1.

Some selected results of edge detections for these test images using the classical methods and proposed scheme are shown in Figures 11-20.

From the results; it has again been observed that the proposed method works well as compare to the previous methods, LOG, Prewitt, Roberts and Sobel (with default parameters in MATLAB).



Fig. 11. Animal Port folio  Image



Fig. 12. Bacteria-Pili Image

Original image | log

Roberts | Sobel

Prewitt | Proposed Alg. T=( 114, 60, 175)

Fig. 13. Boat Image



Original Image | log

Roberts | Sobel

Prewitt | Proposed Alg. T=( 115, 54, 175)

Fig. 15. Gram-negative Bacterial Image



Original Image | log

Roberts | Sobel

Prewitt | Proposed Alg. T=( 67, 30, 133)

Fig. 14. Backbone Image



Original Image | log

Roberts | Sobel

Prewitt | Proposed Alg. T=(155, 91, 214)

Fig. 16. Zebra Image

Fig. 17. Rose Image



Fig. 19. tire Image



Fig. 18. Girl Image



Fig. 20. things Image

## VI. CONCLUSION AND FEATURE WORK

This paper shows the new algorithm based on the Shannon entropy for edge detection using histogram of the image. The objective is to find the best edge representation and minimize the computation time. A set of experiments in the domain of edge detection are presented. An edge detection performance is compared to the previous classic methods, such as, LOG, Prewitt, Roberts and Sobel. Analysis show that the effect of the proposed method is better than those methods in execution time, also is considered as easy implementation. The significance of this study lies in decreasing the computation time with generate suitable quality of edge detection. In this way entropic edge detector presented in this paper uses Shannon entropy with multi threshold values. It is already pointed out in the introduction that the traditional methods give rise to the exponential increment of computational time.

Experiment results have demonstrated that the proposed scheme for edge detection can be used for different gray level digital images. Another benefit comes from easy implementation of this method. An important future investigation will be the study of edge detection in the case of automatic target recognition, medical image applications and detection of video changes.

### REFERENCES

[1] Mohamed A. El-Sayed , "Study of Edge Detection Based On 2D Entropy", International Journal of Computer Science Issues (IJCSI) ISSN : 1694-0814 , Vol. 10, Issue 3, No 1, pp. 1-8, 2013.

[2] Mohamed A. El-Sayed , S. F.Bahgat , and S. Abdel-Khalek "Novel Approach of Edges Detection for Digital Images Based On Hybrid Types of Entropy", International Applied Mathematics & Information Sciences, Vol. 7, No. 5, pp.1809-1817, 2013.

[3] G.C. Anagnostopoulos, SVM-based target recognition from synthetic aperture radar images using target region outline descriptors. Nonlinear Anal.-Theor. Meth. App. 71, 12, e2934–e2939, 2009.

[4] M.T. Doelken, H. Stefan, E. Pauli, A. Stadlbauer, T. Struffert, T. Engelhorn, G. Richter, O. Ganslandt, A. Doerfler, T. Hammen, 1H-MRS profile in MRI positive- versus MRI negative patients with temporal lobe epilepsy. Seizure 17, 6, 490–497, 2008.

[5] Y.-T. Hsiao, C.-L. Chuang, Y.-L. Lu, J.-A. Jiang, Robust multiple objects tracking using image segmentation and trajectory estimation scheme in video frames. Image Vision Comput. 24,10, 1123–1136, 2006.

[6] S. Kresic-Juric, , D. Madej and S. Fadil. Applications of Hidden Markov Models in Bar Code Decoding. Intl. J. Patt. Recog. letters, 27, 1665-1672, 2006.

[7] G. Markus, , Essam A. EI-Kwae and R.K. Mansur. Edge detection in medical images using a genetic algorithm. IEEE Trans. on Medical Imaging, 17, 469-474, 1998.

[8] M. Wang and Y. Shuyuan, "A Hybrid Genetic Algorithm Based Edge Detection Method for SAR Image", In: IEEE Proceedings of the Radar Conference'05, pp. 1503-506, May 9-12, 2005.

[9] A. El-Zaart, "A Novel Method for Edge Detection Using 2 Dimensional Gamma Distribution", J. of Comput. Sc. 6, 2, pp. 199-204, 2010 .

[10] V. Aurich, and J. Weule, "Nonlinear Gaussian filters performing edge preserving diffusion. ", Proceeding of the 17th Deutsche Arbeitsgemeinschaft für Mustererkennung (DAGM) Symposium, Sept. 13-15, Bielefeld, Germany, Springer-Verlag, pp. 538-545, 1995.

[11] M. Basu, "A Gaussian derivative model for edge enhancement.", Patt. Recog., 27:1451-1461, 1994.

[12] G. Deng, and L.W. Cahill, "An adaptive Gaussian filter for noise reduction and edge detection.", Proceeding of the IEEE Nuclear Science Symposium and Medical Imaging Conference, Oct. 31-Nov. 6, IEEE Xplore Press, San Francisco, CA., USA, pp. 1615-1619, 1993.

[13] C. Kang, and W. Wang, "A novel edge detection method based on the maximizing objective function.", Patt. Recog., 40, pp. 609-618, 2007.

[14] Q. Zhu, "Efficient evaluations of edge connectivity and width uniformity.", Image Vis. Comput., 14, pp.21-34,1996.

[15] B. Mitra, "Gaussian Based Edge Detection Methods- A Survey ". IEEE Trans. on Systems, Manand Cybernetics , 32, pp. 252-260, 2002.

[16] R. C. Gonzalez, and R.E. Woods, "Digital Image Processing.", 3rd Edn., Prentice Hall, New Jersey, USA. ISBN: 9780131687288, pp. 954, 2008.

[17] J.F. Canny, "A computational approach to edge detection", IEEE Transactions on Pattern Analysis and Machine Intelligence (IEEE TPAMI), Vol. 8(6), pp. 769-798, 1986.

[18] N. Senthilkumaran and R. Rajesh, "Edge Detection Techniques for Image Segmentation - A Survey", Proceedings of the International Conference on Managing Next Generation Software Applications (MNGSA-08), 2008, pp.749-760.

[19] Bowyer K.W., Kranenburg C., Dougherty S. "Edge Detector Evaluation Using Empirical ROC Curves" IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 354-359, 1999.

[20] Heath M., Sarkar S., Sanocki T., Bowyer K.W. "A Robust Visual Method for Assessing the Relative Performance of Edge Detection Algorithms", IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), Vol.19(12), pp.1338-1359, 1997.

[21] M. Shin, D. Goldgof, K.W. Bowyer, "An Objective Comparison Methodology of Edge Detection Algorithms for Structure from Motion Task", IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 190-195, 1998.

[22] R. Deriche, " Using canny's criteria to derive a recursively implemented optimal edge detector", International Journal of Computer Vision (IJCV), Vol. 1(2), pp. 167-187, 1987.

[23] C.E. Shannon, A mathematical theory of communication. Int. J. Bell. Syst. Technical, 27, 379-423, 1948.

[24] B. Singh and A. P. Singh, " Edge Detection in Gray Level Images Based on the Shannon Entropy", J. Comput. Sci., 4 , 3, 186-191, 2008.

[25] Mohamed A. El-Sayed, Edges Detection Based On Renyi Entropy With Split/Merge. Computer Engineering And Intelligent Systems (CEIS) ISSN: 2222-2863. Vol. 3, No. 9, pp. 32-41 ,2012.

[26] S Jayaraman, S Esakkirajan and T Veerakumar, "Digital Image Processing," Tata McGraw Hill Education ptd. Ltd, New Delhi, 7th ed., 2012, pp.368-393.

# Dynamic Software Architecture for Medical Domain Using Pop Counts

UMESH BANODHA

Department of Computer Applications
Samrat Ashok Technological Institute
Vidisha (M.P.), INDIA

KANAK SAXENA

Department of Computer Applications
Samrat Ashok Technological Institute
Vidisha (M.P.), INDIA

*Abstract*—Over the past few decades, the complexity of software for almost any era has increased significantly. The aim of this paper is to provide an approach which not only feasible but also decision-oriented in medical era. It focus on the careful planning and organizing success in continuous process improvements in software and hardware technology as this brings a lot of trouble to system development and maintenance. We have used the pop count method to develop the dynamic software architecture with the existence of quality attributes in order to find out the level of severity in patients of any diseases on the specialist perception. This is useful for taking decision on priority healing and regular concentration of the patients even in the absence of the specialist. Further the method (model) tested on the 25 symptoms of 100 patients which does not contain any dichotomous data and found with the help of statistical evaluation (that it result almost perfect classification) that the architecture is conformance to the medical software architecture quality requirements.

*Keywords—Software Architecture; Quality Attributes; Pop count; medical process reengineering*

## I. INTRODUCTION

### A. Objective

The medical era consists of two broadly classified data (1) Process Data and (2) Application data. Process data is actually the data which managed by medical system where as the application data is transformed across task units.

In fact, the great efforts is the implementation of good human factors practices in the design of software specially in the domain of medicine either is its concerned with the devices, diagnosis, usage, treatment etc. However, the survey finding indicates that lack of attention to human factors during software development may lead to errors that have the potential for patient's health or even death [11]. Thus, the design principles and participation of human factors especially in Indian scenario are very important, in order to increase the patients safety and authentication of the software in the medical era. Today, also the perception of specialists' doctors in remote areas is not up to the mark for the acceptability of the software in the decision / diagnosis and treatment.

For the above problem, Software Architecture is being widely used to describe a very high level design of large software systems. "The software Architecture of a program or computing system is the structure or structures of the system, which comprise software elements, the externally visible

properties of those elements, and the relationships among them"[1]. Really it will be interesting to find out where software fits in with the software development life cycle especially in medical domain.

Architecture is the structure of the components of a program or system, their interrelationships and the principles and guidelines governing their design and evolution over time. [1]

The software architecture of a program or computer system is the structure or structures of the system, which comprise software components, the externally visible properties of those components, and the relationships among them. [2]

The importance of software architecture arose at very first step during system development especially in medical domain as they virtually affect every later stages of the development process which will direct or indirect impact the mortality rate of the patient at the end. Thus, good software architecture can reduce the risk with building a technical solution and make the system implementation and testing more traceable as well as achieve higher quality attributes. [3]

### B. Need of Software Architecture in Medical Domain

Using software architecture we can present a common abstraction of a medical decision system that most if not all of the system's user can use as a basis for mutual understanding, negotiations, consensus and communication.

Medical architecture manifests the earliest design decisions about a system, and these early binding carry weights far out of proportion to their individual gravity with respect to the system's remaining development, its deployment and maintenance. The structure defines constraints on implementation. It explains the organizational structure and predicts the qualities of system. The architecture makes it easier to reason about and manage changes. It helps in evolutionary prototyping. The architecture enables more accurate cost and schedule estimates.

Software architecture constitutes a relatively small, intellectually graspable medical model for how a system is structured and how its elements work together, and this model is transferable across system. [1]

### C. Medical Process modeling approach

Medical process is defined as the art of healing, i.e., a gradual process of medicine tending to cure. It is a method that

helps to understand the actions, work flow, and tasks of an organization, and how the tasks are executed. The focus in process modeling is on the functional processes which are entities that start with an initial event and end with a result. A process has always an input and an output, input triggers the process and process results in an output [5, 6].

The process consists of four steps (Figure1) in the highest abstraction level. Process begins when the patient arrives to the reception of a hospital/clinic to meet doctor /staff as an initial event. It ends when the patient is discharged as a result. The actor of the processes is a doctor with support unless mentioned otherwise [4, 7]. In this paper we focused on phase 1 and 2 of Process modeling approach i.e. Arrival/ First assessment of patient and planning the care [8, 9, 10].

### D. Requirements of medical domain

As we know that the medical system requires high performance then we can study the behavior of the elements (components) in time frame, the frequency and volume of inter-element communication. The medical model requires the modifiability then we can assign the responsibility to elements so appropriate change done in the whole system efficiently and abruptly.

The medical model need highly secure system; we need to provide protection of each elements and inter-element communication. We also need to provide the access power of each element so that authenticate user can observer the specialized elements of the system.

The model should provide the scalability. The model should provide localize use of resources so all patients are get benefited. It also concerned with higher capacity replacements. The model is combination of components and connector so it represents the incremental subsets of the system and manages inter-component usage. The model should be reusable i.e. each component or sub-components of the system be reusable. We restrict inter-element coupling so that we extract an element it does not come out with too many attachments to its current environment to be useful [1].



Fig. 1.   Process Modeling Approach

## II.   PROPOSED METHOD

To make progress, we focus on the patient's symptoms caused by a specific set of medical conditions; and to generalize as the symptoms of patients of any diseases and the specialist will divide the set of their order. Thus, this paper gives the method which not only useful to the specialist but also a non-specialist i.e. moderate user can also use. For this we have work out on 100 patients' data with 25 attributes as symptoms and found the result in much improved manner.

Our method works on:

- We use the binary values as 0 if that symptom is not present otherwise 1.

- The medical specialist according to the diseases will decide a binary code of 25 values of 25 symptoms. These can be extended to as many as values required for the analysis.

- We proposed the dynamic software architecture using the pop count methods to estimate the level of severity among the set of the patients. Thus, in the absence of the specialist a moderate use can work out. The concept of this paper is to predict the level of severity with pop count method as the only sorting method cannot give the fruitful result

### A. Pop count Method

We have the data of the patients in the form of symptoms sets which range from any number of symptoms to any number in any order except the data of patients should arrange in the same order $P = \{P_1, P_2, P_3,…,P_n\}$ here P is the set of the patients,

$P_i = \{S_1, S_2, S_3,…,S_m\}$ here $P_i$ is the set of individual patient's symptoms with $S_1, S_2, S_3,…,S_m$ symptoms.

Thus $\{P = S_{ij}$ of all symptoms values of $i^{th}$ patients' $j^{th}$ symptoms$\}$

Now, in order to find out the patients' with critical condition we use the method of pop-count with count leading one's operation. We use the function as $f(S_{ij}) = pop (S_{ij} \wedge (\sim(-S_{ij})))$ where '$\wedge$' denotes the XOR bitwise operations. After performing the above function on the set P we indexed the entire data in order to find the list of critical data. Picking up any threshold value perform the preprocessing of the data and categories the data among various clusters in the styles (various types of styles) for various clusters depends upon the values ( standards) provided by the specialists' perception.

Figure 2 shown step by step activity of the pop-count method. We can observe that there are two decision steps by which we can dynamically observe the critical symptoms and take appropriate assessment.

### B. Software Architecture Evaluation

The proposed software architecture will faster the treatment of serious patients. The method can be treated as the base for the software development which can have the following charactrestics. It can be of high performance, secure with less risk and due to its generality nature. It can be used for reusability. One level it gives the fast treatment with effectiveness on health point of view and on the other hand it may also affect medical spending. Though, we cannot do a complete evaluation of the impact of these changes, but we

provide some information. In order to predict the level of severity in general we estimate logistic regression model as a function of severity on the presence of priority symptoms as

$$F(P_i) = e^{\beta_0 + \beta_1 s_{ij}}$$

i.e. the function of patients' equaling a "success" rather than a failure, where $\beta_0$ is the intercept and $\beta_1$ is the regression coefficient multiplied by the symptoms' success ( $S_{ij}$ ) and base e denotes the exponential function.

$$L(\beta_0, \beta_1) = \prod_{i=1}^{n} p(s_{ij})^{C_i}(1-p(s_{ij})^{1-c_i})$$

On the basis of clusters on symptoms, we applied the weka software tool for predicting the assessment regarding patients' symptoms. The evaluation of model is done on kappa statics, mean absolute error, root mean square error, root Absolute error and root relative square error.

*1) The Kappa Statistic:* Interobserver variation can be measured in any situation in which two or more independent observers are evaluating the same thing. The calculation is based on the difference between how much agreement is actually present ("observed" agreement) compared to how much agreement would be expected to be present by chance alone ("expected" agreement).

*2) Mean absolute error (MAE):* The MAE measures the average magnitude of the errors in a set of forecasts, without considering their direction. The mean absolute error is a quantity used to measure how close forecasts or predictions are to the eventual outcomes

*3) Root mean squared error (RMSE):* The difference between forecast and corresponding observed values are each squared and then averaged over the sample. The RMSE is most useful when large errors are particularly undesirable. They are negatively-oriented scores: Lower values are better.

*4) Relative Absolute Error (RAE):* The relative absolute error is very similar to the relative squared error in the sense that it is also relative to a simple predictor, which is just the average of the actual values.

*5) Root Relative Absolute Error (RAE):* By taking the square root of the relative squared error one reduces the error to the same dimensions as the quantity being predicted.

## III. ANALYSIS

The evaluation result of the proposed modal worked out on the 100 patients with 25 symptoms of any diseases. The result as per table 1 shows a perfect performance on the training data set and correctly classified all the data especially of class A, B that is no false negative and false positive in the report. The analysis reported as kappa with 1 i.e. perfect agreement in class A, B and some extent class C is also almost perfect where as in class D & F it is substantial and class E an having moderate. Alternatively, we have used the relative errors to strengthen the result.



Fig. 2.  Flow chart of method using pop-count

It is discovered that in class A and B equal to zero and rest it range from 0.01 to 0.05. In order to strengthen the software architecture evaluation we used MAE which ranges from 0.0 to 0.03 and other error rates. Thus, a model has a lower error rate will be preferred as it has more powerful classes capability and ability in terms of medical and bio-informatics fields. Statistical analysis shows that classes A and B are the perfect. The graphs are the pictorial representation of the various results.

Our proposed model is capable for effectively handle the patients data both on architectural and non-architectural aspect. The resultant is of high performance, even if the specialists change the order of symptoms and hence reduce the ambiguity and confusion in terms of risk. Thus, we can say that the model is having the quality attributes which are required in the architectural evaluation process maintainability, usability, performance, reliability, reusability and availability.

TABLE I.        STATISTICS OF VARIOUS CLUSTERS

|        | F     | A    | B    | C     | D     | E     |
|--------|-------|------|------|-------|-------|-------|
| Kappa  | 0.89  | 1    | 1    | 0.95  | 0.84  | 0.72  |
| MAE    | 0.02  | 0    | 0    | 0.00  | 0.012 | 0.03  |
| RMSE   | 0.03  | 0    | 0    | 0.05  | 0.078 | 0.12  |
| RAE    | 10.10 | 0.00 | 0    | 5.02  | 16.00 | 27.90 |
| RRSE   | 31.78 | 0.03 | 0.01 | 22.41 | 40    | 52.84 |



Graph I        Statistics between clusters and kappa



Graph II        Statistics between clusters, MAE and RMSE



Graph III        Statistics between clusters, RAE and RRSE

IV.        CONCLUSIONS

The paper deals with the analysis of model with the impact of software architecture based on 100 records (data sets) of general patients. As a conclusion, we have proposed a model which is to evaluate and investigate the level of severity of patients based on the pop-count. This will help not only the specialist but also the dichotomous and moderate users to decide about the useless information; useful information about a supposed illness on the evidence suggests that improvements in medical care with higher use experienced considerable reductions in mortality, even with the existence of human factors. The proposed software architecture uses all the patients with equal priorities and symptoms of which the order of their availability and non-availability are to be set by the specialist. This also use the quality attributes among them is of interdependencies.

REFERENCES

[1] Len Bass,Paul Clements, Rick Kazman, "Software Architecture in Practice", Pearson education, Second edition", 2005.

[2] Banani Roy and T.C. Nicholas Graham, "Method for Evaluating Software Architecture: A Survey", Technical Report 2008-545, School of computing queen's university at Kingston Ontario, Canada, April, 2008.

[3] Qiushi Wang, Zhao Yang, "A Method of selecting appropriate software architecture style/pattern: Quality Attributes & Analytic Hiererchy Process", University of Gothenburg, Chalmers University of Technology, Göteborg, Sweden, June 2012

[4] P. Nykanen and J. Makinem, "Integration of medication information in electronics patient record systems", The dementia patient case. Turku School of economics, Research report LTH-1:2007, Turku, 2007.

[5] Wang, "Modelling information architecture for the organization. Information and Management" 32, 6, 1997, pp. 303-315.

[6] R. Weber, "Conceptual modeling and ontology: Possibilities and pitfalls", Journal of Database Management,14, 2, 2003, pp. 1-20.

[7] Makinen,J. Et.Al., "Process Models of medication Information", Procd. Of the 42$^{nd}$ Hawaii International Conf. on System sciences 2009, 1-7.

[8] U.Banodha , K.saxena, "Impact of Pipe and Filter Style on Medical Process Re-engineering", International Journal of Engineering Sciences, October 2011.

[9] U.Banodha , K.saxena,"Usability of Software Architecture design pattern in Medical process reengineering model", International Journal of Application or Innovation in Engineering & Management (IJAIEM), Volume 2, Issue 6, June 2013, ISSN 2319 – 4847.

[10] U.Banodha, K.Saxena, "Comparison of Software Architecture Styles in Medical Process Re-engineering Model", International Journal of Wisdom Based Computing, Vol. 2(1), April 2012.

[11] http://www.ncbi.nlm.nih.gov/pubmed/15306163

AUTHORS PROFILE

**UMESH BANODHA**, Assistant Professor at Samrat Ashok Technological Institute, VIDISHA (M.P.), an Autonomous Institute, affiliated to Rajiv Gandhi Technical University, Bhopal. I did MCA, M.Tech (Honors) and Pursuing Ph.D. My Area of interest Software Engineering / Architecture, Databases, UML, object oriented, Programming Languages etc. I am member of various international / National journals. I published more than 12 research papers in various conferences and journals (National / International).

**KANAK SAXENA**, Ph. D. in computer Science from the Devi Ahilya University, Indore, INDIA. She is professor in the Computer Applications Department at the Samrat Ashok Technological Institute affiliated to Rajiv Gandhi Technical University, Bhopal. Her Current research focuses on Database Systems, Parallel computing, Data Uncertainty and design and other interests include Network security and performance and Software Engineering. She is the member of editorial board of various international journals. She is the member of the international committee of the International Conference on Computer Science and Its Applications. She Published more than 80 research Papers in Various Conferences and Journals National / International).

# OLAWSDS:
# An Online Arabic Web Spam Detection System

Mohammed N. Al-Kabi

Faculty of Sciences & IT
Zarqa University
Zarqa, Jordan

Heider A. Wahsheh

Computer Science Department
College of Computer Science
King Khalid University
Abha, Saudi Arabia

Izzat M. Alsmadi

Information Systems Department
College of Computer & Information
Sciences Prince Sultan University
Riyadh 11586, P. O. Box 66833,
Saudi Arabia

*Abstract*—For marketing purposes, Some Websites designers and administrators use illegal Search Engine Optimization (SEO) techniques to optimize the ranking of their Web pages and mislead the search engines. Some Arabic Web pages use both content and link features, to increase artificially the rank of their Web pages in the Search Engine Results Pages (SERPs).

This study represents an enhancement to previous work in this field. It includes the design and implementation of an online Arabic Web spam detection system, based on algorithms and mathematical foundations, which can detect the Arabic content and link web spam depending on the tree of the spam detection conditions, beside depending on the user's feedback through a custom Web browser. The users can participate in making the decision about any Web page, through their feedbacks, so they judge if the Arabic Web pages in the browser are relevant for their particular queries or not. The proposed system uses the extracted content and link features from Arabic Web pages to determine whether to label each Web page as a spam or as a non-spam. This system also attempts to learn from the user's feedback to enhance automatically its performance.

Statistical analysis is adopted in this study to evaluate the proposed system. Statistical Package for the Social Sciences (SPSS) software is used to evaluate this new system which considers the users feedbacks as dependent variables, while Arabic content and links features on the other hand are considered independent variables. The statistical analysis with the SPSS is used to apply a variety of tests, such as the test of the analysis of variance (*ANOVA*). *ANOVA* is used to show the relationships between the dependent and independent variables in the dataset, which leads to solving problems and building intelligent decisions and results.

*Keywords—Arabic Web spam; content-based; link-based; Information Retrieval*

## I. INTRODUCTION

Arab Internet users suffer from two problems, the first problem is the low percentage of the Internet Arabic content, and the second problem is Arabic Web spam which leads Web search engines to refer to irrelevant Web pages. The success of spamming techniques to deceive a search engine leads the Internet users to lose credibility in the search engine they used, in addition to some other negative aspects of spamming such as wasting the time and efforts of the search engine users.

This study proposes an integrated system to reduce the Arabic content and link Web spam, and filter the search engines from these malicious Arabic web pages. Although this study relies on a set of content and link Arabic Web spam conditions that have been used before, however this study differs from its predecessors by involving the Web search engine users to assess the relevancy of Arabic Web pages rendered by Search Engine Results Pages (SERPs).

The proposed system allows users to use a synchronization technique, in which the users can browse the Arabic Web pages, and give their feedbacks assessment for each visited Web page under some security considerations and confidentiality. The use of a synchronization technique helps the proposed system to ensure that the submitted assessment is conducted by users not agents and robots.

The evaluation of the results of the proposed system is based on the use of Statistical Package for the Social Sciences (SPSS) software, which enables us to conduct a statistical analysis, and confidence predictive method. SPSS software considers Arabic Web spam features as independent variables, while it considers the Search Engine Ranking (SER), TrustRank, and link popularity scores as dependent variables. The statistical analysis in SPSS applies a variety of tests, such as the test of the analysis of variance (*ANOVA*). *ANOVA* has two types (one-way and two-way analysis of variance). In this study we used two-way analysis of variance to show the relationships between the dependent and independent variables in the dataset.

The main aim of this research is the development of a system which can filter the search engines from unwanted and spam Web pages based on the Web pages' features and the users which have a main role in determining the relevancy of SERPs with their different queries.

The rest of the paper is divided as follows: Section two presents selected related work of Web spam studies. Section three presents developed system overview. Section four elaborates experiments and results. Section five summaries the paper and its contribution.

## II. RELATED WORKS

The literature is rich with many studies related to Web spam, where this topic is studied from different perspectives. This section presents few of these studies which are closely

related to the subject of this paper: Detection of Web spam, and those studies dedicated to the evaluation of the correlation between spam and the trust.

The authors of this study enhanced the previous study of [1], which built Arabic content/link Web spam detection system. The study of [1], collected a large data set of Arabic Web pages (spam and non-spam), where various number of content and link based features extracted. Arabic content/link Web spam detection system based on the tree of the decision tree machine learning algorithm to build the rules of the proposed system, which yields the accuracy of 90.10% for Arabic content based, 93.10% for link based, and 89.01% in detecting both Arabic content and link Web spam detection system.

Content trust is essential to determine the quality of Web content, and a hot topic of research. The task of determining trustworthy information from inaccurate or untrustworthy information is becoming a hard task. The study of [2] shows how to adopt content trust to detect Web spam and rank each Web page accordingly. Text feature attributes, and information quality are used their novel content trust learning algorithm. They also developed a system to detect Web spam which shows its effectiveness relative to other alternative ways.

The study of [3] presents a new ranking algorithm for Web search engines which capable to eliminate spam Web pages from their results. A small blacklist of classified spam web pages is used by their algorithm. This ranking algorithm is based in its identification of spam on two aspects tendency and authority. Therefore a high quality Web pages with low or no spam tendency will get high ranks by their ranking algorithm. They conducted tests which show the effectiveness of this ranking algorithm relative to PageRank algorithm.

The paper of [4] studied the feedback of the users and converted it to the query log. For each user, a query log file was assigned. This log file contains: query words, document returned to the search engine, Web documents that users triggered within clicked date and time, and the rank of retrieved documents. The researchers applied two approaches: Web spam detection, and query spam detection. Web spam detection removes spam link and content features from the query log graphs, while query spam detection eliminates all queries that gain a high number of spam Web pages.

In [5] a language model approach was proposed, which extracted a combination of content-based and link-based features from two popular spam datasets (Webspam-UK2006 and Webspam-UK2007). Kullback-Leibler (KL) divergence was applied on the spam Web pages to characterize the relation between the two linked Web pages. The proposed model has improved the F-measure of Webspam-UK2006, and Webspam-UK2007 to about 6% and 2% respectively.

The study of [6] presented the influence of cloaking techniques to increase the rank of Web pages. Lin study proposed three techniques: TagDiff2, TagDiff3, and TagDiff4 to determine if the URLs are cloaked [6]. The proposed techniques are based on discovering differences in the copies

(HTML tags) of a specified Web page when it is sent to the Web crawler and to the Web browser. The conducted tests showed that tag-based methods exceed the link-based and content-based results in precision and recall. The Decision tree J48 uses the integration of content-based and tag features to yield an accuracy of 90.48%.

The study of [7] proposed a new methodology to detect spam Web pages based on the Qualified-Link (QL) analysis, and content-based features with the language-model (LM). Kullback-Leibler (KL) divergence was applied on the spam Web pages to find the relation between two linked Web pages based on both the content-based and link-based features. An automatic classifier was built to combine QL and LM features. The conducted results were applied on WEBSPAM-UK2006, and WEBSPAM-UK2007 datasets and showed an accuracy of 89.4% and 54.2% respectively.

Cloaking is a known Web spam technique which is used to deceive Web search engines, where the content of a Web page presented to Web search engine crawlers is different from the content of a web page provide to a Web browser. Therefore the study of [8] presents three tag-based methods to identify cloaked URLs. The effectiveness of their methods is compared against the effectiveness of term- and link-based methods, and the results prove that these three tag-based methods are more effective than term- and link-based methods. Also he presents in his paper a taxonomy to classify various cloaking detection methods. Lin study described and discussed dynamic cloaking.

The combined usage of trust and distrust propagations by semi-automatic anti-spam algorithms proved its effectiveness, but little work is done in this field. Therefore the study of [9] presents a framework to assign for each Web page a GoodRankscore (trustworthy score) and BadRank score (untrustworthy score). Afterward they propose a novel Good-Bad Rank (GBR) algorithm, where the propagation of a page's trust/distrust is based on probability of the Web page being trusted/distrusted. Tests conducted by those researchers show the effectiveness relative to link-based anti-spam algorithms that propagates only trust or distrust.

## III. SYSTEM OVERVIEW

This study aims to develop and improve the techniques used in [1], and proposed new system called Online Arabic Web Spam Detection System (*OLAWSDS*).

The authors of [1] built an Arabic content/link Web Spam Detection System, which mainly consists of the following main parts:

*1) Built in Web crawler: The role of the crawler to automate fetching the content of Arabic Web pages.*

*2) Arabic Web Spam data collections: It contains 23,000 labeled Arabic spam and non-spam Web pages.*

*3) Arabic content/link Web pages Analyzer: This is a customized tool that analyzes and measures content and links of Arabic Web pages features, in order to evaluate their optimized features. Table 1 summarizes the main extracted content/link features.*

TABLE I.  ARABIC CONTENT/LINK WEB SPAM FEATURES [1].

| Arabic Content Web spam features | | Arabic Link Web spam features | |
|---|---|---|---|
| 1. | Meaningless keyword stuffing (Arabic/English/Symbol) (in Web pages, Meta tags). | 2. | Number of image links. |
| 3. | Compression ratio for Web pages. | 4. | Number of internal links. |
| 5. | Number of images. | 6. | Number of external links. |
| 7. | Average length of Arabic/English words inside the Web pages. | 8. | Number of redirected links. |
| 9. | URL lengths. | 10. | Number of empty link text. |
| 11. | Size of compression ratio (in Kilobytes). | 12. | Number of empty links. |
| 13. | Web page size (in Kilobytes). | 14. | Number of broken links (which refers to null destinations). |
| 15. | The maximum Arabic/English word length. | 16. | The total number of links (the internal and external). |
| 17. | Size of hidden text (in Kilobytes). | | |
| 18. | Number of Arabic/English words inside <Title tag>. | | |

The authors of [1] labeled the Web pages as either spam or non-spam pages in the data collections depending on their judgments and previous Arabic/ non-Arabic Web spam studies.

In this study we improved the Arabic Web Spam Detection System, by computing more derived features and benefits of the client/server model.

*OLAWSDS* computes the following derived features:

- Search Engine Ranking score.
- TrustRank score.
- Link popularity score.

Client/server model is a distributed system architecture, which divides the work between the server that provides the hosting of the services, and the clients which request the services. Our proposed *OLAWSDS* used the client/server model, by considering the improved Arabic Web Spam Detection System as a server. *OLAWSDS* built a custom Web browser used to explore the Web pages through the clients' computers, and the clients used it as judgments area by including their decisions, which send to the server. Figure 1 presents the main parts of *OLAWSDS* and the flow of work.



Fig. 1.  OLAWSDS Architecture

In the proposed *OLAWSDS* system, we used the Arabic Web spam dataset of [1] as a black list dataset. Every user in client using *OLAWSDS* can browse the Web pages, and check and identify any spammed Web page(s).

The clients send their feedbacks to *OLAWSDS*, which reside in the server. *OLAWSDS* considered the valid users decision when their metrics exceed the thresholds. The thresholds depend on how many users send their decisions for particular Web page(s). *OLAWSDS* system saves the client server IP address and computes the Web spam features including the user's decision features. *OLAWSDS* sends the user final decision about a particular Web page either as a spam or non-spam, then update the Arabic Web spam black list dataset.

Visual and audio code is used by *OLAWSDS* to avoid spammers and robots from getting into this system, and make a fake decision of the type of Web pages. To avoid the participant of the spammers or the robots as user's clients, *OLAWSDS* requests from any client before sending the decision to fill the visual or audio code. Figure 2 presents an example of used visual and audio code.

Fig. 2.   Example of used Visual and audio code

*OLAWSDS* also saves the client IP address and prevent any decision from the same user for the same Web pages for one day.

## IV.   EXPERIMENT AND RESULTS

In this study we used the SPSS software to evaluate our proposed *OLAWSDS*. Analysis of Variance (*ANOVA*) was computed for 10,000 Arabic Web pages.

Statistical Package for the Social Sciences (SPSS) is a Statistical Software Package acquired by IBM in 2009, used for data mining, text analytics, and statistical analysis. SPSS divides the variable in the dataset to the main two types; dependent and independent variables, then applies the regression analysis [10].

- Dependent variables: It is defines as the results, operational outputs of the independent variables [11].

In this study, our proposed *OLAWSDS* computes the search engine ranking score. Which considered as dependent variables, where the computed score depends on the content and link features of the collected Web pages.

*OLAWSDS* also computes the TrustRank score, which is a link analysis technique which identify useful Web pages which are linked in most cases to other good Web pages, while spam Web pages point to the spam Web pages [12]. As search engine ranking score, the TrustRank depends on the many features of content and links features it is considered as a dependent variable.

*OLAWSDS* computes the link popularity score based on both external and internal links, so link popularity score considered as dependent variable.

- Independent variables: It is defined as the inputs of the independent variables (the values that determine the values of other variables) [11]. All Arabic web spam features are considered as independent variables.

*ANOVA* has two main types (one-way and two-way analysis of variance). It is used to show the relationships between the dependent and independent variables in the dataset. In this study we used *ANOVA* with two-way analysis of variance to determine the effect of independent variables on two or more continuous dependent variables [11].

The null hypothesis is that there is no relationship between two measured variables, or that the independent variable has no effect on the dependent variables (i.e. means are same). The alternative hypothesis states that there is relationship and effect between two measured variables (i.e. means are different). The goal of *ANOVA* is to accept or reject the null hypothesis [11].

We applied two-way analysis of variance on the three groups of the independent variables which affects the dependent variables (search engine ranking score, TrustRank, and link popularity).

Table 2 presents the results of two-way analysis of variance on the independent variables which affects search engine ranking score. Where the number of independent variables belong to table 1.

$f = F$ test; which is used to compare the statistical models to find the best fit population from which the data were sampled.

Degrees of freedom (*df*); if there are *N* observations in total, $df$total $= N - 1$.

*P*-value; It is a value used to evaluate the statistical standards, when *P*-value$< 0.05$, we reject the null hypothesis [11].

Kappa statistic (*KS*); computes the percentage of error reduction compared to all errors in the classification sample. When there is no agreement between two raters *KS* is zero or close to zero, and when *KS* value is close to 1 this mean we have a perfect agreement between two raters [1].

Mean absolute error (*MAE*); measures the average of the errors in a set of the estimation, to show how much the estimation relatives to the actual outcomes [1].

Root Mean Squared Error (*RMSE*); measures the average of the errors, through measuring the difference between estimates and corresponding observed values. The range of *RMSE* from 0 to ∞; low values are better than high values [1].

TABLE II.          ANOVA RESULTS FOR SEARCH ENGINE RANKING SCORE
MODEL

| *Independent Variable Number* | *df* | *f* | *P-value* | *KS* | *MAE* | *RMSE* | *Accuracy* |
|---|---|---|---|---|---|---|---|
| 1-7, 9, 11, 13, 15, 17, and 18 | 12 | 278.11 | 0.006 | 0.92 | 0.03 | 0.19 | 0.96 |

According to table 2, we can find that *f* test is 278.1157 with *P*-value equals to 0.006751, which is less than significant value of *P*-value (0.05), so we reject the null hypothesis, and accept the alternative hypothesis which asserts the relationship between the fifteen content and link Web spam features and the search engine ranking score. This model yields an accuracy of 96%. Other performance measurements computed include; *KS, MAE,* and *RMSE* were close to optimal values.

Table 3 shows *ANOVA* results of the independent variables which affect TrustRank scores, with the same statistical measurements, that used before.

TABLE III.        ANOVA RESULTS FOR TRUSTRANK MODEL

| Independent Variable Number | df | f | P-value | KS | MAE | RMSE | Accuracy |
|---|---|---|---|---|---|---|---|
| 1-7, 9-12, 13, 15, 17, and 18 | 14 | 5467.138 | 0.000 | 0.99 | 0.007 | 0.037 | 0.99 |

Table 3 shows that *f* test is equal to 5467.138 with *P*-value is equal to 0.000, which leads to the rejection of the null hypothesis, and accept the alternative hypothesis. The TrustRank model yields an accuracy of 99%.

Table 3 shows improvement results of *KS* and other performance measurements than the results of table 2, which are very close to the optimal values. Table 4 shows the *ANOVA* results for the link popularity model.

TABLE IV.        ANOVA RESULTS FOR THE LINK POPULARITY MODEL

| Independent Variable Number | df | f | P-value | KS | MAE | RMSE | Accuracy |
|---|---|---|---|---|---|---|---|
| 2, 4, 6, 8, 10, 12, 14, and 16 | 7 | 2205.75 | 0.000 | 0.77 | 0.14 | 0.28 | 0.88 |

According to Table 4, *f* test is equal to 2205.754 with *P*-value is equal to0.000, so the null hypothesis has to be rejected. The link popularity model yields an accuracy of 88%. *KS* yields an accuracy of 0.77 which considered a significant value, since it is close to 1.

The performance measurement; *MAE* and *RMSE* yields accepted values (close to zero). The comparison of the results of the three previous tables reveal that TrustRank and search engine ranking score models yields nearly the same results, followed by the link popularity model.

## V.    CONCLUSION

Website masters and developers struggle to improve their Websites' visibility; such actions help to increase the value of the search engine ranking score, trust rank, and link popularity and give them better opportunities in e-commerce marketing and advertisement campaigns.

In this paper, we proposed an Online Arabic Web Spam Detection System (*OLAWSDS*), which uses features extracted from content and links and benefits from client/server models. SPSS package is used to evaluate our proposed system using the test of the analysis of variance (*ANOVA*). Results showed improved results in comparison with other approaches in terms of prediction accuracy or performance.

REFERENCES

[1]  H. A. Wahsheh, M. N. Al-Kabi, and I. M. Alsmadi, "A link and Content Hybrid Approach for Arabic Web Spam Detection, "International Journal of Intelligent Systems and Applications (IJISA), vol. 5, no. 1, pp. 30-43, 2013.

[2]  W. Wang, G. Zeng, D. Tang, "Using evidence based content trust model for spam detection, "Expert Systems with Applications, vol. 37, pp. 5599-5606, 2010.

[3]  H. Wang, Y. Lia, K. Guo, "Countering Web Spam of Link-based Ranking Based on Link Analysis," Procedia Engineering, vol. 23, pp. 310–315, 2011.

[4]  C. Castillo, C. Corsi, D. Donato, "Query-log mining for detecting spam," Proceedings of the 4th international workshop on Adversarial information retrieval on the Web Pages AIRWeb '08, ACM, pp. 17-20, 2008.

[5]  J. Martinez-Romo, L. Araujo, "Web spam Identification Through Language Model Analysis," Fifth International Workshop on Adversarial Information Retrieval on the Web AIRWeb '09, Madrid, Spain, pp. 21-28, 2009.

[6]  J. Lin, "Detection of cloaked Web spam by using tag-based methods," Expert Systems with Applications, vol. 36, pp. 7493-7499, 2009.

[7]  L. Araujo, J. Martinez-Romo, "Web spam Detection: New Classification Features Based on Qualified Link Analysis and Language Models," IEEE Transactions on Information Forensics and Security, vol. 5, pp. 581-590, 2010.

[8]  J. Lin, "Detection of cloaked web spam by using tag-based methods," Expert Systems with Applications, vol. 36, pp. 7493–7499, 2009.

[9]  X. Liu, Y. Wang, S. Zhu, H. Lin, "Combating Web spam through trust-distrust propagation with confidence," Pattern Recognition Letters, vol. 34, no. 13, pp. 1462-1469, 2013.

[10] SPSS software, Retrieved October, 18, 2013, fromhttp://www-01.ibm.com/software/analytics/spss/

[11] R. N. Cardinal, "Graduate-level statistics for psychology and neuroscience ANOVA in practice, and complex ANOVA designs," 2004, fromhttp://egret.psychol.cam.ac.uk/psychology/graduate/Guide_to_ANOVA.pdf

[12] Z. Gyongyi, H. Garcia-Molina, and J. Pedersen, "Combating web spam with trustrank,"  VLDB '04 Proceedings of the Thirtieth international conference on Very large data bases, vol. 30, pp. 576-587, 2004

AUTHORS PROFILE

Mohammed Naji Al-Kabi obtained his Ph.D. degree in Mathematics from the University of Lodz/Poland (2001), his masters degree in Computer Science from the University of Baghdad/Iraq (1989), and his bachelor degree in statistics from the University of Baghdad/Iraq (1981). Mohammed Naji AL-Kabi is an assistant Professor in the Faculty of Sciences and IT, at Zarqa University. Prior to joining Zarqa University, he worked many years at Yarmouk University in Jordan, Nahrain University and Mustanserya University in Iraq. He also worked as a part time lecturer at Jordan University of Science and Technology, Princess Sumaya University for Technology, and Sunderland university. AL-Kabi's research interests include Information Retrieval, Web search engines, Data Mining, Social media, Natural Language Processing and Software Engineering. He is the author of more than 66 peer reviewed articles in these topics. His teaching interests focus on information retrieval, Web programming, data mining, DBMS (ORACLE & MS Access).

Heider Wahsheh, born in Jordan, in August 1987, he obtained his Master degree in Computer Information Systems (CIS) from Yarmouk University, Jordan, 2012. Since 2013 Mr. Wahsheh starts working as a lecturer in the college of Computer Science at King Khalid University, Saudi Arabia. His research interests include: Information Retrieval, Data Mining, and Mobile Agent Systems.

Izzat Alsmadi. An associate professor in software engineering. Born in Jordan 1972, Izzat Alsmadi has his master and PhD in software engineering from North Dakota State University (NDSU), Fargo, USA in the years 2006 and 2008 respectively. His main areas of research include: software engineering, testing, metrics, and information retrieval.

# RPOA Model-Based Optimal Resource Provisioning

Noha El.Attar

Ras-Elbar High Institute for Specific Studies and
Computer,
Ras El- Bar, Damietta, Egypt.

Wael Awad

Department of Mathematics and Computer Science
Faculty of Science, Port Said University
PortSaid, Egypt

Samy Abd El -Hafeez

Department of Mathematics and Computer Science
Faculty of Science, Port Said University
PortSaid, Egypt

Fatma Omara

Departement of Computer Science,
Faculty of Computers and Information, Cairo University,
Cairo, Egypt

*Abstract*—**Optimal utilization of resources is the core of the provisioning process in the cloud computing. Sometimes the local resources of a data center are not adequate to satisfy the users' requirements. So, the providers need to create several data centers at different geographical area around the world and spread the users' applications on these resources to satisfy both service providers and customers QoS requirements. By considering the expansion of the resources and applications, the transmission cost and time have to be concerned as significant factors in the allocation process.**

**According to the work of our previous paper, a Resource Provision Optimal Algorithm (RPOA) based on Particle Swarm Optimization (PSO) has been introduced to find the near optimal resource utilization with considering the customer budget and suitable for deadline time. This paper is considered an enhancement to RPOA algorithm to find the near optimal resource utilization with considering the data transfer time and cost, in addition to the customer budget and deadline time, in the performance measurement.**

*Keywords*—*Cloud Computing; Resource Provision; Data Communication; Particle Swarm Optimization*

## I. INTRODUCTION

One of the main definitions of cloud is that it is a cluster of distributed computers providing on-demand computational resources or services to the remote users over a network [1]. In other words, the cloud computing is a new paradigm for hosting and delivering services on demand over the internet where the users could access services depending on their Quality of services (QoS) requirements regardless where these services are hosted, they only care about how much they will pay, and how much time is expected to provide the required hardware and software resources [2] [3].

Cloud computing has the growing popularity and adoption feature [4]. This means that there is no data center with limited capacity. In other words, if a data center is become overloaded, it may be possible to relocate some workloads to another data center [5]. The workload sharing makes the cloud system has the ability to expand the resource pool and provides more flexible and cheaper resources [1].

Cloud computing delivers three application layers as services that are Infrastructure as a service (IaaS) for delivering IT infrastructure, Platform as a service (PaaS) for software development and deployment, and software as a service (SaaS) for providing software product [6]. depending on the, frameworks, or a final . IaaS rule is to provide Hardware as a Service by offering basic storage, memory and computing capabilities as standardized services over the network. To acquire computing resources, a user launches a server instance on the infrastructure of the cloud provider, thereby specifying the instance's characteristics such as the available processing power, main memory and I/O capacity. Most commonly, the notion of an instance materializes into a virtual machine that is launched on the physical IT infrastructure of the provider. Examples of successful IaaS providers are; Amazon EC2, Joyent, Rackspace. PaaS is a way to rent hardware, operating systems, storage and network capacity over the internet. It delivers a computing platform or software stack as a service to develope applications. This can broadly be defined as application development environment which is offered as a 'service' by the vendors. AppEngine, Bungee Connect, LongJump, Force.com, WaveMaker are examples of PaaS [7]. Instead of executing an application locally, this application could be located on the cloud and can be accessed online via web interfaceby SaaS. The applications are already-created as fully or partially remote services. Sometimes they are in the form of web-based applications and other times they consist of standard non remote applications with Internet-based storage or other network interactions [8]. Yahoo mail, Google docs applications, Salesforce.com CRM apps, Microsoft Exchange Online, Facebook are examples of SaaS [7].

Other important characteristics that distinguish Cloud Computing from other distributed systems are dynamic provisioning, Geo-distribution, ubiquitous network access, and Utility-based pricing [9].

One of main criteria affects data sharing in resource provision process under the cloud environment is the network performance, especially network bandwidth consumption. The network bandwidth can be stated as the amount of data shared over the network in a given time [11]. The relationship between

the characteristics of the offered traffic, the link speed and the resulting Quality of Service (QoS) should be concerned for calculating the overall cost of the resource provisioning. On the other hand, the users ask about the needed resources that achieve their services and deal with the providers with a pay-as-you-go strategy without concerning about the network communication characteristic because the providers do not offer a guaranteed network bandwidth to the users. But sometimes, achieving the requirements of one cloud user may make an influence on another cloud user in the same cloud environment due to the network latency which increases response time and decreases throughput. That may cause time delay in service providing [11]. The possibility of growing the Cloud environments by spreading the resources across wide geographical area using Wide Area Network (WAN), causes a delay for providing the services because of the limitation of the communication infrastructure, which is provided by the Internet [12].

The work in this paper is considered an enhancement of RPOA algorithm which is based on Particle Swarm Optimization to find the near optimal resource utilization with considering the data transfer time and the cost in the performance measurement [3].

The rest of this paper is organized as follows; the related works about the resource provision problem in cloud computing are presented in section 2. The principles of the resource provision in the cloud environment are discussed in Section 3. Section 4 states the influence of data communication on provisioning process. Section 5 represents the enhancement RPOA algorithm and its implementation and evaluation using CloudSim toolkit. The conclusions is presented in section 6.

## II. RELATED WORKS

Over the years, the demand of services provisioning through the Internet has been increased . In order to handle these huge numbers of users' applications, Cloud infrastructure providers (i.e., IaaS providers) have established data centers in multiple geographical locations to achieve availability and ensure reliability in case of site failures [13].

Many researchers tried to meet the conditions of Quality of Services (QoS) and agreements of Service Level Agreement (SLA) in the resource provisioning process.

Hu; Y., et. al. [14] have merged two resource allocation policies; Shared Allocation policy (SA) and Dedicated Allocation policy (DA). These allocation policies are evaluated by heuristic algorithm on basis of the smallest number of servers required to meet the negotiated SLA. Chaisiri; S., et. al. [15] have introduced an algorithm for resource provisioning based on the cost optimization by considering the demand price uncertainty in the provisioning stages to manage virtual machines. Byun; E., et. al. [16] have suggested a Partitioned Balanced Time Scheduling algorithm to estimate the minimum number of computing hosts required to execute a workflow within a user-specified finishing time. The main goal of this algorithm is to minimize the resource cost rather than the makespan of the workflow. This algorithm is for the automatic execution of large scale workflow-based applications on dynamically and elastically provisioned computing resources.

Jung; G., Sim; K. M., [17] have proposed the agent-based adaptive resource allocation algorithm. This algorithm aimed to satisfy customer needs of service in terms of fast allocation time and execution response time where the provider try to allocate each customer request to an appropriate data center among the distributed data centers by offering his resources under the infrastructure as a service model. The authors have proposed an adaptive resource allocation model that took into consideration both the geographical distance between the location of consumer and data centers and the workload of data center.

There is a proposed algorithm for resource provisioning by considering the communication link between the customer and the provider with minimum cost, as illustrated by Gaurav Raj [18]. He has proposed an Efficient Broker Cloud Management as a new provider cloud communication paradigm, explaining communication mechanism between the provider and the cloud using cloud exchange. He obtained optimum route on the basis of cost factor considering hops count, bandwidth and network delay using Optimum Route Cost Finder algorithm.

Gaurav Raj, Ankit Nischal [19] are tried to find an efficient way to utilize the resources within the cloud, and to create virtual machines with consideration to optimum cost based on a performance factor. This performance factor depends upon the overall cost of the resource, communication channel cost, and reliability and popularity factor. They proposed a framework for communication between the resource owner and the cloud using Resource Cloud Communication Paradigm (RCCP) and extended the CloudSim by adding provisioning policies and Efficient Resource Allocation (ERA) algorithm in VMM allocation policy as a decision support for resource provisioning [20].

El-Attar; N., et. al [3] have introduced a Resource Provision Optimal Algorithm (RPOA) based on Particle Swarm Optimization PSO to find a Workload Resource map WM that is commensurate with customer budget and suitable for deadline time. According to this algorithm, it is found that maximizing the performance of computing resource could be achieved by allocating its capacity for the maximum number of workloads using (PSO) algorithm with the utilization function. *let list*)

## III. RESOURCE PROVISIONING PROBLEM

Identifying the proper resource to the required service is considered a fundamental problem in Cloud computing. Cloud providers have to achieve the availability by allocating the appropriate resources to the required services without any confliction in the resource demands and with determining the right amount of resources required for the execution of services in an optimal way. This optimization aims to achieve the defined Quality of Service (QoS) that is affected by various parameters related to the application (i.e. security, accessibility, availability and reliability), the resources (i.e. availability, reliability, throughput and utilization) ,as well as, the user-defined requirements (i.e. cost and time) and the provider profit (i.e. resource utilization and revenue) [21].

In other words, the main goal of the provisioning policy is how to spread the applications load on convenient Cloud

resources to achieve the optimization objective of satisfying both service providers and customers QoS requirements (i.e. minimizing both response time and cost of resource utilization and, in the same time, maximizing the provider profit with the available customer's budget) [3]. The Cloud computing service provider's profit is achieved by providing high-quality services to the users through the efficient allocating of the resources on demand [22].

In [3], a Resource Provision Optimal Algorithm (RPOA) has been proposed based on particle swarm optimization to find the near optimal scheduling map of available resources with minimization of user response time and resource usage cost. According to the work in this paper, the RPOA has been enhancement and modified by considering the influence of communication parameters.

## IV. Data Communication In Resource Provisioning

One of the main features of Clouds is to spread workloads over distributed virtualization infrastructure and to cover larger geographical areas. The resource providers in the cloud environment rely on the concept of virtualization to supply their computational resources. The aim of virtualization technologies is to hide the underlying infrastructure by introducing a virtual layer between the physical infrastructure and the computational resources. In fact, the large clouds are based on exploitation of distributed virtualization infrastructures of other clouds to provide new types of services "Distributed IaaS, PaaS, and SaaS". On the other hands, the distributed Cloud service is composed of a set of VMs spread over a wide geographical area, coordinated in order to achieve a specific service which is provided on demand to the user in order to meet his requirements [12].

Another scenario of resource provision is to use the federation perspective that allows the Cloud providers to use virtualization infrastructures of other federated Clouds, trying to meet more and more kinds of required services [23, 24]. The federation makes the small and the medium Clouds able to increase their virtualization capabilities by using the virtualization infrastructures of other federated Clouds to achieve the required service.

By using the federation perspective, the resource provision is become more elastic by increasing the virtualization capabilities of Clouds and, on the other hand, to enable the Clouds providers to rent their computational and storage resources when their virtualization infrastructures are unused [12].

Any Cloud user cares about two main criteria; the service response time in the data center, i.e., turnaround time [25], and the cost of the available resources that can be allocated to the service. Generally, the allocation process of any cloud service can be divided into three consecutive stages; scheduling, computation and transmission [25].

On the other hand, the potential growth of Cloud environments due to the distribution of Cloud infrastructures over large scale geographical area makes the provisioning of services face some delay problems due to the latency that happened in transmission through the Internet [12].

Therefore, the transmission time has to be considered as a significant criterion in the allocation process specially when there is a need for transferring a huge dataset.

## V. The Enhancement of RPOA

According to the RPOA algorithm, the data centers are distributed over different resource pools. Each pool includes a specific resource type (e.g., computing power or memory capacity), and each computing resource is limited with the available number of hosts and both certain memory capacity and computing power consumption. On the other hand, each task workload is associated with a number of subtasks which require a certain amount of computing power or memory capacity depending on the workload type. In the same time, the RPOA algorithm is considered that the work loads are independent and have the same priority. By these conceptions, the measurement of the implementation's cost and time are become more accurate [3].

The work in this paper is considered a new version of RPOA which has been proposed in [3] with considering the influence of communication bandwidth and time latency beside the cost and time of the resources. This is because the network bandwidth and latency have a critical effect on the network performance, and accordingly on the provisioning process.

The network bandwidth can be defined as the rate at which the number of messages can be transferred from one point to another point in a given amount of time and it can be measured as baud rate, rate of data transfer and bit rate or throughput. The bandwidth that is used in cloud computing is the network bandwidth [11]. On the other hand, the network latency refers to any delay typically incurred in the processing of network data. This delay occurs because of the resources distribution over large number of computing nodes on the network as in the web based cloud applications [26].

### A. The Enhancement of RPOA Environment

According to RPOA algorithm, there are 'm' number of available resources, and 'n' workloads that contain '$j$' of subtasks, with consideration of some principle assumptions (fixed quantity of resources $'Q'$, with a definite price $'p_j'$ and a default execution time $'t_j'$). On the other hand, every workload $'i'$ has a set of tasks $'j'$ that need a specific resource quantity denoted by '$q$', and every customer has the availability of deciding the price of each task that can be paid '$bp_j^i$'. All of these consideration with a very important constraint which is ensuring that the total number of all workloads have to be not more than the available resources [3].

On the other hand, the cost and time of data transfer have to be defined, i.e. the communication link type and the available data packet size on the communication channel. The Transfer cost can be calculated according to equations (1, 2, 3) [27]:

*Cost of Transfer (CoT)= TC +RT+LC+Thr_c*     (1)

Where,

*Transmission Cost (TC) = TT* Cost per sec*     (2)

*Throughput Cost ( Thr_C)=*

*Total Bandwidth Price / Packet Size*    (3)

Routing Protocol (*RT*) depends on the network distance as shown in Table.1, and both Loading Cost (*LC*) and total bandwidth price are applied from the resource provider.

TABLE I.    DYNAMIC ROUTING PROTOCOL DEFAULT ADMINISTRATIVE DISTANCES [28]

| Route Source | Default Distance |
|---|---|
| Connected interface | 0 |
| Static route | 1 |
| Enhanced IGRP (EIGRP) summary route | 5 |
| Exterior Border Gateway Protocol (BGP) | 20 |
| Internal EIGRP | 90 |
| IGRP | 100 |
| OSPF | 110 |
| IS-IS | 115 |
| RIP | 120 |
| EIGRP external route | 170 |
| Interior BGP | 200 |
| Unknown | 255 |

Similarly, the transfer latency time will be calculated by equation (4) [29].

*Latency (Packet Delivery Time(PDT)) = TT +PD*    (4)

Where Transmission Time (TT), and Propagation Delay (PD) are calculated by equations (5, 6):

*TT=Packet Size /Bit Rate,*    (5)

which are stated in table 2, and

*PD=Distance/ Propagation Speed (PS)*    (6)

where Propagation Speeds are known as follows [30]:

a ) for copper wires= 2*108 meter/sec.

b) for wireless = 3*108 meter/sec.

TABLE II.    List of Network Technology bit rate [31]

| Technology | Rate bit/sec |
|---|---|
| Modem 56k (8000/8000 baud) (V.92) | 56.0/48.0 kbit/s |
| ISDN Basic Rate Interface (single/dual channel) | 64/128 kbit/s |
| IDSL (dual ISDN + 16 kbit/s data channels) | 144 kbit/s |
| HDSL ITU G.991.1 aka DS1 | 1,544 kbit/s |
| MSDSL | 2,000 kbit/s |
| SDSL | 2,320 kbit/s |
| SHDSL ITU G.991.2 | 5,690 kbit/s |
| ADSL (G.Lite) | 1,536/512 kbit/s |
| ADSL (G.DMT) | 8,192/1,024 kbit/s |
| ADSL2 | 12,288/1,440 kbit/s |
| ADSL2+ | 24,576/3,584 kbit/s |
| DOCSIS v1.0 (Cable modem) | 38,000/9,000 kbit/s |
| DOCSIS v2.0 (Cable modem) | 38,000/27,000 kbit/s |
| DOCSIS v3.0 (Cable modem) | 160,000/120,000 kbit/s |
| Uni-DSL | 200,000 kbit/s |
| VDSL ITU G.993.1 | 52,000 kbit/s |
| VDSL2 ITU G.993.2 | 100,000 kbit/s |
| BPON (G.983) fiber optic service | 622,000/155,000 kbit/s |
| GPON (G.984) fiber optic service | 2,488,000/1,244,000 kbit/s |
| Classic WaveLAN (Wireless) | 16,384 kbit/s |
| IEEE 802.11 (Wireless) | 16,384 kbit/s |
| RONJA (Wireless) | 81,920 kbit/s |
| IEEE 802.11a (Wireless) | 442,368 kbit/s |
| IEEE 802.11b (Wireless) | 90,112 kbit/s |
| IEEE 802.11g (Wireless) | 442,368 kbit/s |
| IEEE 802.16 (WiMAX) (Wireless) | 573,440 kbit/s |
| IEEE 802.11g with Super G by Atheros (Wireless) | 884,736 kbit/s |
| IEEE 802.11g with Nitro by Conexant (Wireless) | 1,146,880 kbit/s |
| IEEE 802.11n (Wireless) | 4,915,200 kbit/s |
| IEEE 802.11ac (maximum theoretical speed) (Wireless) | 58,1333,053kbit/s |

So, depending on RPOA implementation [3], there are *N* users ask for services that are defined

with a matrix $q_{nm}$ subtasks that have to be achieved by the available resources that have a vector $Q_m$ of resources fixed capacity $Q=[Q_1,Q_2,Q_3,.......Q$m*]*. The Total processing time $t_j^i$ and cost $c_j^i$ of every workload *i* on the suitable resource *j* is calculated according to equations 7, 8 [3] based on equation 1,4.

$$t_j^i = \sum \left( \frac{q_i}{Q_m} * t_j \right) + PDT_j^i \qquad (7)$$

$$c_j^i = \sum (p_j^i * t_j^i) + CoT_j^i \qquad (8)$$

RPOA uses the elasticity definition to find the existing relation between time and cost of both workloads and resources [32].

### B. The Enhancement of RPOA Implementation and Evaluation

The difficulty of obtaining workloads for real applications makes the evaluation of resource provisioning algorithm to be not an easy process. So, the CloudSim simulator is used to overcome this problem.

CloudSim with CloudAnalyst interface is used as Simulation framework to evaluate the Enhancement RPOA algorithm performance. CloudSim is a framework developed by the GRIDS laboratory of University of Melbourne which enables seamless modelling, simulation and experimenting on designing Cloud computing infrastructures [33]. CloudSim is a self-contained platform which can be used to model data centers, service brokers, scheduling and allocation policies of a large scaled Cloud platform. It provides a virtualization engine with extensive features for modelling the creation and life cycle management of virtual engines in a data center. CloudSim framework is built on the top of GridSim framework also is developed by the GRIDS laboratory, and it is written in java 6.1, because Java is a programming language designed for the distributed environment of the internet [34] [35]. CloudAnalyst has an easy graphical user interface that provides a high degree of configurability and flexibility by giving the ability of quickly and easily changing of the parameters that needed to be assumed in the simulation. Also CloudAnalyst gives a graphical output in the form of tables and charts to summarize the potentially large amount of statistics that is collected during the simulation with the ability of repeating the experiments under similar environment. Also, it gives more accurate evaluation by allowing the comparison between different scheduling algorithms [34].

RPOA is applied for memory provisioning; the user requests the memory on the available data centers. To evaluate the enhancement RPOA algorithm, some resource allocation requests are applied at different intervals of peak hours with different average of peak users. The datacenters that are used in this simulation have (x86) architecture, working on Linux operating system, and with Xen virtual machines that are based on time shared allocation policy.

The simulation is done over a various number of user bases requests on different available data centers. Every user base has its own budget, and asks for a specific amount of memory also it can be divided to independent subtasks; everyone has its own request and budget. On the other hand, every data center has a certain amount of memory, and has a predefined execution time. Also, the available bandwidth and its transmission cost are defined with the type of network communication and the available networks bit rate.

As shown in Figure 1, CloudAnalyst divides the worlds into regions to make the cost and time calculation more accurate.



Fig. 1. Region Boundaries Details in CloudAnalyst

Tables 3 and 4 show the variant parameters of a ten sample from hundred simulated experiments; respectively. Table 3 shows the user bases services request and contains all of the request details as, user region, request size, number of subtasks; if it contains subtasks; and user budget. Other important parameters are that the start and end peak hours that influence the time of server responding.

TABLE III. USER BASES DETAILS

Table 4, also shows a sample of the data centers details that contains data center region, Virtual machine's cost for every

| | Region | No. of Subtask per Request | Request Size | User Budget | Peak Hours (start) | Peak Hours (End) |
|---|---|---|---|---|---|---|
| UB1 | 3 | 20 | 8.678532 | 0.7 | 3 | 9 |
| UB2 | 1 | 10 | 8.583069 | 0.4 | 4 | 10 |
| UB3 | 5 | 5 | 7.271767 | 0.5 | 2 | 7 |
| UB4 | 0 | 40 | 3.290176 | 0.9 | 1 | 8 |
| UB5 | 2 | 8 | 23.97537 | 0.2 | 3 | 9 |
| UB6 | 2 | 8 | 8.678532 | 0.3 | 4 | 10 |
| UB7 | 5 | 25 | 8.583069 | 0.9 | 2 | 7 |
| UB8 | 3 | 30 | 7.271767 | 0.4 | 5 | 10 |
| UB9 | 1 | 20 | 3.290176 | 0.3 | 3 | 9 |
| UB10 | 0 | 15 | 23.97537 | 0.8 | 1 | 5 |

data center, data transfer cost, and Resource size per G. Byte.

TABLE IV. USER BASES DETAILS

| | Region | Cost per Proc $/Hr | Memory Cost | Trans. Cost Byte per Sec | Resource Size G. Byte |
|---|---|---|---|---|---|
| DC1 | 0 | 0.2 | 0.05 | 0.01 | 10.742188 |
| DC2 | 1 | 0.08 | 0.06 | 0.09 | 29.296875 |
| DC3 | 2 | 0.15 | 0.07 | 0.08 | 9.765625 |
| DC4 | 4 | 0.19 | 0.04 | 0.12 | 39.0625 |
| DC5 | 5 | 0.3 | 0.03 | 0.13 | 48.828125 |

Figure 2 is the simulation result interface from CloudAnalyst. It displays the location of every user base on the selected suitable datacenter.

Fig. 2.   Simulation Results

Figures 3 and 4 show the user base and data center configuration, in order. Figure 3 displays the details of user bases configuration (with their subtasks), and virtual machines that used to every user base.



Fig. 3.   UserBase Configuration

Also, Figure 4 displays the data center configurations with showing the physical hardware details of every data center.



Fig. 4.   Data Center Configuration

Figure 5 shows the chart of statistical results of the executed simulation. Only 50 requests are displayed from the hundred simulated user bases. This is to make the comparison results to be visible and clear. It displays the required amount of every user base's task and its relation with the whole capacity of the available resources.

This relation is to illustrate the importance of using the whole capacity of a data center to one or more of user bases that is because leaving a small amount of capacity in every datacenter will not be useful in another allocation process. But still there is another main condition in the allocation process. This condition is to search for the closest data center to the user base from both the latency and region, for trying to reach the goal of cost and time optimization by using RPOA.

Table 5 displays the overall cost that is paid for every datacenter containing the processing and the transfer cost.

TABLE V.        GRAND TOTAL COST FOR SOME DATACENTERS

| | | VM cost | Data Transfer Cost | Total Cost |
|---|---|---|---|---|
| DC1 | | 1.00 | 10784.16 | 10785.2 |
| DC2 | | 0.40 | 706.05 | 706.447 |
| DC3 | | 0.75 | 1401.43 | 1402.18 |
| DC4 | 0.95 | 936.79 | 937.744 | |
| DC5 | 1.50 | 14991.03 | 14992.5 | |

The main objective of the Enhancement RPOA algorithm is that find the near optimal allocation of the user requests on available data centers with the user budget restriction and deadline time with putting the influence of communication factors in consideration.

As shown in table 5, the data transfer cost has a significant effect on the total cost, this leads to take the distance between request region and resource region in consideration, and trying to allocate the request to the nearest data center which achieves the user demands.

The Enhancement RPOA algorithm provides a result by taking budget and deadline time into consideration as it can do; the error rate with the proposed sample of a hundred user base is nearly 18% for cost and 24% for time.

## VI.    COMPARATIVE STUDY

The first algorithm which we can use to compare with RPOA is Closest Data Center Algorithm (CDCA) [34]. This algorithm is considered the simplest service broker policy. It is based on creating an index table of all data centers indexed by their region. When the user base sends a message with his requirements, the algorithm looks for the data centers which are in the same region of the user base and picks from the data center located at the earliest/highest region in the proximity list and puts the request in the requests' queue of this data center [34].

The comparison between RPOA and CDCA in terms of cost reduction is illustrated in Figure 8. The CDCA's working mechanism is to provide resource with minimizing the cost of the provided service that is because the CDCA searches for the resource in the nearest region of user base that reduces the transfer cost, and so reduces the total cost.

Fig. 5.   Relation between the size of Resource and Request.



Fig. 6.   Comparison between Execution Cost and User Budget In RPOA



Fig. 7.   Comparison between Deadline Time and Execution Time In RPOA

Fig. 8.   Comparison between Cost and budget in RPOA and CDCA



Fig. 9.   Comparison between Execution and Deadline Time in RPOA and CDCA

On the other hand, RPOA gives good results in time optimization comparing with CDCA, as shown in Figure 10. RPOA tries to reduce execution time to fit deadline time, but CDCA does not care about time. It may waste time by letting the task in a long queue in the nearest resource.

The above comparisons can be summarized in that, RPOA algorithm can enhance the time optimization more than CDCA with 65%, but it has a weakness in reducing execution time with 13% than CDCA.

Fig. 10. Comparison between Cost and Budget in RPOA and BRTA



Fig. 11. Comparison between Execution and Deadline Time in RPOA and CDCA

The second Comparison will be done between RPOA algorithm and Best Response Time Algorithm (BRTA) [34]. This policy locates all the available data centers and indexes them. After the service broker receives the user's request, it identifies the closest data center (i.e. in terms of latency) and iterates through the list of all data centers and estimates the response time at each data center by querying the last recorded processing time using Internet Characteristics. If this time is recorded before a predefined threshold, the processing time for that data center is reset to 0. This means that the data center has been idle for duration of at least the threshold time. If the least estimated response time is for the closest data center, the broker selects the closest data center. Else, the broker picks either the closest data center or the data center with the least

response time with a 50:50 chance (i.e. latency load balanced 50:50) [34].

The comparison between cost and user budget in both RPOA and BRTA algorithms is displayed in Figure 10. BRTA algorithm cares only about how to provide the service with minimizing execution time as much as possible, but without considering the user budget, so it tries only to fit a deadline time without giving a suitable execution cost. Concluding form the BRTA characteristics, RPOA can be considered better than BRTA with 55% of cost reduction, but its performance in the execution time minimization is less than BRTA algorithm with 12% as shown in Figure 11.

## VII. CONCLUSIONS

The Enhancement RPOA is an algorithm that is implemented to optimize both response time and cost of

resource provisioning process. On the other hand, the problem is the limitation of the resources capabilities which needs to spread the workloads on different geographical areas around the world. This leads to consider the communication factors, (i.e., communication time and cost) into consideration, which influence the performance of the allocation process. CloudSim is the simulation toolkit which is used to implement and evaluate the Enhanced RPOA algorithm and find the workload –resource map that achieves the user requirements.

## REFERENCES

[1] J. Li, et.al., "Online Optimization for Scheduling Preemptable Tasks on Iaas Cloud Systems", J. Parallel Distrib. Comput., Vol. 72, Issue 5, 2012, pp. 666–677.

[2] Q. Zhang, et. al., "Cloud Computing: State-Of-The-Art and Research Challenges", J. of Internet Services and Applications, Vol. 1, No. 1, 2010, pp. 7–18.

[3] N. El.Attar, et. al., "Resource Provision for Services Workloads based on (RPOA)", Int. J. of Computer Science Issue (IJCSI), Vol. 9, Issue 3, 2012, pp. 553- 560.

[4] M. Armbrust, et. al., "Above the Clouds: a Berkeley View of Cloud Computing", Communications of the ACM, vol. 53, Issue 4, 2009, pp. 50-58.

[5] T. Forell, et. al, "Cloud Management Challenges and Opportunities", in: IEEE Inte. Symposium on Parallel and Distributed, 2011, pp. 881–889.

[6] R. Bossche, et. al, "Online Cost-Efficient Scheduling of Deadline-Constrained Workloads on Hybrid Clouds", Future Generation Computer Systems, Vol. 29, Issue 4, 2013, pp. 973–985.

[7] R. Grewal, P Pateriya, "A Rule-based Approach for Effective Resource Provisioning in Hybrid Cloud Environment", Int. J. of Computer Science and Informatics ISSN, Vol. 1, Issue. 4, 2012, pp. 2231 –5292.

[8] Y. Chen, S. Tai, "Optimal Provisioning of Resource in a Cloud Service", IJCSI, Vol. 7, No. 6, 2010, pp. 95-99.

[9] J. lmeida, et.al. , "Joint Admission Control And Resource Allocation In Virtualized Servers", Journal of Parallel and Distributed Computing, Vol. 70, No. 4, 2010, pp. 344-362.

[10] F. Silva, et al., "Application Execution Management on the InteGrade Opportunistic Grid Middleware", J. of Parallel and Distributed Computing, Vol. 70, No. 5, 2010, pp. 573- 583.

[11] K. Chandrakanth, S. Gayathri, "Examining Bandwidth Provisioning on Cloud Computing and Resource Sharing via Web", Int. J. of Engineering Science & Advanced Technology, Vol. 2, Issue-5, 2012, pp. 1372 – 1376.

[12] A. Celesti, et. al, "Virtual Machine Provisioning through Satellite Communications in Federated Cloud Environments", Future Generation Computer Systems, Vol. 28, Issue 1, 2012, pp. 85–93.

[13] R. Buyya, et. al., InterCloud: Utility-Oriented Federation of Cloud Computing Environments for Scaling of Application Services, in 10th Int.. Conf. on Algorithms and Architectures for Parallel Processing (ICA3PP), pp 13- 31, May 21-23, 2010.

[14] H. Wong, et. al., "Resource Provisioning for Cloud Computing", Conference of Center for Advanced Studies on Collaborative Research, 2009, pp. 101- 111.

[15] S. Chaisiri, et. al., "Optimization of Resource Provisioning Cost in Cloud Computing", IEEE Transactions on Services Computing , Vol. 5, Issue 2, 2012,    pp. 164-177 .

[16] E. Byun, et. al, "Cost Optimized Provisioning of Elastic Resources for Application Workflows", Future Generation Computer System, Vol. 27, No. 8, 2011, pp. 1011-1026.

[17] G. Jung, K. Sim, "Agent-based Adaptive Resource Allocation on the Cloud Computing Environment", International Conference on Parallel Processing Workshops (ICPPW), 2011,  pp. 345-351.

[18] G. Raj, et. al., "An Efficient Broker Cloud Management System", Int. Conf.  on Advances in Computing and Artificial Intelligence, 2011, pp. 72-76.

[19] G. Raj, A. Nischal, "Efficient Resource Allocation in Resource provisioning policies over Resource Cloud Communication Paradigm", Int. Journal on Cloud Computing: Services and Architecture (IJCCSA),Vol.2, No.3, Issue 3, 2012, pp. 11-18.

[20] R. Calheiros, et. al, "CloudSim: A Novel Framework for Modeling and Sim9ulation of Cloud Computing  Infrastructures and Services", http:/arxiv.org/ftp/arxiv/papers/0903/0903.2525.pdf, 2009.

[21] T. Varvarigou, et. al., "A Study on the Effect of Application and Resource Characteristics on the QoS in Service Provisioning Environments", Int. J. of Distributed Systems and Technologies, Vol. 1, Issue 1, 2010, pp. 55-75.

[22] H. Zhou, "Dynamic Resource Provisioning for Interactive Workflow Applications on Cloud Computing Platform", proceeding of: Methods and Tools of Parallel Programming Multicomputers - Second Russia-Taiwan Symposium, MTPP 2010, Vladivostok, Russia, May 16-19, 2010, pp. 116-125.

[23] J. Goiri, J. Torres,  "Characterizing Cloud Federation for Enhancing Providers' Profit", in: Cloud Computing, IEEE International Conference on Cloud Computing, 2010, pp. 123–130.

[24] R. Ranjan, R. Buyya, "Decentralized Overlay for Federation of Enterprise Clouds", Handbook of Research on Scalable Computing Technologies, 2010, pp. 191–217.

[25] X. Nan, et. al, "Optimal Resource Allocation for Multimedia Cloud Based on Queuing Model", IEEE 13th International Workshop on Multimedia Signal Processing (MMSP), 2011, pp. 1-6.

[26] Arista Whitepaper, Architecting Low Latency Cloud Networks, 2007, http://www.computerlinks.co.uk/FMS/22489.cloud_network_latency.pdf , visited on Sunday 16-June-2013 at 11 pm.

[27] H. Lehpamer, Microwave Transmission Networks: Planning, Design, and Deployment, chapter 1- Transmission Network Media, 2004, McGrow-HillCompanies, USA.

[28] Cisco, Configuring IP Routing Protocol-Independent Features, http://www.cisco.com/en/US/docs/ios/12_2/ip/configuration/guide/1cfin dep.html, visited on Friday, 14- June- 2013 at 10 am.

[29] B. Forouzan,  Data Communication and networking, fourth edition, Chapter 3 Data And Signals, McGraw-Hill Forouzan Networking Series, 2007.

[30] S. Alshaban, et. al., "Measurement of Transmission Time Delay and Efficiency of ATM LANE", Research Journal of Applied Sciences, Engineering and Technology, Vol. 2, Issue 2, pp.176-179, 2010

[31] L. Surhone, et. al., List of Device Bit Rates, Paperback Edition, Chapter-1, pp. 1-18, 2010, VDM Verlag Dr. Mueller AG & Co. Kg.

[32] I. Png, C. Cheng, Managerial economics, http://www.comp.nus.edu.sg /~ipng/mecon/sg/03elas_sg.pdf, 2001, visited on June - 8- 2013, 12:30 am.

[33] R. Buyya, et. al, Modeling and Simulation of Scalable Cloud Computing Environments and the CloudSim Toolkit: Challenges and Opportunities, in the 7th High Performance Computing and Simulation (HPCS 2009) Conf., June 21 - 24, 2009

[34] B. Wickremasinghe, "CloudAnalyst: A CloudSim-based Tool for Modelling and Analysis of Large Scale Cloud Computing Environments", Distributed Computing Project, Csse Dept., University Of Melbourne, 2009, pp. 433-659.

[35] R. Calheiros, et. al., "CloudSim: A Toolkit for Modeling and Simulation of Cloud Computing Environments and Evaluation of Resource Provisioning Algorithms", J. of Software: Practice and Experience (SPE), Vol. 41, Issue 1, 2011, pp. 23-50.

# Investigating the combination of structural and textual information about multimedia retrieval

Sana FAKHFAKH

Computer Science Department
Laboratory MIRACL, Institute of
Computer Science and Multimedia
University of Sfax, Tunisia

Mohamed TMAR

Computer Science Department
Laboratory MIRACL, Institute of
Computer Science and Multimedia
University of Sfax, Tunisia

Walid MAHDI

Computer Science Department
Laboratory MIRACL, Institute of
Computer Science and Multimedia
University of Sfax, Tunisia

*Abstract*—**The expansion of structured information in different applications introduces a new ambiguity in multimedia retrieval in semi-structured documents. We investigate in this paper the combination of textual and structural context for multimedia retrieval in XML document thus we present a indexing model which combines textual and structural information. We propose a geometric method who use implicitly of textual and structural context of XML elements and we are particularly interested by improve the effectiveness of various structural factors for multimedia retrieval. Using a geometric metric, we can represent structural information in XML document with a vector for each element.**

**Given a textual query, our model lets us combine scores obtained from each sources of evidence and return a list of relevant retrieved multimedia element. Experimental evaluation is carried out using the INEX Ad Hoc Task 2007 and the Image CLEF Wikipedia Retrieval Task 2010. The results show that combination of scores of textual modality and structural modality significantly improves compared results of using a single modality.**

*Keywords—Geometric distance; multimedia retrieval; element; structure; document modeling*

## I. INTRODUCTION

This paper falls under the context of multimedia retrieval in XML documents. The need with this kind of information is justified by quick change of scopes of application which use structural documents (format HTML or XML) what imposes new challenges in the field of search for information. Indeed, nowadays XML document passed a simple tool for exchanging data to a new storage medium. XML document includes textual element and multimedia element such as image, audio and video. These elements are organized according to structure which includes information notably although there is not only one manner to organize contents. However, the choice of structure depends greatly on the context of use of the textual contents.

Mainly in the literature, there are two main classes of approaches in the field of multimedia retrieval: retrieval methods based on multimedia content (MR-content) and multimedia methods to retrieval based on context (MR-Context).

The approaches of the multimedia retrieval based on content use specific features of low level according to type of media [1][2]. We can cite for example image retrieval that exploits visual features (color, texture, forms…). These methods have proven effective with media "image" in well defined fields such as medical field this is due to requirement for thorough knowledge of distinctive media. This type of research can be applied to only one type of media in system due to lack of semantic representation in media content.

The approaches of the multimedia retrieval based on context do not depend on type of media in question [3] [4]. Indeed, these methods rely on information surrounding the multimedia element representing its semantic description. Multimedia retrieval based on textual context is most used, although the structural context remains an obvious source which plays a part paramount in understanding of structured documents.

In this article, we focus on techniques for multimedia retrieval based on textual and structural context in XML documents. This type of document includes textual information and structural constraints. So, XML document cannot be effectively exploited by classical techniques of IR, which regard document as a plane source of information.

The implicit incorporation of multimedia elements in XML documents requires the exploitation of textual context for multimedia retrieval. However, the textual context remains insufficient in most of time. The idea is to calculate the relevancy score of media element based on information from the textual and structural context to answer a specific information needs of user, expressed as query composed of set of keywords.

Let us take for example an image media. If we exploit the image context which is composed by description of its contents such as its title, name, descriptive texts which surround it, title of XML document ... In following figure, we present document extracted from "WIKIPEDIA" encyclopedia describing lion. We notice the existing simultaneous textual and multimedia information. For image retrieval from time after "Pleistocene", we extract information from the textual description, not from title (figure 1).

Fig. 1.   Example of a multimedia object context.

## II.   RELATED WORKS

In our work, we will be interested by media "image". Most existing work in this area uses the information from textual description of image. There are other sources of evidence that were used as visual descriptors, information from link around the image [5], and structure of XML document. To resolve difficulties in multimedia retrieval field, you must define adequate source of evidence for representation a multimedia element and defining appropriate indexing model.

In this context, we present our structural indexing system combining conceptual information for semi-structured documents dedicated to approximate retrieval data. We begin with an overview of existing work in multimedia retrieval. Then we turn to the presentation of our approach while detailing the preprocessing, extraction of textual and structural and phase calculation relevance of multimedia element in information a better response to needs expressed by user. Finally, we present the results of applying our method on two universal bases "INEX 2007" and "ImageCLEF 2010".

The advent of structured documents has caused new problems in information retrieval world, and more specifically in multimedia elements retrieval. These problems are strongly related to nature of these documents that provide the structure as a new source of evidence. Thus, nowadays, XML documents include multimedia elements of different types (audio, video and image) implicitly embedded in the textual elements. These multimedia elements (such as physical objects) do not contain enough information to be able to answer a given query. Therefore, the calculation of relevance score of multimedia element must be linked to textual and structural information provided by other nodes XML [5].

Several works deal XML document as a flat source of information and ignore the structure of XML documents. In this context, [6] say: "Ignore the document structure is to ignore its semantics". Indeed, XML document is used to describe a set of data by a structure that provides a semantic

lexicon. Thus, it facilitates the presentation of information in terms of interpretation and exploitation. Replying to this need, new works appear in the field of multimedia retrieval that takes in account the structure as source of relevant information. Existing work in structured retrieval of multimedia elements is decomposed in two classes.

The first class includes some works which proceed to adopt some traditional technical of retrieval information as language model. In this context, the team CWI/UTwente performs a step of filtering results to keep the fragments containing at least one multimedia element [7][8].

The second class includes the specific work to be structured multimedia retrieval. This class uses the structure as a source of evidence in the process of selection of multimedia elements. As first step, [9] proposed a method which combines structure of XML document (XPath) with the use of links (XLink). This method consists to divide XML document into regions. Each region represents an area of ancestors of the multimedia element. His score is calculated in function of the scores of each region. This method exploits vertical structure only. In a second time, [10] have used the addition of horizontal structure to the notion of hierarchy. [10] use a method called "CBA" (Children, Brothers, Ancestors), which takes into consideration the information carried by the children , brothers and fathers nodes for calculate the relevance of multimedia elements. The authors propose an alternative method "OntologyLike" which is based on the identification of XML document to ontology. To calculate the similarity between nodes the authors use similarity measures that are mainly based on the number of edges to calculate the distance between nodes.

There are other approaches to multimedia retrieval are based on exploitation of links in XML document [11]. This work was improved by proposing a hybrid approach that combines structure with using of links that is consider as semantic links [12]. This method above consists to divide the document into regions according the hierarchical structure and the location of image in document. This factor plays a role in the weighting of links for compute the score of image.

In this paper, we propose a new metric for multimedia retrieval in XML documents which involves the use of geometric distances to calculate the relevance of each node from the multimedia node. This method consists of placing the nodes of XML document in Euclidean space and defines each node by a vector of coordinates to calculate then the distance between each pair of nodes. This distance will play a beneficial role in to calculate the score of multimedia element.

## III.   FROM XML ELEMENT TO GEOMETRIC CHARACTERISTIC

We focus on techniques for multimedia retrieval based on textual and structural context in XML documents. XML documents cannot be effectively exploited by classical techniques of IR, which regard document as a bog of words. Therefore, the calculation of relevance score of multimedia element must be linked to textual and structural information provided by other nodes XML [5]. Thus, it facilitates the presentation of information in terms of interpretation and

exploitation. Replying to this need, we propose a new method in the field of multimedia retrieval that takes into account the structure as a source of evidence and its impact on search performance. We present a new source of evidence dedicated to multimedia retrieval based on the intuition that each textual node contains information that describes semantically a multimedia element. And the participation of each text node in the score of a multimedia element varies with its position in there XML document. To compute the geometric distance, we initially place the nodes of each XML document in a Euclidean space to calculate the coordinates of each node by algorithm 1. Then, we compute the score of a multimedia element depending on the distance between each textual node [15].



Fig. 2. The steps of passing an XML document to geometric representation.

Figure 2 shows the steps of passing an XML document to a geometric representation of the XML elements in a Euclidean space. The first step consists to present a XML document as XML tree to take into account XML document properties.

An XML tree is described by a set of relationships between nodes. Formally an XML tree is a pair A = (E, *R*) where E is a set of XML elements and $\subset E^2$ , ( *(p,q)* $\in R$ if *p* is the parent of *q* ) is a set of relations satisfying:

$$\exists! \boldsymbol{r} \in \boldsymbol{E}, \forall \boldsymbol{q} \in \boldsymbol{E}, (\boldsymbol{r}, \boldsymbol{q}) \in \boldsymbol{R} \tag{1}$$

With *r* is the root of the tree.

$$\forall \boldsymbol{p} \in \boldsymbol{E} - \{\boldsymbol{r}\}, \exists! \boldsymbol{q} \in \boldsymbol{E}, (\boldsymbol{p}, \boldsymbol{q}) \in \boldsymbol{R} \tag{2}$$

Each node has a parent except the root *r*.

In second step, we will spend to presentation of XML tree in a geometric representation. This step is mainly based on equalities extraction in XML tree according to our proposed hypotheses.

The XML tree representation allowed us to unveil certain relationships of neighboring, brotherhood and offspring. Indeed, the distance *d* which separate two or more brothers with their common ancestors iteratively is the same. And brothers of the same hierarchical level are equidistant.

These distances are defined according to the relationship of contiguity and semantic similarity between nodes. These distances are not quantized but will be extracted in function of the position of each textual node in XML tree.

All these properties result in: For all $q_i = (x_{i1}, x_{i2} \cdots x_{im})$ and $q_j = (x_{j1}, x_{j2} \cdots x_{jm})$ where Q is a set of vectors in $\mathbb{R}^m$

- In the same hierarchy, if there are more than two brothers then their adjacent nodes are equidistant:

**Property 1**

$$\forall q_i, q_j, q_k \in Q, if \; A_1(q_i) = A_1(q_j) = A_1(q_k)$$

$$d(q_i, q_j) = d(q_j, q_k)$$

- The distance between any node and its descendants is the same:

**Property 2**

$$\forall q_i, q_j, q \in Q, \qquad n \in \mathbb{N},$$

$$if \; A_n(q_i) = A_n(q_j) = q$$

$$d(q_i, q) = d(q_j, q)$$

With $n \in \mathbb{N}^*$ , we define function $A_n$ by:

$$\forall q \in E,$$

$$A_n(q) = \begin{cases} \{q\} \; if \; n \; = \; 0 \\ A_{n-1}(p) \; if \; \exists \; p \; \in \; E, (p,q) \; \in \; R \; and \; n \; > \; 0 \\ \phi \; else \end{cases}$$

From these relationships, we can generate system of equations taking into account for kinship relationships nodes based on hierarchy and adjacency. These relationships are decried by equalities in this order (these equations are only examples):

$$d(n_1, n_2) \; = \; d(n_1, n_3)$$

$$d(n_1, n_2) \; = \; d(n_1, n_4)$$

$$d(n_1, n_7) \; = \; d(n_1, n_8)$$

$$d(n_1, n_7) \; = \; d(n_1, n_9)$$

These distances are defined according to the relationship of contiguity and semantic similarity between nodes. They are not quantized but will be extracted in function of the position of each textual node in the XML tree. The resulting system is nonlinear, its resolution requires the use of an approximate resolution iteratively method where we used iterative solution method (see Algorithm 1).

The process begins by assigning to each XML node a random vector followed by tries to improve the coordinate values of each node according to an error value (the sum of the squared deviations). At each iteration, the coordinates are improved together with the minimization of this error. The algorithm stops when the error reaches its minimum value (no improvement is possible). Let *Q* the set of vectors obtained at a given iteration during the running of the algorithm, the error is defined by:

$$error(Q) = \sum_{\substack{\forall q_i, q_j, q_k \in Q, \\ A_1(q_i) = A_1(q_j) = A_1(q_k)}} (d(q_i, q_j) - d(q_j, q_k))^2$$

$$+ \sum_{\substack{\forall q_i, q_j, q \in Q, \ n \in \mathbb{N} \\ A_n(q_i) = A_n(q_j) = q}} (d(q_i, q) - d(q_j, q))^2$$

---

**Algorithm 1** Resolution algorithm approximate nonlinear system of equations

---

**Require:** $(Q = (q_1, q_2 \cdots q_{|Q|}), R)$ : an XML tree as $q_i = (q_{i1}, q_{i2} \cdots q_{im}) \quad \forall i \in [1, |Q|]$

$m$:dimension

**for** $(i, j) \in [1, Q]^2$ **do**
$q_{ij} \leftarrow random\ value$
**end for**

$Q_1 \leftarrow (q_1, q_2 \cdots q_{|Q|})$

**Repeat**
$P \leftarrow Q_1$
**for** $(i, j) \in [1, Q]^2$ **do**
$Q_2 \leftarrow (q_1, q_2 \cdots q_{i-1}, q_i + d_j(1), q_{i+1} \cdots q_{|Q|})$
$Q_3 \leftarrow (q_1, q_2 \cdots q_{i-1}, q_i + d_j(\varepsilon), q_{i+1} \cdots q_{|Q|})$
$Q_4 \leftarrow (q_1, q_2 \cdots q_{i-1}, q_i + d_j(1 - \varepsilon), q_{i+1} \cdots q_{|Q|})$
$t \leftarrow 0$

**While** *error(Q1)> error(Q2)> error(Q3)> error(Q4)* **do**
$Q_4 = (q_1, q_2 \cdots q_{i-1}, q_i + 2^t d_j(1), q_{i+1} \cdots q_{|Q|})$
$t = t + 1$
**end while**

$t \leftarrow 0$

**While** *error(Q1)< error(Q2)<error(Q3)< error(Q4)* **do**
$Q_1 = (q_1, q_2 \cdots q_{i-1}, q_i - 2^t d_j(1), q_{i+1} \cdots q_{|Q|})$
$t = t + 1$
**end while**

**While** $|error(Q1) - error(Q2)| > \varepsilon$ **do**
$Q_5 \leftarrow \dfrac{Q_1 + Q_2}{2}$
let $Q_5 = (p_1, p_2 \cdots p_{|Q|})$

**if** $error(p_1, p_2 \cdots p_{i-1}, p_i - d_j(\varepsilon), p_{i+1} \cdots p_{|Q|})$ $>$ $error(p_1, p_2 \cdots p_{i-1}, p_i + d_j(\varepsilon), p_{i+1} \cdots p_{|Q|})$ **then**

---

$Q_1 \leftarrow Q_5$
**else**
$Q_2 \leftarrow Q_5$
**end if**

**end while**

**end for**

**until** $P = Q_1$

---

Where $m$ is the dimension of the Euclidean space and $\forall v \in \mathbb{R}, D_j = (d_1, d_2 \cdots d_m)$ is such as:

$$d_k = \begin{cases} 0 \ if \ k \neq j \\ v \ otherwise \end{cases}$$

### A. INDEXING SYSTEM

We propose an indexing system ***MXS-index*** composed by two parties: party of textual indexing and party of structural indexing. In first party, our approach uses NLP (Natural Language Processing) techniques to extract the candidate XML nodes of the resulting indexing. The weight of these nodes is depending on the frequency of each of these terms and the number of elements in the corpus according to the number of elements containing the term. In Second party, we built structural index using information extract from XML tree and geometric metric.

Each XML node will be presented by a characteristic vector (figure 3). We start by extract geometric proprieties. And we compute coordinates of each XML nodes. This party is accompanied by generating XML data model which processes ancestor, descendant and proximity relationships (figure 4).



Fig. 3. Geometric characteristic vector of XML node

Figure 5 schematize the process of textual and structural indexing XML documents with our indexing system. Well as the transition of XML document as a tree presentation to geometric presentation in Euclidean space.

Fig. 4.    Architecture of our indexing model MXS – index.



Fig. 5.    Treatment process of XML document.

## B. ADDING STRUCTURAL INFORMATION IN THE RELEVANCE SCORE OF MULTIMEDIA ELEMENT

A multimedia element (e.g. image) does not contain textual content. Its score is based on textual nodes in its neighborhood. The transition from the XML tree structure representation of elements in a Euclidean space, where we exploit the dissimilarity distances separating a multimedia node and other textual nodes, is performed by extracting the equations satisfying the properties defined earlier and the application of algorithm 1. To calculate the distance between a node $n$ and multimedia element $H$, we calculate the Euclidean distance between their respective feature vectors $q_n$ and $q_H$:

$$dist(n, H) = \sqrt{\sum_{i=1}^{m}(q_n - q_H)^2} \tag{3}$$

With $m$ is the dimension of the Euclidean space. $q_n$ is defined by: $q_n = (xn_{i1}, xn_{i2} \cdots xn_{im})$ with $xn$ are the vector characteristics of node $n$. And $q_H$ is defined by: $q_{Hn} = (xH_{i1}, xH_{i2} \cdots xH_{im})$ with $xH$ represent the coordinates compose the vector characteristics of a node $H$. We calculate the score for each textual node depending on the frequency of each term *(tf)* and the number of elements in the corpus according to the number of elements containing the term *(idf)*.

A textual node is presented by: $n = (n_1, n_2 \cdots n_{|v|})$ where $n_i$ is the weight of the term $t_i$, $v$ is the set of indexing terms:

$$ni = tf(t_i, n) \times idf(t_i) \qquad (4)$$

With

$$idf(t_i) = log\left(\frac{N}{N_i}\right) \qquad (5)$$

Where $N$ is the total number of XML elements in the corpus, $N_i$ is the number of elements that contain the term $t_i$ and $tf(t_i, n)$ is the frequency of the term $t_i$ in node $n$. The score of textual node depends on the weight of each indexing term. A query is made by the list $v = (v_1, v_2 \cdots v_{|v|})$ where $v_i \in \{0, 1\}$ (0: not exist, 1: exist) according membership $t_i$ at the query. The score of textual node $n$ for the query $q$ is defined by:

$$rsv(q, n) = q \times n^T = \sum_{i=1}^{|V|} q_i \times n_i \qquad (6)$$

Where $\mu$ is the set of textual elements. The score of multimedia node $H$ is defined by:

$$rsv(q, H) = \sum_{n \in \mu} \frac{rsv(q, n)}{dist(n, H)} \qquad (7)$$

With *dist (n, H)* is the distance between feature vectors corresponding to the nodes $n$ and $H$. This equation leads to assign the importance of contribution of all nodes in computing the score of multimedia element that shows its beneficial impact in multimedia retrieval.

## IV. EVALUATION AND RESULTS

We evaluate our system into two databases extracted from two collections: INEX 2007 (Initiative for the Evaluation of XML Retrieval) Ad Hoc task [13] and ImageCLEF 2010 Wikipedia image retrieval task [14]. These databases are composed by XML documents extracted from Wikipedia (Table I).

TABLE I.     INEX 2007 AND IMAGECLEF 2010 COLLECTIONS

| Company | INEX 2007 | ImageCLEF 2010 |
|---|---|---|
| Task | *Collection XML Ad Hoc* | *Wikipedia Retrieval* |
| Number of XML document | 659388 | 237434 |
| Number of image | 246730 | 237434 |
| Topics | 19 | 70 |

We evaluate our method with using only textual context (TC). The XML structure is not taken account. For INEX 2007 and ImageCLEF 2010 test set, we respectively obtain the following MAP values: 0.2376 and 0.1674. In the second time, we use XML structure will determine the image relevance score and will differentiate between images (using textual and structural context TC and TS). The evaluation results show

that this method provides a MAP which is equal to 0.2572 as MAP with using "ImageCLEF 2010" collection. The result has been improved significantly with the "INEX 2007" collection to 0.3102 as MAP. This increase is due to nature of "INEX 2007" collection that includes XML documents with heterogeneous structure.

So in "INEX 2007" collection we find documents with high depth. This factor highlights structural information and amplifies effect textual information based on computed distances. For against, our system is more stable with "ImageCLEF 2010" collection, this is due to rapid convergence of results. With our measure, we have shown that combined use of textual and structural context can properly determine the relevance of multimedia element, and the structure plays a primordial role in multimedia retrieval (Figure6).



Fig. 6.    Results of the impact our approach on INEX 2007 and ImageCLEF 2010 based in MAP(Mean Average Precision).

## V. CONCLUSION

In this paper, we propose a novel approach for multimedia retrieval in XML documents. This method consists to calculate the score of multimedia element according the textual context provided by nodes in proximity and structural context from distance between nodes and multimedia element. Thank to geometric metric, we could assign a weight to each textual node in the XML document. Although all textual parts are useful, they should not all be taken into account with the same importance degree.

Experiments show the interest of our method on INEX 2007 and ImageCLEF 2010 collections. Our work is focused on media image but it can be used with any other media, since the visual context of multimedia objects is not used.

In the future, we want to exploit another factor to calculate the relevance of multimedia element such as the title of image, the weighting of the links in XML document ... As well as another source of evidence as visual descriptors and the study parameters combination of using of structural, textual and visual context.

REFERENCES

[1] M. S. Lew, "Content-based multimedia information retrieval: State of the art and challenges", *ACM Trans. Multimedia Comput. Commun. Appl*, vol. 2, pp. 1–19, 2006.

[2] A. W. M. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-based image retrieval at the end of the early years", IEEE Trans. Pattern Anal. Mach. Intell., vol. 22, no. 12, pp. 1349–1380, Dec. 2000. [Online]. Available: http://dx.doi.org/10.1109/34.895972

[3] H. Elghazel, K. Idrissi, A. Baskurt, and C. Ben Amar, "Approche textuelle pour la recherche d'image", *in 3rd International Conference on Sciences of Electronic, Technologies of Information and Telecommunications SETIT 2005*, Mar. 2005. [Online]. Available: http://liris.cnrs.fr/publis/?id=2153

[4] D. Tjondronegoro, J. Zhang, J. Gu, A. Nguyen, and S. Geva, "Integrating text retrieval and image retrieval in xml document searching", in *INEX*, 2005, pp. 511–524.

[5] A. Hliaoutakis, G. Varelas, E. Voutsakis, E. G. M. Petrakis, and E. Milios, "Information retrieval by semantic similarity*", in Intern. Journal on Semantic Web and Information Systems (IJSWIS).Special Issue of Multimedia Semantics*, 2006, pp. 55–73.

[6] T. Schlieder and M. Holger, "Querying and ranking xml documents", *Journal of the American Society for Information Science and Technology*, vol. 53, pp. 489–503, 2002.

[7] T. Tsikrika, P. Serdyukov, H. Rode, T. Westerveld, R. Aly, D. Hiemstra, and A. P. V. de, "Structured document retrieval, multimedia retrieval, and entity ranking using pf/tijah", *in 6th Initiative on the Evaluation of XML Retrieval*, INEX 2007, ser. Lecture Notes in Computer Science, vol. 4862. London: Springer Verlag, March 2008, pp. 306–320. [Online]. Available: http://doc.utwente.nl/64734/

[8] T. Westerveld, H. Rode, R. O. van, D. Hiemstra, G. Ramirez, V. Mihajlovic, and A. V. de, "Evaluating structured information retrieval and multimedia retrieval using pf/tijah", *in Comparative Evaluation of XML Information Retrieval Systems,* ser. Lecture Notes in Computer Science, N. Fuhr, M. Lalmas, and A. Trotman, Eds., vol. 4518. Berlin, Germany: Springer Verlag, 2007, pp. 104–114. [Online]. Available: http://doc.utwente.nl/64261/

[9] Z. Kong and M. Lalmas, "Xml multimedia retrieval", *in SPIRE*, 2005, pp. 218–223.

[10] M. Torjmen, K. Pinel-Sauvagnat, and M. Boughanem, "Using textual and structural context for searching multimedia elements", *IJBIDM*, vol. 5, no. 4, pp. 323–352, 2010.

[11] H. Awadi and M. Torjmen, "Exploitation des liens pour la recherche d'images dans des documents xml", *CORIA* ,March 2010.

[12] H. Aouadi, M. Torjmen-Khemakhem, and M. B. Jemaa, "Combination of document structure and links for multimedia object retrieval", *Journal of Information Science*, vol. 38, no. 5, pp. 442–458, October 2012.

[13] N. Fuhr, J. Kamps, M. Lalmas, S. Malik, and A. Trotman, "Overview of the inex 2007 ad hoc track", *in INEX, 2007*, pp. 1–23.

[14] A. Popescu, T. Tsikrika, and J. Kludas, "Overview of the wikipedia retrieval task at imageclef 2010", *in CLEF (Notebook Papers/LABs/Workshops),* 2010.

[15] S. Fakhfakh, M. Tmar, and W. Mahdi, "A new metric for multimedia retrieval in structured documents", *in ICEIS (2)*, 2013, pp. 240–247.

# Audio Search Based on Keyword Spotting in Arabic Language

Mostafa Awaid
Biomedical Engineering Department
Higher Technological Institute
10th of Ramadan City,
Egypt

Sahar A. Fawzi
Systems and Biomedical
Engineering Department
Faculty of Engineering, Cairo
University Giza, Egypt

Ahmed H. Kandil
Systems and Biomedical
Engineering Department
Faculty of Engineering, Cairo
University Giza, Egypt Systems and
Biomedical Engineering Department
Higher Institute of Engineering
El-Shorouk, Egypt

*Abstract*— **Keyword spotting is an important application of speech recognition. This research introduces a keyword spotting approach to perform audio searching of uttered words in Arabic speech. The matching process depends on the utterance nucleus which is insensitive to its context. For spotting the targeted utterances, the matched nuclei are expanded to cover the whole utterances. Applying this approach to Quran and standard Arabic has promising results. To improve this spotting approach, it is combined with a text search in case of the existence of a transcript. This can be applied on Quran as there is exact correspondence between the audio and text files of each verse. The developed approach starts by text search to identify the verses that include the target utterance(s). For each allocated verse, the occurrence(s) of the target utterance is determined. The targeted utterance (the reference) is manually segmented from an allocated verse. Then Keyword spotting is performed for the extracted reference to the corresponding audio file. The accuracy of the spotted utterances achieved 97%. The experiments showed that the use of the combined text and audio search has reduced the search time by 90% when compared with audio search only tested on the same content. The developed approach has been applied to non-transcribed audio files (preaches and News) for searching chosen utterances. The results are promising. The accuracy of spotting was around 84% in case of preaches and 88% in case of the news.**

*Keywords—Speech Recognition; Keyword Spotting; Template Matching*

## I. INTRODUCTION

Keyword spotting (KWS) is a technique used to allocate and identify target words/utterances in continuous speech. Keyword spotting systems can be classified into two categories: speaker dependent and speaker independent. For speaker dependent systems, models are developed for a specific speaker. While speaker independent systems need to be more generic and hence need more complex design.

Arabic language is the official language in more than twenty countries with population of more than one billion persons. Since it is the language of Islam religion, more people need use it and to learn its proper pronunciation. Arabic syllables begins with a consonant (c) followed by a vowel (v) or long vowel (v:) and may include one or two extra consonants. Syllables are classified according to the length of

the vowels, which also known as Harakatt [1]. The five basic syllable structures in classical Arabic are: CV, CV: , CVC , CV: C , and CVCC.

Audio keyword spotting systems are difficult due to the huge variability of pronunciations between different speakers or even between repetitions of the same word by the same speaker. There exist different approaches to implement audio keyword spotting systems. Template matching approaches are used in small-scale systems and may result in accurate results when exact matching is needed [2, 3].

For audio files with corresponding text available, as in the case of the holly Quran, a text search can be used to help allocate the sentences (verses) containing the requested utterance.

This paper is organized as follows. Section two includes a description of the system. The results and discussion are presented in section Three. The conclusions are given in section four.

## II. DESCRIPTION OF THE SYSTEM

The proposed audio keyword spotting system is supposed to search and allocate requested speech segments (word, connected words, sentences) within continuous speech using a Template Matching based approach with the help of text search. There is no restriction concerning the number of occurrences of the word.

The system is divided into two successive phases. The first phase is the text search in which the target utterance is given as a text for the system. The text search results in a set of sentences that include the target text. A target utterance is segmented from one of allocated audio files. Then, all allocated audio files is searched for the targeted utterance (Keyword spotting phase), as shown in Fig. 1.

The system can be used to extract the matched utterances of a given word or phrase independent of their contextual sensitivity. The developed approach has overcome this difficulty by extracting the nucleus of the utterance defined as the reference speech segment (excluding the peripheral syllables). After allocating the targeted utterances, a reconstruction procedure is applied to accomplish full matching

with the original utterance. The allocation process was implemented through Pre-processing, Features Extraction, and Classification (Cross Correlation or Minimum Mean Square Error) [4].



Fig. 1.   General Block Diagram of the system [4]

### A. Text Search:

The text file is searched for the target text ignoring the vowelization differences between the targeted word(s) and the matched one(s) in the sentence.  The audio file is segmented into shorter audio files knowing that each silence corresponds to a sentence separator. The allocated sentences and their corresponding audio files are the new search domain.

### B. Targeted utterance preparation:

The first step to create feature vectors representing the acoustic signal is pre-processing. A high pass filter is used to decrease the noise and to flatten the speech signal spectrum (Pre-emphasis), using (1).

$$H(z) = 1 - 0.95z^{-1} \tag{1}$$

Since the vocal tract changes relatively slow, speech is considered a random process with slowly varying properties [5]. So, the speech utterance is divided into a number of overlapping frames having durations of around 10 msec. A Hamming window is applied to minimize the discontinuities at the beginning and end of each frame. The Hamming window is given by (2).

$$W(n) = 0.54 - 0.46 \, \cos\left(\frac{2\pi n}{N-1}\right) \tag{2}$$

In order to omit the co-articulation effect, frames corresponding to the first and last syllables of the utterance are ignored. The remaining frames represent the utterance nucleus.

### C. Features Extraction

The speech features techniques used are Mel-Frequency Cepstral Coefficient (MFCC) and Linear Predictive Cepstral Coefficient (LPCC).

#### 1)   Mel-Frequency Cepstral Coefficient (MFCC)

MFCC is a popular feature set, used in speech recognition, and based on the frequency domain of Mel scale for the human ear scale [3]. Mel-scale is based on filter bank processing. The Mel-frequency scale formula is based on mathematical equation given by (3).

$$f_{Mel} = 1127.01 \ln\left(\frac{f}{700} + 1\right) \tag{3}$$

Steps to derive MFCC:

- Fourier transform of each frame of the signal is obtained.

- Mel scaling is applied using triangular overlapping windows.

- Calculate the log of the power spectrum at each of the Mel frequencies.

- Compute the discrete cosine transform (DCT).

- The MFCCs are the amplitudes of the resulting spectrum.

#### 2)   Linear Predictive Cepstral Coefficient (LPCC)

LPCC is one of the most powerful speech analysis techniques for extracting good quality features [6].   The process for obtaining the LPCC features vectors is shown in Fig. 2



Fig. 2.   Block diagram of the computation steps of LPCC.

The notion behind LPCC is to model the human vocal tract by an all-pole filter. The LPC Coefficients $a_i$, are the coefficients of the all pass transfer function $H(z)$ modeling the vocal tract [7], as shown by (4):

$$H(z) = \frac{1}{1 - \sum_{i=1}^{p} a_i z^{-i}} \tag{4}$$

### D. Classification:

Two classification methods are used.

#### 1)   Mean Square Error (MSE)

Mean Square Error (MSE) is a signal fidelity measure used to compare two signals. Such a measure provides a quantitative score describing the degree of fidelity/ similarity [8]. Since $\mathbf{x} = \{xi | i = 1, 2, \ldots, N\}$ and $\mathbf{y} = \{yi | i = 1, 2, \ldots, N\}$ are two finite-length, discrete signals representing two distinct utterances, where $N$ is the number of signal frames and $xi$, $yi$ are the features vectors of frames constituting $x$ and $y$, respectively. The MSE between the two signals is determined be (5).

$$MSE(x, y) = \frac{1}{N} \sum_{i=1}^{N} (x_i - y_i)^2 \tag{5}$$

Finally, the minimum square error computed is compared with the predetermined threshold, to accept or reject the tested pattern.

#### 2)   Cross-Correlation

Cross-correlation is a measure of similarity of two waveforms. This process is considered as a Template Matching (TM) based approach. The cross-correlation coefficient should be high between the target and the reference utterances.

*3)* *Recovery of the full utterance:*

In this phase, the frames corresponding to the first and last syllables are concatenated to the nucleus in order to restore the complete target utterance. This is achieved through applying the template matching explained above.

The hybrid technique described in this paper was fully implemented in MATLAB 7.12.

### III. RESULTS AND DISCUSSION

The system performance is tested in two tracks. The first track is by applying it on Quranic words uttered by professionals readers, in this case a text search is accompanying the audio search. The rules of recitation of the Quran lead to consistent pronunciations. The system accuracy is evaluated by detecting correct and complete target words. In the second track, only audio search is applied on standard Arabic audio files representing lecture of Quran explanation, BBC Arabic news and a BBC Arabic interview.

#### A. Experimental setup

For the first track: More than three hundred utterances were used as reference patterns. Selected words or syllables were used as keywords and 15 hours' of audio files of Quran data were used for evaluation. Feature parameters used were 8 and 12 MFCCs (Mel-Frequency Cepstral Coefficients) and 12 LPCCs (Linear predictive Cepstral Coefficients). The algorithm was applied on words/ phrases from Quran context.

For the second track: Fifty utterances were used as reference patterns. Selected words were used as keywords and 60 minutes of audio files of Quran explanation "Tafseer", BBC Arabic news and a BBC Arabic interview were used for evaluation. Feature parameter used was 8 MFCCs (Mel-Frequency Cepstral Coefficients) because; it's the best feature according to the previous experiments. The algorithm was applied on words/ phrases from episodes of "Tafseet" and BBC News.

#### B. Experiments results of the first track

- In order to measure the value added by using text search before Audio Keyword Spotting, a Quranic audio file representing the first 42 verses from "*EL-Rahman*" Surah was searched to allocate the repeated utterance "آلَاء ربكما". The Audio Keyword Spotting allocated the utterances in time duration of 300 seconds. When text search was added, a search time of 30 seconds has been achieved with the same accuracy. So a reduction of 90% of the search time was achieved.

- Silence detection was performed in order to divide large audio files into smaller ones which make the audio search more efficient, as shown in Fig. 3.



Fig. 3. Detecting the duration of silences Inter-Verse in "سورة العصر"

- In another experiment the Keyword ["الحمد لله"; Al-hamdulellah] was allocated in a long utterance as shown in Fig. 4.



Fig. 4. Detect a Keyword in utterance.

The results obtained by applying the hybrid text/audio search on different Quranic utterances, using different sets of parameters, are summarized in the charts represented in Fig.5 and Fig.6.



Fig. 5. Percentage of Accuracy for Words.

| WORDS | إِبْرَاهِيمَ | الشمس | السموات | أحمد | الجبال | الصلاة | الفجرْ | باخع | الأخدود | محمد |
|---|---|---|---|---|---|---|---|---|---|---|
| 8 MFCC | 100 | 87 | 88 | 100 | 81 | 96 | 60 | 100 | 100 | 100 |
| 12 MFCC | 96 | 100 | 92 | 100 | 81 | 100 | 60 | 100 | 100 | 100 |
| 12 LPCC | 92 | 93 | 81 | 100 | 88 | 100 | 60 | 100 | 100 | 100 |



Fig. 6. Percentage of Accuracy for Phrases.

| PHRASES | أسألكمْ عليه من أجرْ | رب العَالمين | السمواتِ والأرضْ | خَالدينَ فيها | ذي القُربى |
|---|---|---|---|---|---|
| 8 MFCC | 100 | 96 | 84 | 100 | 100 |
| 12 MFCC | 100 | 100 | 84 | 100 | 100 |
| 12 LPCC | 100 | 100 | 84 | 92 | 100 |

Results obtained from the first set of experiments are presented in tables I, II, III. The uttered word search was performed using audio reference uttered by the same reader (*El-Hossary*) with different features and different coefficients.

TABLE I.  FIRST TRACK RESULTS OF (8 MFCC) FEATURE VECTOR FOR QURAN

| KEYWORD | TRANSCRIPT | TEMPLATES | ACCURACY% |
|---|---|---|---|
| إبْرَاهِيَم | Ebrahiem | 26 | 100 |
| ربِ العَالمين | Rab-Elalameen | 25 | 96 |
| الشمس | A-Shams | 30 | 87 |
| السموات | A-Samawat | 26 | 88 |
| السمواتِ والأرضْ | A-Samawat w-Al-Ard | 25 | 84 |
| الجبال | Al-jebal | 26 | 81 |
| الصلاة | Al-Salah | 25 | 96 |
| خَالدينَ فيها | Khaldeena-Feha | 25 | 100 |
| أسألكمْ عليه من أجرْ | Asalokom Alyeh men-Agr | 7 | 100 |
| ذي القُربى | Ze-lqurba | 5 | 100 |
| الفجرْ | Al-Fajr | 5 | 60 |
| أحمد | Ahmad | 1 | 67 |
| باخع | Bakhea | 2 | 100 |
| الأخدود | Al-Okhdod | 1 | 100 |
| محمد | Muhammad | 4 | 100 |
| TOTAL ACCURACY | | | 90.6% |

TABLE II.  FIRST TRACK RESULTS OF (12 MFCC) FEATURE VECTOR FOR QURAN

| KEYWORD | TRANSCRIPT | TEMPLATES | ACCURACY% |
|---|---|---|---|
| إبْرَاهِيَم | Ebrahiem | 26 | 96 |
| ربِ العَالمين | Rab-Elalameen | 25 | 100 |
| الشمس | A-Shams | 30 | 100 |
| السموات | A-Samawat | 26 | 92 |
| السمواتِ والأرضْ | A-Samawat w-Al-Ard | 25 | 84 |
| الجبال | Al-jebal | 26 | 81 |
| الصلاة | Al-Salah | 25 | 100 |
| خَالدينَ فيها | Khaldeena-Feha | 25 | 100 |
| أسألكمْ عليه من أجرْ | Asalokom Alyeh men-agr | 7 | 100 |
| ذي القُربى | Ze-lqurba | 5 | 100 |
| الفجرْ | Al-Fajr | 5 | 60 |
| أحمد | Ahmad | 1 | 100 |
| باخع | Bakhea | 2 | 100 |
| الأخدود | Al-Okhdod | 1 | 100 |
| محمد | Muhammad | 4 | 100 |
| TOTAL ACCURACY | | | 94.2% |

TABLE III.  FIRST TRACK RESULTS OF (12 LPCC) FEATURE VECTOR FOR QURAN

| KEYWORD | TRANSCRIPT | TEMPLATES | ACCURACY% |
|---|---|---|---|
| إبْرَاهِيَم | Ebrahiem | 26 | 92 |
| ربِ العَالمين | Rab-Elalameen | 25 | 100 |
| الشمس | A-Shams | 30 | 93 |
| السموات | A-Samawat | 26 | 81 |
| السمواتِ والأرضْ | A-Samawat w-Al-Ard | 25 | 84 |
| الجبال | Al-jebal | 26 | 88 |
| الصلاة | Al-Salah | 25 | 100 |
| خَالدينَ فيها | Khaldeena-Feha | 25 | 92 |
| أسألكمْ عليه من أجرْ | Asalokom Alyeh men-agr | 7 | 100 |
| ذي القُربى | Ze-lqurba | 5 | 100 |
| الفجرْ | Al-Fajr | 5 | 60 |
| أحمد | Ahmad | 1 | 100 |
| باخع | Bakhea | 2 | 100 |
| الأخدود | Al-Okhdod | 1 | 100 |
| محمد | Muhammad | 4 | 100 |
| TOTAL ACCURACY | | | 92.7% |

Table IV shows the best feature vector obtained from the first set of experiments presented in tables I, II, III.

TABLE IV.  FIRST TRACK RESULTS OF BEST FEATURE VECTOR FOR QURAN

| KEYWORD | TRANSCRIPT | TEMPLATES | BEST FEATURE | ACCURACY% |
|---|---|---|---|---|
| إبْرَاهِيَم | Ebrahiem | 26 | 8-MFCC | 100 |
| ربِ العَالمين | Rab-Elalameen | 25 | 12-MFCC | 100 |
| الشمس | A-Shams | 30 | 12-MFCC | 100 |
| السموات | A-Samawat | 26 | 12-MFCC | 92 |
| السمواتِ والأرضْ | A-Samawat w-Al-Ard | 25 | 8-MFCC | 84 |
| الجبال | Al-jebal | 26 | 12-LPCC | 88 |
| الصلاة | Al-Salah | 25 | 12-LPCC | 100 |
| خَالدينَ فيها | Khaldeena-Feha | 25 | 8-MFCC | 100 |
| أسألكمْ عليه من أجرْ | Asalokom Alyeh men-agr | 7 | 12-LPCC | 100 |
| ذي القُربى | Ze-lqurba | 5 | 8-MFCC | 100 |
| الفجرْ | Al-Fajr | 5 | 8-MFCC | 80 |
| أحمد | Ahmad | 1 | 8-MFCC | 100 |
| باخع | Bakhea | 2 | 8-MFCC | 100 |
| الأخدود | Al-Okhdod | 1 | 12-LPCC | 100 |
| محمد | Muhammad | 4 | 12-LPCC | 100 |
| TOTAL ACCURACY | | | | 97% |

Another set of experiments were performed to evaluate the effect of changing the reference reader, for the same utterance, which is referred to as cross-reader. Results obtained from the first set of experiments are presented in table V. The utterance to be allocated is pronounced by one reader (in this case *El-Hossary*) and the reference was recorded by another reader (in this case *El-Menshawy*) and vice versa. Promising results reached 72%.

TABLE V.    FIRST TRACK CROSS-READER BETWEEN UTTERANCE (EL-MENSHAWY) AND REFERENCE (EL-HOSSARY)

| KEYWORD | TRANSCRIPT | TEMPLATES | BEST FEATURE | ACCURACY % |
|---|---|---|---|---|
| ذي القُربى | Ze-lqurba | 5 | 8-MFCC | 80 |
| أسألكمْ عليه من أجرْ | Asalokom alyh men-agr | 7 | 12-LPC | 57 |
| الفجرْ | Al-fajr | 6 | 12-LPCC | 50 |
| الأخدود | Al-Okhdod | 1 | 12-MFCC | 100 |
| البروج | Al-Brouj | 1 | 8-MFCC | 100 |
| باخع | Bakhea | 2 | 12-LPCC | 100 |
| أحمد | Ahmad | 1 | 12-MFCC | 100 |
| محمد | Muhammad | 4 | 12-LPCC | 25 |
| الطامة | Al-Tamaa | 1 | 12-LPCC | 100 |
| TOTAL ACCURACY | | | | 71.2 % |

## C. Experiments results of the second track

In this track, experiments were performed on general Arabic episodes such as lecture of Quran explanation "Tafseer", BBC Arabic news and a BBC Arabic interview. In this case the audio search was conducted over the whole record, since there were no text scripts available. Results obtained from the first set of experiments are presented in tables VI, VII.

TABLE VI.    SECOND TRACK RESULTS OF (8 MFCC) FEATURE VECTOR FOR TAFSEER KEYWORDS

| KEYWORD | TRANSCRIPT | TEMPLATES | BEST FEATURE | ACCURACY% |
|---|---|---|---|---|
| الوسواس | Al-Waswas | 3 | 8- MFCC | 67 |
| شهر | Shahr | 4 | 8 -MFCC | 100 |
| الناس | A-Nas | 8 | 8 -MPCC | 100 |
| القدر | Al-qdr | 6 | 8 -MFCC | 67 |
| TOTAL ACCURACY | | | | 83.5% |

TABLE VII.    SECOND TRACK RESULTS OF (8 MFCC) FEATURE VECTOR FOR ARABIC NEWS KEYWORDS

| KEYWORD | TRANSCRIPT | TEMPLATES | ACCURACY% |
|---|---|---|---|
| الأردن | Al-Ordon | 4 | 100 |
| اللاجئين السوريين | Al-ajean A-Soreen | 2 | 100 |
| المتظاهرين | Al-motazahreen | 5 | 60 |
| وزارة الداخلية | Wzart Al-dakhelea | 4 | 50 |
| بورسعيد | Por-Saeed | 2 | 100 |
| المجلس العسكري | Al-Magles Al-Askary | 2 | 100 |
| ميدان التحرير | Medan Al-Tahrer | 2 | 100 |
| الأخبار السعيدة | Al-Akhbar Al-Saeda | 5 | 80 |
| هولاندية | Holandya | 3 | 100 |
| TOTAL ACCURACY | | | 87.8% |

## IV. CONCLUSIONS

In this work, a keyword spotting approach based on Template Matching was used to perform audio search for words/ phrases in audio files. The Audio Keyword Spotting is also used to allocate silence periods in audio files which results in dividing larger audio files into smaller ones. This division process improves the search process as it is performed in smaller audio files. Considering the audio files in Quran that have corresponding text files, a hybrid technique depending on both text search and audio keyword spotting is developed. This hybrid technique results is 97% accuracy when performing audio searching using audio reference of the same reader. 90% time reduction is achieved when compared with audio search only. This was accomplished using the sets of features (MFCCs and LPCCs). It was shown that the MFCC with an order 8 results in the best spotting accuracy. The accuracy of the developed spotting reached 72% when testing cross-readers utterances( the reference reader is tested against a different one).

Using the same audio keyword spotting technique to search in general audio files such as "News" and "preaches(Tafseer)" episodes, the recognition rates reached around 84% for preaches and around 88% for the same speaker in each test without the help of text search and with no recitation rules to control the speaker's pronunciation.

Despite the simplicity of the technique, it proves to be very efficient and shows high robustness to obtain high recognition rates under all circumstances.

## V. FUTURE WORK

Record a larger evaluation standard database, for different speakers and different environments, to get more test cases.

The system can be expanded to cover the whole Quran by complete the implementation for acoustical database of the rest recitation rules. This can achieve by manually segmenting the phonetic units of each rule from various referenced readers sounds.

In the updates of this system, we may use the resulting automatically detected Keywords in online process such as News and Arabic dialog programs for different speakers.

### REFERENCES

[1] A. Youssef, O. Emam." An Arabic TTS based on the IBM Trainable Speech Sythesizer." Department of Electronics & Communication Engineering, Cairo University, Giza, Egypt, 2004.

[2] Yung-Hwan Oh, Jeong-Sik Park and Kyung-Mi Park" Keyword Spotting in Broadcast News." Department of Electrical Engineering & Computer Science Korea Advanced Institute of Science and Technology, Daejeon, Korea, 2007.

[3] J. Junkawitsch, L. Neubauer, H. Höge, and G. Ruske "A New Keyword Spotting Algorithm with Pre-Calculated Optimal Thresholds" Technical University of Munich, Germany, 1996.

[4] A. Kandil, A. Bialy, S. Fawzi, M. Awaid. "Towards Speech Corpus for the Quran Using Keyword Spotting." M.S. thesis, Biomedical and systems Engineering Department, Cairo University, Giza, Egypt, 2013.

[5] Rabiner, L. and Juang, B. -H., Fundamentals of Speech Recognition, PTR Prentice Hall, San Francisco, NJ, 1993.

[6] Grangier D. and Bengio, S., "Learning the Inter-frame Distance for Discriminative Template-based Keyword", International Conference on Speech Communication and Technology (INTERSPEECH), 2007.

[7] Octavian Cheng, Waleed Abdulla, Zoran Salcic "Performance Evaluation of Front-end Processing for Speech Recognition Systems.", Electrical and Computer Engineering Department, Auckland University, New Zealand, 2005.

[8] Zhou Wang and Alan C. Bovik. "Mean Squared Error: Love It or Leave It? [A new look at signal fidelity measures]" IEEE Signal Processing Magazine, pp. 98-117, January 2009.

# A Tentative Analysis of the Rectangular Horizontal-slot Microstrip Antenna

Md. Tanvir Ishtaique ul Huque[1] and Md. Imran Hasan[2]

Department of Electronics and Telecommunication Engineering, Rajshahi University of Engineering & Technology

Rajshahi 6204, Bangladesh

*Abstract*—**In this paper, we have presented a new type of microstrip antenna mentioned as rectangular horizontal-slot patch antenna. Our main motto is to design a novel antenna which has the simplicity in structure and higher return loss. We have followed a tentative approach which leads us to an exceptional result, better than conventional one and the experimental outcomes result some guidelines for further practice. Here all of these antennas were analyzed by using GEMS (General Electro-Magnetic Solver) commercial software from 2COMU (Computer and Communication Unlimited).**

*Keywords—GEMS; Microstrip antenna; Rectangular horizontal slot antenna; Return loss; Slot antenna; antenna*

## I. Introduction

Microstrip patch antennas are popular, because they have some advantages due to their conformal and simple planar structure. They allow all the advantages of printed-circuit technology. A vast number of papers are available on the investigation of various aspects of microstrip antennas [1, 6, 7, 8, 10, 11, 12, 13, 14, 15, and 16].

The term "Microstrip" comes because the thickness of this metallic strip is in micro meter range. The key features of a microstrip antenna are relative ease of construction, light weight, low cost and either conformability to the mounting surface or, at least, an extremely thin protrusion from the surface [2, 5]. These criteria make it popular in the field of wireless communication. Microstrip antennas are the first choice for this high frequency band due to its light weight, low cost, and robustness.

Microstrip patch elements are available in various configurations [2, 3, 5]. But the most common is the rectangular patch element. In this paper we have presented a new type of rectangular patch antenna mentioned as rectanguar slot patch antenna. We have introduced slots in it and examined the effect of it. Here we have only considered the return loss as the performance parameter. Return loss is a measure of how well the antenna is matched or how much power is going to be used effectively [5].

A high return loss is always desirable. This paper focus the design procedure, characteristic and the corresponding performance analysis of both the conventional rectangular antenna and newly introduced rectangular horizontal-slot antenna and provides a mean to choose the effective one based on their performance parameter.

Here all of these antennas have been designed simulated by using the GEMS (General Electro-Magnetic Solver) version 7.71.01 simulator. GEMS package includes a time domain solver based on the parallel conformal FDTD.

The proposed antennas are designed by using Taconic TLY-5 dielectric substrate with permittivity, $\varepsilon_r$ =2.2 and height, h =1.588 mm. These designed antennas are promising to be a good candidate for the wireless applications due to the simplicity in structure, ease of fabrication and higher return loss.

## II. Antenna Configuration And Design

Microstrip patch antenna, illustrated in Figure 1, consists of very thin metallic strip (patch) placed on ground plane where the thickness of the metallic strip is restricted by $t << \lambda_0$ and the height is restricted by $0.0003\lambda_0 \leq h \leq 0.05\lambda_0$ [9, 13, 14]. The microstrip patch is designed so that its radiation pattern maximum is normal to the patch. For a rectangular patch, the length L of the element is usually $\lambda_0/3 < L < \lambda_0/2$ [9, 13, 14].

There are numerous dielectric substrates that can be used for the design of microstrip antennas and their dielectric constants are usually in the range of $2.2 \leq \varepsilon_r \leq 12$ [5, 9]. To implement the microstrip antennas usually Fr-4 ($\varepsilon_r$=4.9), Rogers TMM 4($\varepsilon_r$=4.5), Taconic TLY-5 ($\varepsilon_r$=2.2), Alumina (96%) ($\varepsilon_r$=9.4), Teflon(PTFE) ($\varepsilon_r$=2.08), Arlon AD 5 ($\varepsilon_r$=5.1) dielectric materials are used as the substrate [ 2, 5].



Fig. 1. Single element Rectangular microstrip patch antenna.

The Performance of the microstrip antenna depends on its dimension. Depending on the dimension the operating frequency, radiation efficiency, directivity, return loss and other related parameters are also influenced [3, 4]. Here, in this paper, a tentative analysis is made to find out the effect of introducing slots in the rectangular microstrip antenna and that leads to a comparative investigation between the rectangular microstrip antenna and rectangular horizontal-slot microstrip antenna.

For an efficient radiation a practical width of the Rectangular patch element becomes [2, 3, and 5]

$$w = \frac{1}{2 f_r \sqrt{\mu_0 \varepsilon_0}} \times \sqrt{\frac{2}{\varepsilon_r + 1}} \qquad (1)$$

And the length of the antenna becomes [2, 3, 5]

$$L = \frac{1}{2 f_r \sqrt{\varepsilon_{eff}} \sqrt{\varepsilon_0 \mu_0}} - 2\Delta L \qquad (2)$$

Where

$$\Delta L = 0.41 h \frac{\varepsilon_{eff} + 0.3}{\varepsilon_{eff} - 0.258} * \frac{\left(\frac{w}{h} + 0.264\right)}{\left(\frac{w}{h} + 0.8\right)} \qquad (3)$$

$$\varepsilon_{eff} = \frac{\varepsilon_r + 1}{2} + \frac{\varepsilon_r - 1}{2\sqrt{1 + 12\frac{h}{w}}} \qquad (4)$$

Where, λ is the wave length, $f_r$ (in Hz) is the resonant frequency, L and W are the length and width of the patch element, in mm, respectively and $\varepsilon_r$ is the relative dielectric constant.

In the following Fig. 2, Fig. 3 and Fig. 4 the antennas, mentioned as Antenna Type I, Antenna Type II and Antenna Type III respectively, are the single element rectangular microstrip antenna, where the quarter wavelength transformer method [2, 5, 15, 16] is used to match the impedance of the patch element with the transmission line. Here equation no 1 and 2 have been used to design all of these antennas.

Theoretically Antenna Type I, Antenna Type II and Antenna Type III have been designed to be operated in 8 GHz, 11 GHz and 12 GHz operating frequency respectively. In each Antenna type, there are two diffreent models- one is conventional rectangular microstrip patch antenna and another one is the rectangular-slot microstrip patch antenna. In the rectangular-slot microstrip patch antenna, each slot has 0.5 mm width, 0.5 mm separation distance from the both ends of the patch element and two successive slots are also separated by a distance of 0.5 mm.

## II. ANTENNA ANALYSIS AND SIMULATION RESULT

In this paper, it is considered that the substrate permittivity of the antenna is $\varepsilon_r$ = 2.2 (Taconic TLY-5) and substrate thickness is 1.588 mm.



Fig. 2. Single element rectangular microstrip patch antenna (a) without slot and with (b) one slot (c) two slots (d) three slots (e) four slots (f) five slots (g) six slots (h) seven slots (i) eight slots (j) nine slots (k) ten slots (l) eleven slots. (Antenna Type I).

Fig. 3. Single element Rectangular microstrip patch antenna (a) without slot and with (b) one slot (c) two slots (d) three slots (e) four slots (f) five slots (g) six slots (h) seven slots. (Antenna Type II).



Fig. 4. Single element Rectangular microstrip patch antenna (a) without slot and with (b) one slot (c) two slots (d) three slots (e) four slots (f) five slots (g) six slots. (Antenna Type III).

After simulation, as shown in Fig. 4, we found that, return loss is -33.06 dB at 7.5 GHz for the single element rectangular patch antenna. Simulation shows that the return loss is increased with introducing slot in the microstrip antenna and we get the maximum return loss -48.86 dB at 12.7 GHz operating frequency having three slots in the rectangular patch antenna.

Fig. 5 and Fig. 6 also show the same kind of manner as the Fig. 4. According to Fig. 5, return loss is -16.48 dB at 10.5 GHz for the single element rectangular patch antenna and because of introducing slot in the microstrip antenna; we get the maximum return loss -27.25 dB at 16.3 GHz operating frequency having two slots in the rectangular patch antenna.

Fig. 5.   Return loss pattern variation with respect to the frequency at different number of horizontal slots of Antenna Type I.



Fig. 6.   Return loss pattern variation with respect to the frequency at different number of horizontal slots of Antenna Type II.

Fig. 7. Return loss pattern variation with respect to the frequency at different number of horizontal slots of Antenna Type III.

Return loss pattern variation with respect to the frequency at different number of horizontal slots of Antenna Type III.

Fig. 6 states that return loss is -15.66 dB at 11.5 GHz for the single element rectangular patch antenna. Simulation also shows that the return loss is increased with introducing slot in the microstrip antenna and we get the maximum return loss -21.3 dB at 17.6 GHz operating frequency having three slots in the rectangular patch antenna.

In this paper, we have tried to find out a microstrip antenna having maximum return loss to get maximum efficiency and our investigation show that by introducing slot in the rectangular microstrip antenna we get greater return loss than the conventional rectangular patch antenna.

A comparative analysis of Fig. 4, Fig. 5 and Fig. 6 lead us to the Fig. 7, Fig. 8 and Fig. 9 which depict the effect of number of slots of the Antenna Type I, Antenna Type II and Antenna Type III respectively, from which we can easily understand that how many number of slots are necessary in each type of antenna to get the maximum return loss.

From Fig. 7, it is clear that for the Antenna Type I, having the slot of 0.5 mm width with the prescribed configuration; to get the maximum return loss three (3) slots are required. Fig. 8 and Fig. 9 also enunciate that for the Antenna Type II and Antenna Type III, having the slot of 0.5 mm width with the prescribed configuration, to get the maximum return loss two (2) slots and one (1) slot are required respectively.



Fig. 8. Variation of the Maximum Return Loss and the Operating frequency with respects to the number of slots of Antenna Type I



Fig. 9. Variation of the Maximum Return Loss and the Operating frequency with respects to the number of slots of Antenna Type II

Fig. 10. Variation of the Maximum Return Loss and the Operating frequency with respects to the number of slots of Antenna Type III

### III. MEASUREMENTS AND RESULT DISCUSSION

Our whole experimental outcomes conduct us to to the following Table I, Table II and Table III for three different types of antennas known as Antenna Type I, Antenna Type II and Antenna Type III respectively. These tables are the fruits, in brief, of our full tentative analysis.

Here we have followed a tentative approach to find out the

effect of slot in the single element rectangular patch antenna. Our experimental investigation leads us to the following outcomes depending on the data available in the Table I, Table II and Table III. All of these observations give us an optimum result about the selection of the slot in the patch element.

- Return loss increases with increasing antenna area or in turn slot length. Such as

For (11.6mm×14.8mm) Antenna Type I having slot length of 12.8mm Maximum Return loss is -48.86 dB.

For (8.1mm×10.8mm) Antenna Type II having slot length of 8.8mm Maximum Return loss is -27.25 dB.

For (7.3mm×9.9mm) Antenna Type III having slot length of 7.9mm Maximum Return loss is -21.3 dB.

- After observation, we get that the operating frequency of the microstrip antenna (without slot) is always 5.5 GHz less than the estimated one. Such as

For Antenna Type I, Maximum return loss is obtained at 7.5 Ghz, but it was designed for 8 GHz.

For Antenna Type II, Maximum return loss is obtained at 10.5 Ghz, but it was designed for 11 GHz.

For Antenna Type IIII, Maximum return loss is obtained at 11.5 Ghz, but it was designed for 12 GHz.

TABLE I. GEOMETRIC GESTALT AND SIMULATION RESULT OF THE "ANTENNA TYPE I"

| Antenna Length (mm) | Antenna Width (mm) | Slot Width (mm) | Slot Length (mm) | Number of Slots | Theoretically desired Frequency (GHz) | Obtained Frequency (GHz) in Simulation | Max. Return Loss. in Simulation (dB) |
|---|---|---|---|---|---|---|---|
| 11.6 | 14.8 | 0.5 | 12.8 | 0 | 8 | 7.55 | -33.06 |
| | | | | 1 | 8 | 12.6 | -45.17 |
| | | | | 2 | 8 | 12.7 | -37.83 |
| | | | | 3 | 8 | 12.7 | -48.86 |
| | | | | 4 | 8 | 12.8 | -33.9 |
| | | | | 5 | 8 | 12.5 | -19.09 |
| | | | | 6 | 8 | 12.5 | -27.08 |
| | | | | 7 | 8 | 12.6 | -30.11 |
| | | | | 8 | 8 | 12.6 | -31.3 |
| | | | | 9 | 8 | 12.5 | -34.77 |
| | | | | 10 | 8 | 12.5 | -34.6 |
| | | | | 11 | 8 | 12.5 | -33.6 |

TABLE II. GEOMETRIC GESTALT AND SIMULATION RESULT OF THE "ANTENNA TYPE II"

| Antenna Length (mm) | Antenna Width (mm) | Slot Width (mm) | Slot Length (mm) | Number of Slots | Theoretically desired Frequency (GHz) | Obtained Frequency (GHz) in Simulation | Max. Return Loss. in Simulation (dB) |
|---|---|---|---|---|---|---|---|
| 8.1 | 10.8 | 0.5 | 8.8 | 0 | 11 | 10.5 | -16.48 |
| | | | | 1 | 11 | 16.4 | -20.6 |
| | | | | 2 | 11 | 16.3 | -27.25 |
| | | | | 3 | 11 | 16.7 | -20.69 |
| | | | | 4 | 11 | 16.4 | -18.8 |
| | | | | 5 | 11 | 16.3 | -19.49 |
| | | | | 6 | 11 | 16.3 | -20.05 |
| | | | | 7 | 11 | 16.3 | -19.68 |

TABLE III.        GEOMETRIC GESTALT AND SIMULATION RESULT OF THE "ANTENNA TYPE III"

| Antenna Length (mm) | Antenna Width (mm) | Slot Width (mm) | Slot Length (mm) | Number of Slots | Theoretically desired Frequency (GHz ) | Obtained Frequency (GHz) in Simulation | Max. Return Loss. in Simulation (dB) |
|---|---|---|---|---|---|---|---|
| 7.3 | 9.9 | 0.5 | 7.9 | 0 | 12 | 11.5 | -15.66 |
| | | | | 1 | 12 | 17.6 | -21.3 |
| | | | | 2 | 12 | 17.8 | -19.29 |
| | | | | 3 | 12 | 17.5 | -17.31 |
| | | | | 4 | 12 | 17.5 | -18.92 |
| | | | | 5 | 12 | 17.5 | -19.75 |
| | | | | 6 | 12 | 17.6 | -20.11 |

- When we add slot in an Antenna, we find that the Maximum Return Loss (MRL) is shifted to a different operating frequency, depending on the number of slots. The shifted operating frequency becomes

$$F_M = F_I + F_S$$

Where
$F_M$ = Shifted operating frequency having the MRL of the Slotted Antenna.
$F_I$ = Estimated operating frequency for which the Antenna has been designed.
$F_S$ =The amount of frequency shifting which is maximum up to 5.8 GHz.
As an example, for the Antenna Type III, MRL -21.3dB was found in 17.6 Ghz operating frequency whereas the antenna was designed to operate at 12 GHz. Here, $F_M$=17.6 Ghz, $F_I$= 12 Ghz and $F_S$=5.6 Ghz.

- Here we have introduced the concept of slot in the microstrip antenna to get the maximum return loss and our experiment shows that we need to use maximum four(4) slots to get the maximum return loss. Such as

For Antenna Type I, when the number of slot is 3 the maximum return loss becomes -48.86 dB.

For Antenna Type II, when the number of slot is 2 the maximum return loss becomes -27.25 dB.

For Antenna Type III, when the number of slot is 1 the maximum return loss becomes -21.3 dB.

- Experimental statistics state that the difference of Maximum Return Loss between the slot and without slot antenna decreases with decreasing order of the microstrip antenna area. Such as,

For (11.6mm×14.8mm) Antenna Type I, Maximum Return loss is -33.06 dB and after introducing slot when there are three(3) slots in the same antenna the maximum return loss becomes -48.86 dB. The difference in Return Loss between two conditions is 15.8 dB.

For (8.1mm×10.8mm) Antenna Type II, Maximum Return loss is -16.48 dB and after introducing slot when there are two(2) slots in the same antenna the maximum return loss becomes -27.25 dB. The difference in Return Loss between two conditions is 10.77 dB.

For (7.3mm×9.9mm) Antenna Type III, Maximum Return loss is -15.66 dB and after introducing slot when there is one (1) slot in the same antenna the maximum return loss becomes

-21.3 dB. The difference in Return Loss between two conditions is 5.64 dB.

## IV. CONCLUSION

The unique feature of this microstrip antenna is its simplicity to get higher performance. In many applications, basically in wireless communication, it is necessary to design antennas with very high Return loss to meet the demand of long distance communication and the most common configuration to satisfy this demand is the microstrip antenna.

In our ongoing investigation, we have tried to find out the effect of slot in the microstrip antenna and that lead us to a remarkable result. Here we have invented a novel antenna, mentioned as rectangular horizontal slot patch antenna, which shows much higher return loss than the conventional rectangular patch antenna. In this paper, we have only focused some simulated result of introducing slots in the rectangular patch antenna where the slot width is fixed and the separation distance between two slots are also remained constant. We have experimented in three different cases and their outcomes are about similar.

We are still working on the rectangular horizontal as well as vertical slot patch antenna and investigating the effect of the variation of the slot size, separation distance between two slots in the microstrip antenna, with that inspiration to open a new horizon in the field on Antenna Technology.

REFERENCES

[1] R. J. Mailloux, J. F. McIlvenna, N. P. Kernweis, "Microstrip array technology", IEEE Trans. Antenna Propagation Magazine, vol. 29, no. 1, pp. 25-27, 1981.

[2] C. A. Balanis, Antenna Engineering, 2nd ed., Willey, 1982.

[3] T. A. Millikgan, Modern Antenna Design, 2nd ed., IEEE Press, John Wiley & Sons inc., 2007.

[4] M. I. Skolnik, Introduction to RADAR System, 3rd ed., McGraw Hill Higher Education, 2000.

[5] R. Garg, P. Bhartia, I. Bahl, A. Ittipiboon, Microstrip Antenna Design Handbook, Artech House inc., 2001.

[6] W. L. Stutzman, "Estimating directivity and gain of antennas", IEEE Antennas and Propagation Magazine, Vol. 40, No. 4,pp 7-11, August, 1998.

[7] H. J. Visser, Array and Phased Array Antenna Basics, John Wiley & Sons Ltd., 2005.

[8] Muhammad Mahfuzul Alam, Md. Mustafizur Rahman Sonchoy, and Md. Osman Goni, "Design and Performance Analysis of Microstrip Array Antenna", Progress In Electromagnetic Research Symposium Proceedings, Moscow, Russia, Aug. 18-21, 2009.

[9] Md. Shihabul Islam and Md. Tanvir Ishtaique-ul Huque, "Design and Performance Analysis of Microstrip Array Antenna", B.Sc. Engineering thesis, Dept. of ETE, Rajshahi University Of Engineering & Technology(RUET), Rajshahi, Bangladesh, April, 2010.

[10] K. Shambavi, C. Z. Alex, T. N. P. Krishna, "Design and Analysis of High Gain Milimeter Wave Microstrip Antenna Array for Analysis of High Gain Millimeter Wave Microstrip Anteanna Array for Wireless Application", J. of Applied Theoretical and Information Technology(JATIT), 2009.

[11] Asghar Keshtkar, Ahmed Keshtkar and A. R. Dastkhosh, "Circular Microstrip Patch Array Antenna for C-Band Altimeter System", Int. J. of Antenna and Propagation, article ID 389418, Nov., 2007. (doi:10.1155/2008/389418)

[12] M. F. Islam, M. A. Mohd. Ali, B. Y. Majlis and N. Misran, "Dual Band Microstrip Patch Antenna for Sar Applications", Australian Journal of Basic and Applied Sciences, 4(10): 4585- 4591, 2010.

[13] Md. Tanvir Ishtaique-ul Huque, Md. Al-Amin Chowdhury, Md. Kamal Hosain, Md. Shah Alam, "Performance Analysis of Corporate Feed Rectangular Patch Element and Circular Patch Element 4x2 Microstrip Array Antennas", Int. J. of Advanced Computer Science and Applications(IJACSA), vol. 2, no.8, pp. 16-21, 2011.

[14] Md. Tanvir Ishtaique-ul Huque, Md. Kamal Hosain, Md. Shihabul Islam, Md. Al-Amin Chowdhury, "Design and Performance Analysis of Microstrip Array Antennas with Optimum Parameters for X-band Applications", Int. J. of Advanced Computer Science and Applications (IJACSA), vol. 2, no. 4, pp. 81-87, 2011.

[15] Md. Tanvir Ishtaique-ul Huque, Md. Kamal Hosain, Mst. Fateha Samad, Muhammed Samsuddoha Alam , "Design and Simulation of a Low-cost and High Gain Microstrip Patch Antenna Arrays for the X-band Applications", Int. Conf. on Network Communication and Computer (ICNCC 2011), India, pp. 548- 552, March 19-20, 2011.

[16] Md. Tanvir Ishtaique-ul Huque, Md. Shihabul Islam, Mst. Fateha Samad, Md. Kamal Hosain, "Design and Performance Analysis of the Rectangular Spiral Microstrip Antenna and Its Array Configuration", 9th Int. Symp. on Antenna Propagation & EM Theory (ISAPE 2010), China , pp. 219-221, Nov 29 – Dec 2, 2010.

# TCP I-Vegas in Mobile-IP Network

Nitin Jain
Electronics & Communication Engineering
BBDESGI
Lucknow, India

Dr. Neelam Srivastava
Electronics Engineering
IET
Lucknow, India

*Abstract*—Mobile Internet Protocol (Mobile-IP or MIP) provides hosts with the ability to change their point of attachment to the network without compromising their ability to communicate. However, when TCP Vegas is used over a MIP network, its performance degrades because it may respond to a handoff by invoking its congestion control algorithm. TCP Vegas is sensitive to the change of Round-Trip Time (RTT) and it may recognize the increased RTT as a result of network congestion. This is because TCP Vegas could not differentiate whether the increased RTT is due to route change or network congestion. This paper presents a new and improved version of conventional TCP Vegas, which we named as TCP I-Vegas (where "I", stands for Improved). Vegas performs well when compared to Reno but when sharing bandwidth with Reno its performance degrades. I-Vegas has been designed keeping in mind that whenever TCP variants like Reno has to share the bandwidth with Vegas then instead of using Vegas, if we use I-Vegas then the loss which Vegas would have to bear will not be more. We compared the performance of I-Vegas with Vegas in MIP environment using Network Simulator (NS-2). Simulation results show that I-Vegas performs better than Vegas in terms of providing better throughput and congestion window behavior.

*Keywords—TCP Vegas; Mobile-IP; NS-2*

## I. INTRODUCTION

A large number of heterogeneous computer networks interconnected together using TCP/IP protocol suite (Transmission Control Protocol/Internet Protocol) forms Internet. With the fast prevalence of Internet users demand the mobility of hosts, i.e., they expect that the hosts can change their locations continuously without interrupting current communication sessions.

TCP is a reliable, connection-oriented protocol that ensures in-order delivery of a byte stream supplied by an application. It provides reliable service by implementing flow control, error detection, error recovery, in-order delivery, and removing duplicate data. Both the sending and the receiving node must keep state to support reliable delivery, therefore a connection is setup before data are transferred.

MIP provides hosts with the ability to change their point of attachment to the network without compromising their ability in communications. The mobility support provided by MIP is transparent to other protocol layers so as not to affect those applications which do not have mobility features. MIP introduces three new entities required to support the protocol: the Home Agent (HA), the Foreign Agent (FA) and the Mobile Node (MN). The MIP Working Group of the Internet Engineering Task Force (IETF) has compiled a series of Internet Drafts and Request for Comments (RFC) to define

MIP for providing an economical solutions which implements mobility support over the existing Internet infrastructure.

There are several problems of using TCP Vegas in a MIP network. Since TCP Vegas is tuned to perform well in traditional wired networks in which most packet losses are due to congestion. However, in a wireless mobile network, packet losses usually occur due to either random loss or handoff. After a handoff, the throughput of TCP Vegas may be decreased due to a longer BaseRTT of the new routing path, which is usually caused by either triangular routing or route optimization.

In this paper, we present a new and improved version of conventional TCP Vegas which we named as TCP I-Vegas (where "I" stands for Improved). I-Vegas proves to be better in terms of throughput and congestion window behavior, when compared with conventional Vegas. Simulation results proved that our proposed new and improved I-Vegas performs better than Vegas in MIP wired-cum-wireless network.

The rest of paper is organized as follows: Section II presents background of TCP Vegas and Mobile-IP networks. Section III provides issues related with TCP Vegas. Section IV gives algorithm of TCP I-Vegas which we have made in order to improve the performance of TCP Vegas. Section V presents simulation results and discussions. We conclude in Section VI.

## II. BACKGROUND: MOBILE-IP & TCP VEGAS

### A. Mobile-IP

In order to achieve the mobility function, the Internet Protocol (IP) has extended to become the Mobile Internet Protocol (Mobile IP or MIP). MIP provides hosts with the ability to change their point of attachment to the network without compromising their ability to communicate. The mobility support provided by MIP is transparent to the other protocol layers so as not to affect the operation of applications which do not have the mobile capability. Among various IP mobility proposals, Mobile IPv4 [1] & [2] is the oldest and probably the most widely known mobility management proposal with IP. MIPv4 introduces three new entities required to support the protocol: the Home Agent (HA), the Foreign Agent (FA) and the Mobile Node (MN). HA and FA are introduced for mobility management. The basic idea is to use an authenticated registration procedure between a MN and a HA in its home network, and via a FA while MN is visiting a foreign network. Each time a mobile host connects to a network at a new location, it will obtain a temporary address, called Care-of Address (COA) from a foreign agent in the local network. Then the mobile host must inform its home

agent of the new address by a registration procedure, which begins when the mobile host, possibly with the assistance of the foreign agent, sends a registration request with the COA. When the home agent receives this request, it may typically add the necessary information to its routing table, approve the request, and send a registration reply back to the mobile host. Further information on MIP functionality can be found in [1] & [2].

### B. TCP Vegas

TCP Vegas, a conservative algorithm, which is delay based, first proposed by L.S. Brakmo and L.L. Peterson [3] ensures end-to-end integrity of data transfer, while IP performs datagram routing and internetworking functions. It achieves 37 to 71 percent higher throughput than most used TCP version called TCP Reno [4], which is loss based. S. Ahn, P.B. Danzig, Z. Liu and L. Yan [5] have evaluated the performance of Vegas and shown that it does achieve higher efficiency than Reno and causes much less packet retransmissions. However, they have also observed that Vegas when competing with other TCP variants like Reno, it does not receive a fair share of bandwidth, i.e., TCP Reno connections receive about 50 percent higher bandwidth. This incompatibility property is analyzed also by J. Mo and J. Walrand [6]. They show that due to the aggressive nature of Reno, when the buffer sizes are large, Vegas loses to Reno that fills up the available buffer space, forcing Vegas to back off. Hence, there is a need to improve the performance of Vegas, which is a conservative algorithm, so that whenever it shares the bandwidth with other TCP variants like Reno or New Reno [7], the loss which conventional Vegas bears should not be more.

TCP throughput is inversely related to RTT, Vegas measure the difference between the expected and the actual throughput. The idea is that the actual throughput should match the expected throughput if there is no congestion along the network path. A lower actual throughput indicates increased delay, and hence congestion, on the network path. Similar to Reno, Vegas has slow start and congestion avoidance modes.

### C. 1) Slow-Start

During slow-start, Vegas maintains the threshold $\gamma$ (the value of $\gamma$ is generally set to 1). As long as *diff*, when comparing *expected_thruput* and *actual_thruput* is less than $\gamma$ it increases the congestion window by 1 packet every other round trip time, rather than every RTT. Hence, during slow start the Vegas congestion window grows exponentially, though at a slower rate than in TCP Reno. At this point, Vegas needs correction so that it can be made somewhat aggressive.

When either the congestion window reaches the *slow start threshold (ssthresh)* or *diff* is larger than $\gamma$, Vegas enters the congestion avoidance. Upon exiting slow-start, Vegas decreases the congestion window by one eighth of its current size in order to ensure that the network does not remain congested.

### D. 2) Congestion-Avoidance

During congestion-avoidance, Vegas maintains two threshold values $\alpha$ and $\beta$ (the value of $\alpha$ and $\beta$ are usually set as 1 and 3 respectively). The adjustment of congestion window is done based on the value of *diff* given as follows:

$$cwnd = \begin{cases} cwnd + 1 & \text{if diff} < \alpha \\ cwnd - 1 & \text{if diff} > \beta \\ cwnd & \text{otherwise} \end{cases}$$

Where,
$diff = (expected\_thruput - actual\_thruput).base\_RTT$

$expected\_thruput = cwnd/base\_RTT$, where *cwnd* is the current congestion window size and *base_RTT* is the minimum round trip time of that connection.

$actual\_thruput = cwnd/RTT$, where RTT is the actual round trip time

Vegas tries to keep at least $\alpha$ packets but no more than $\beta$ packets in the queues. Roughly speaking, $\alpha$ and $\beta$ in Vegas represent respectively the minimum and the maximum number of packets the source can pipe in the network buffers; therefore $\alpha$ and $\beta$ represent the aggressiveness degree of the TCP Vegas sources. The higher their value, the more Vegas approaches the behavior of Reno. Vegas always attempts to detect and utilize the extra bandwidth whenever it becomes available without congesting the network. This mechanism is fundamentally different from that used by Reno. It always updates its window size to guarantee full utilization of available bandwidth, leading to constant packet losses, whereas Vegas does not cause any oscillation in window size once it converges to an equilibrium point.

In congestion avoidance phase, two changes can be made in the algorithm of Vegas. Firstly, the values of $\alpha$ and $\beta$ can be increased, because the aim is to make the algorithm of Vegas more aggressive. Secondly, when $\alpha < diff < \beta$ the size of the congestion window instead of keeping same, can be increased so that it will share the bandwidth more fairly as compared to other variants of TCP.

### E. 3) Loss Recovery

A packet loss can be detected via time out expiration or via three duplicated ack*s*. In the first case, the *ssthresh* is set to half of the current congestion window value, the congestion window is set to 2, and Vegas performs again the *slow-start*. In second case, when Vegas source receives three duplicate acks, it performs *Fast Retransmit* and *Fast Recovery* as Reno does. Actually, Vegas develops a more refined fast retransmit mechanism based on a fine-grain clock. After fast retransmit Vegas sets the congestion window to ¾, instead of ½ of the current congestion window and performs again the congestion avoidance algorithm.

## III. ISSUES WITH TCP VEGAS

### A. Fairness

Vegas uses a conservative algorithm to decide how and when to vary its congestion window. Reno, in an effort to fully utilize the bandwidth, continues to increase the window size until a packet loss is detected. Thus, when TCP Vegas and Reno connections shares a bottleneck link, Reno uses up most

of the link and router buffer space. Vegas, interpreting this as a sign of congestion, decreases its congestion window, which leads to an unfair sharing of available bandwidth in favor of Reno. This unfairness worsens when router buffer sizes are increased. G. Hasegawa, K. Kurata, M. Murata [8] proposed TCP Vegas$^+$ as a method to tackle Vegas's fairness issue. However, Vegas$^+$ assumes that an increase in the RTT value is always due to the presence of competing traffic and rules out other possibilities like rerouting. We feel that this is not a reasonable assumption. Furthermore, performance of Vegas$^+$ depends on the choice of optimal value for the new parameter $Count_{max}$ introduced in the protocol, which is an open question. G. Hasegawa, K. Kurata, M. Murata [6] and Raghavendra and Kinicki [9] showed that by using RED routers in place of the tail-drop routers, the fairness between Vegas and Reno can be improved to some degree. But there exists an inevitable trade-off between fairness and throughput, i.e. if the packet dropping probability of RED is set to a large value, the throughput share of Vegas can be improved, but the total throughput is reduced. In [10-11] Feng, Vanichpun and Weigle showed that choosing values of $\alpha$ and $\beta$ as a function of the buffer capacity of the bottleneck router could improve the fairness condition. However, they do not propose any mechanism to measure this buffer capacity and to set appropriate values for $\alpha$ and $\beta$.

### B. Rerouting

In Vegas, the parameter *baseRTT* denotes the smallest round-trip delay the connection has encountered and is used to measure the expected throughput. When rerouting occurs in between a connection, the RTT of a connection can change. When the new route has a longer RTT, the Vegas connection is not able to deduce whether the longer RTTs experienced are caused by congestion or route change. Without this knowledge, TCP Vegas assumes that the increase in RTT is due to congestion along the network path and hence decreases the congestion window size [12].

This is exactly opposite of what the connection should be doing. When the propagation delay increases, the bandwidth–delay product ($bw*d$) increases. The expression ($cwnd-bw*d$) gives the number of packets in the buffers of the routers. Since the aim of Vegas is to keep the number of packets in the router buffer between $\alpha$ and $\beta$, it should increase the congestion window to keep the same number of packets in the buffer when the propagations delay increases. In [12] the authors also proposed a modification to the Vegas to counteract the rerouting problem by assuming any lasting increase in RTT as a sign of rerouting. Besides the fact that this may not be a valid assumption in all cases, several new parameters $K$, $N$, $L$, $\delta$ and $\gamma$ were introduced in this scheme and finding appropriate values for these variables remain an unaddressed problem.

### IV. TCP I-VEGAS

The algorithm of Vegas required making it little bit aggressive from conservative so that when compared with other TCP variants like Reno it should perform better than the conventional Vegas.

Modifications in Vegas has been confined to the sender side only because of this our I-Vegas with proposed changes is easy to implement.

Modifications does not introduce any further thresholds, generally hard to set, since it is completely adaptive to the status of the network; in this prospect our I-Vegas with proposed changes appears to be more efficient.

I-Vegas, behavior is not much different from that of the original Vegas in presence of other Vegas sources; so it is able to preserve all the nice features of the original Vegas: good throughput and goodput performance and ability in network congestion avoidance.

### A. Algorithm

Following changes we have made in the algorithm of Vegas in order to make it more aggressive so that its performance get improved as compared to Vegas and it will fairly share the bandwidth when competing with other TCP variants like Reno.

During Slow-Start, we change the *cwnd* of Vegas more aggressively as Reno does.

In the case of rerouting, Vegas should not decrease its *cwnd*, rather to increase the thresholds $\alpha$ and $\beta$ to 3 and respectively.

During RTO and on reception of Three dup ACKs, $\alpha$ and $\beta$ are again set to 1 and 3 respectively.

During congestion avoidance, when *diff* lies between $\alpha$ and $\beta$, instead of keeping *cwnd* unchanged, Vegas should change its *cwnd* as it is changing when $diff < \alpha$.

### V. SIMULATION RESULTS AND DISCUSSION

We have created wired-cum-wireless MIP environment in NS-2 [13] and compared the parameters like throughput and congestion window behavior at different packet error probabilities.

### A. Network Topology

Fig. 1 shows the network topology which is a wired-cum-wireless MIP network. In fig. 1, node 0 and node 1 are W(0) and W(1) wired nodes respectively, node 2 and node 3 are base station nodes behaves like a HA and FA respectively and node 4 behaves like MN that moves between its HA and FA. Table I gives the details. We set up a TCP flow between node 0 to node 4 i.e. between W(0) and MH. As MH moves out from the domain of its HA, into the domain of FA, we observe how packets destined for MH is redirected by its HA to the FA as per MIP protocol definitions.

Fig. 1.    Wired-cum-Wireless MIP Network

TABLE I.          NODE DETAILS

| Node | Nature | TCP Connection |
|------|--------|----------------|
| Node 0 | Wired Node, W(0) (Source Node) | |
| Node 1 | Wired Node, W(1) | |
| Node 2 | Base Station Node Home Agent (HA) | Vegas/I-Vegas |
| Node 3 | Base Station Node Foreign Agent (FA) | |
| Node 4 | Mobile Node (MN) (Sink Node) | |

### B.  Comparison Curves for TCP Vegas and TCP I-Vegas

Fig. 2 to 9 shows the comparison curves in terms of congestion window behavior for TCP Vegas and TCP I-Vegas at different error probabilities. Similarly, fig. 10 to 17 shows the comparison curves in terms of throughput of TCP Vegas and TCP I-Vegas at different error probability.

Figure shows that I-Vegas performs better than Vegas in terms of both congestion window behavior and throughput at different error probabilities.



Fig. 2.    Congestion Window for TCP Vegas at 0% Error



Fig. 3.    Congestion Window for TCP I-Vegas at 0% Error



Fig. 4.    Congestion Window for TCP Vegas at 1% Error



Fig. 5.    Congestion Window for TCP I-Vegas at 1% Error



Fig. 6.    Congestion Window for TCP Vegas at 5% Error

Fig. 7.   Congestion Window for TCP I-Vegas at 5% Error



Fig. 8.   Congestion Window for TCP Vegas at 10% Error



Fig. 9.   Congestion Window for TCP I-Vegas at 10% Error



Fig. 10. Throughput of TCP Vegas at 0% Error



Fig. 11. Throughput of TCP I-Vegas at 0% Error



Fig. 12. Throughput of TCP Vegas at 1% Error



Fig. 13. Throughput of TCP I-Vegas at 1% Error



Fig. 14. Throughput of TCP Vegas at 5% Error

Fig. 15. Throughput of TCP I-Vegas at 5% Error



Fig. 16. Throughput of TCP Vegas at 10% Error



Fig. 17. Throughput of TCP I-Vegas at 10% Error

## VI. CONCLUSION & FUTURE SCOPE

In this paper, we have proposed a modified algorithm of Vegas and named it as I-Vegas, where "I" stands for "improved". We have also shown that making the algorithm of Vegas from conservative to somewhat aggressive, the performance of I-Vegas becomes much better than conventional Vegas. Simulation results proved that performance of I-Vegas in terms of av. throughput and congestion window behavior becomes better than Vegas in MIP network.

Mobile IP is a newly defined protocol which supports mobile users but also is compatible with the current IP. It is still in the process of being standardized, and there are still many items that need to be worked on and enhanced, such as the security issue and the routing issue. We are working on these issues.

REFERENCES

[1] C. Perkins, 'IP Mobility Support', IETF RFC 3220, January 2002

[2] C. Perkins, 'IP Mobility Support for IPv4', IETF RFC 3344, August 2002

[3] L. S. Brakmo, L. L. Peterson, TCP Vegas: end-to-end congestion avoidance on a global Internet, IEEE Journal on Selected Areas in Communications,Vol.13, No.8, October 1995.J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3$^{rd}$ ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.

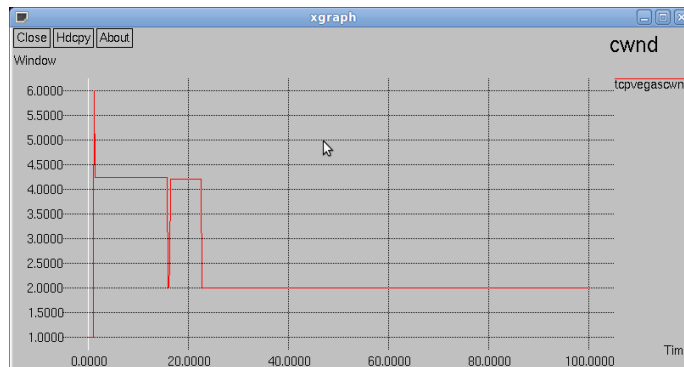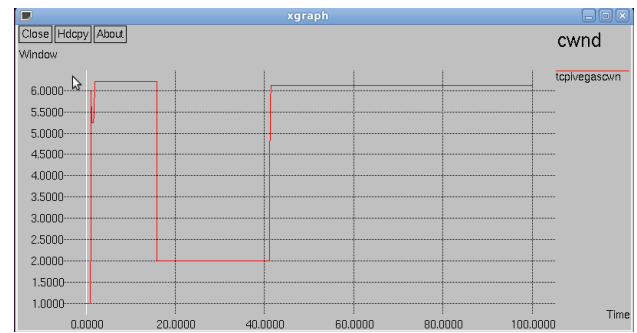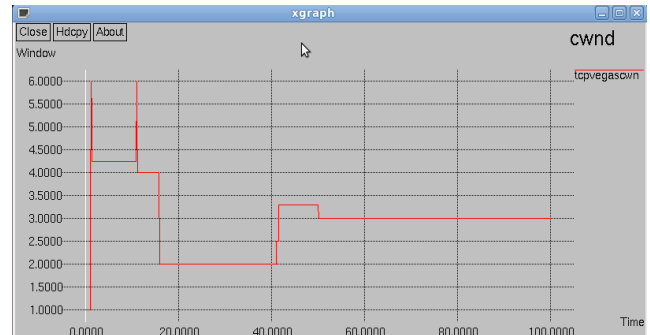[4] R. W. Stevens, TCP/IP Illustrated, Vol. I, The Protocols, Addison-Wesley, U.S.A., 1994.

[5] S. Ahn, P.B. Danzig, Z. Liu, and L. Yan, Evaluation of TCP Vegas: Emulation and Experiment. IEEE Transactions on Communications, 25(4):185-95, Oct 1995.

[6] J. Mo and J. Walrand, Fair End-to-end Window-based Congestion Control, SPIE '98 International Symposium on Voice, Video, and Data Communications, Nov. 1998.

[7] S. Floyd, T. Henderson, The NewReno modification to TCP's fast recovery algorithm, RFC 2582 April 1999.

[8] G. Hasegawa, K. Kurata, M. Murata, Analysis and improvement of fairness between TCP Reno and Vegas for deployment of TCP Vegas to the internet, Proceedings of the IEEE International Conference on Network Protocols (ICNP 2000) November 2000.

[9] A.M. Raghavendra, R.R. Kinicki, A simulation performance study of TCP Vegas and random early detection, Proceedings of IPCCC'99, February 1999 pp. 169–176.

[10] E. Weigle, W. Feng, A case for TCP Vegas in high-performance computational grids, Proceedings of Ninth International Symposium on High Performance Distributed Computing August 2001.

[11] W. Feng, S. Vanichpun, Enabling compatibility between TCP Reno and TCP Vegas, IEEE   Symposium on Applications and the Internet (SAINT 2003) January 2003.

[12] R.J. La, J. Walrand, V. Anantharam, Issues in TCP Vegas, July 1998.

[13] The Network Simulator - NS-2. URL: http://www.isi.edu/nsnam/ns/index.html.

# Complexity of Network Design for Private Communication and the P-vs-NP Question

Stefan Rass

Alpen-Adria University, System Security Group, Department of Applied Informatics,
Universitaetsstrasse 65-67, 9020 Klagenfurt, Austria
stefan.rass@aau.at

*Abstract*—We investigate infeasibility issues arising along network design for information-theoretically secure cryptography. In particular, we consider the problem of communication in perfect privacy and formally relate it to graph augmentation problems and the P-vs-NP-question. Based on a game-theoretic privacy measure, we consider two optimization problems related to secure infrastructure design with constraints on computational efforts and limited budget to build a transmission network. It turns out that information-theoretic security, although not drawing its strength from computational infeasibility, still can run into complexity-theoretic difficulties at the stage of physical network design. Even worse, if we measure (quantify) secrecy by the probability of information-leakage, we can prove that approximations of a network design towards maximal security are computationally equivalent to the exact solutions to the same problem, both of which are again equivalent to asserting that $P = NP$. In other words, the death of public-key cryptosystems upon $P = NP$ may become the birth of feasible network design algorithms towards information-theoretically confidential communication.

*Index Terms*—Complexity; NP; Privacy; Security; Game Theory; Graph Augmentation

## I. INTRODUCTION

Encryption is a standard mean to establish private communication channels. Mostly, security rests on intractability assumptions (as for public-key cryptography) or empirical investigations (as for many symmetric encryptions). This intractability-based paradigm is opposed by techniques that use properly designed communication infrastructures to provide confidential data transmission channels. Notable examples of the latter are quantum key distribution (QKD) [1], [2] or multipath transmission (MPT) [3], [4], [5], [6], [7]. Contrary to conventional cryptography, these techniques do not hinge on computational intractability, whose related assumptions may become invalidated by increasing computational power, novel computer architectures (such as quantum- or DNA-computing [8], [9]), or new scientific discoveries (e.g., if $P = NP$, then most public-key cryptography is essentially insecure). Such resilience is the main motivation to look at quantum- or MPT techniques. However, the price for independency on intractability is often expensive infrastructure design, whose complexity-theoretic quantification is our goal in this work. Specifically, we investigate the (in)tractability of network graph design for the sake of running secure multipath transmission (which QKD also requires to achieve end-to-end

security from point-to-point unless quantum repeaters become reality [10]).

### A. Related Work and Contribution

In the quantum cryptography area, the problem network topology design to optimally support QKD has received attention in [11], [12], [13], [14]. Such considerations are justified and substantiated by previous findings [3] that multipath transmission is actually a *necessity* for confidential conversation (cf. theorem II.4) in the absence of classical cryptography or special-purpose channels (say quantum or wire-tap [15]). On the pure classical road, [4], [5] have identified graph connectivity as a necessary and/or sufficient criterion for secure communication. Related protocols like [6] then simply *presume* multiple paths to be available in a network infrastructure; a luxury that hardly any real-life network will provide. More importantly, most of the prior literature on MPT neglects complexity issues that arise in the necessary network construction. That gap motivated this work, as it poses the question for the minimal extension to a given graph to permit MPT in the sense as [6], [5], [7] and others attempt it. References [12], [13], [14] studied and classified the problem as at least NP-hard, which in turn motivates our search for *approximations* rather than exact solutions.

The contribution of this article is the unfortunate observation that even finding an approximate network design is already equivalent to proving that $P = NP$. While the problem of whether one can build a secure cryptosystem on the assumption that $P \neq NP$ is still unanswered ([16] provides an interesting discussion, unfortunately leaving the initial question essentially open), the confidence in the strength of nowadays public-key encryption seems well justified, based on the evidence at hand. Still, the work of [17] presents evidence *against* the well-established conjecture that one-way permutations (based on computational intractability) alone would suffice to set up a secret key agreement. We approach the same problem here via graph-connectivity based techniques (i.e., multipath transmission).

Hence, insofar secure communication avoids intractability by switching from encryption to multipath transmission based techniques (which also covers some implementations of quantum networks), intractability arises again, yet only in a different form. The good news, detailed in the concluding

section, is nevertheless the observation that for secure communication, we can safely use encryption if we assume $P \neq NP$, or otherwise construct network infrastructures for perfectly secure multipath transmission, which is feasible if ultimately $P = NP$ is proven.

### B. Organization

In order to make this work as self-contained as possible, we use Section II to introduce the notation, network, adversary and security models. Subsection II-D sketches the general approach to private communication by MPT, upon which the game-theoretic privacy measure is defined in Section II-E. The network design problems are stated in Section III, with the analysis and main results following in Section IV.

## II. MODELS AND NOTATION

Vectors are printed as bold-face letters, complexity classes are written in small caps, sets are denoted by upper-case letters, matrices are upper-case bold-printed. For a discrete set $X$, we write $|X|$ for its cardinality. Whenever $x$ is a string representation (encoding) of a problem, we write $|x|$ to denote its length, and whenever $x$ is a real variable, then $|x|$ is its absolute value. The distinction will always be clear from the context. The symbol $\text{poly}(n)$ denotes an existing yet not further specified polynomial in the given variable (or expression) $n$.

### A. Network Model

Let the network infrastructure consist of a set of $V$ *devices*, and a set $E \subseteq V \times V$ of (bidirectional) communication channels between these devices. Without loss of generality, we can assume that channels cannot be attacked, because a vulnerable channel $u-v$ can be emulated by adding an intermediate vulnerable device $w$ and inserting the two (invincible) channels $u-w$ and $w-v$ to the network model. Our representation for a network infrastructure is thus an undirected graph $G(V, E)$, where $V$ is the set of nodes (devices) and $E$ is the set of edges (point-to-point connections).

Let $s, t$ be two distinct nodes in the graph $G$. An *s–t-path* $\pi$ in $G$ is a set of consecutive vertices starting at $s$ and ending in $t$. We denote the set of vertices in $\pi$ as $V(\pi)$. Two *s–t*-paths $\pi_1, \pi_2$ are said to be *node-disjoint*, if their only common points are $s, t$, i.e. if $V(\pi_1) \cap V(\pi_2) = \{s, t\}$. The *s–t-vertex connectivity* of $G$ is the cardinality of the smallest set of nodes whose removal renders $s$ unreachable from $t$ in $G$. The *vertex connectivity* of $G$ is the size of the smallest set of nodes such that after deletion, the graph becomes either disconnected or trivial [18]. We write $G(V \setminus U, E)$ to denote the subgraph induced by $V \setminus U$ and the remaining edges in $E$. We say that a graph is *k-connected*, if its vertex connectivity is $k$. The vertex-connectivity number is directly linked to the existence of node-disjoint paths:

**Theorem II.1** ([18, Thm.5.17])**.** *A nontrivial graph $G(V, E)$ is k-connected for some integer $k \geq 2$ if and only if for each pair $s, t \in V$ of distinct nodes, there are at least $k$ node-disjoint s–t-paths in $G$.*

This justifies calling a graph *biconnected* if it is 2-connected, or as equivalently used in [19], $G$ cannot be disconnected by removing a single vertex.

### B. Adversary Model

In many practical environments, flaws in some security system might concern a whole set of devices rather than only a single machine (e.g. exploits found in the firmware of a particular router might apply to a set of routers throughout the infrastructure, or also a buffer-overflow exploit in the operating system (OS) might apply to many machine running on the same OS in the same version). As we are after perfectly private communication, we must not assume any bound on the adversary's computational capabilities. Following the common practice in information-theoretic security, we model computationally unbounded adversaries via *monotonous adversary structures*.

Motivated by the above considerations, we represent an adversary $\mathcal{A}$ by a family of subsets $\mathcal{A} \subset \mathcal{P}(V)$, where $\mathcal{P}(V)$ denotes the power-set of $V$. Such sets within $\mathcal{A}$ may, for example, be characterized by common vulnerabilities. The family $\mathcal{A}$ thus is a collection of potentially compromised sets of devices within the network, each of which represents another possible attack scenario. The set $\mathcal{A}$ is called an *adversary structure*.

We call $\mathcal{A}$ *monotonous* if $Y \in \mathcal{A}$ implies $Z \in \mathcal{A}$ for any $Z \subseteq Y$. This captures the adversary's option to compromise less than the maximal number of nodes, or equivalently, covers situations in which not all of the adversary's servant nodes deliver useful information. A *threshold adversary* is a special case of a monotonous structure, in which all entries have equal cardinality $k$. Taking a *fixed* such threshold $k$, the structure has to no more than $|\mathcal{A}| = \binom{|V|}{k} \in O(|V|^k) = \text{poly}(|V|)$ elements, hence is polynomial. On the contrary, assuming that the adversary can conquer up to, say any fraction of $\lceil p \cdot |V| \rceil$ nodes for $0 < p < 1$ makes $|\mathcal{A}| = \binom{|V|}{\lceil p|V|\rceil} = 2^{O(|V|\log|V|)}$, which is exponential. In the following, we will exclusively deal with *polynomial size monotonous adversary structures*.

It should be noted that a threshold adversary might not always be an appropriate model. As [3] points out, the assumed threshold might yield a gross overestimation of the required graph connectivity, hence working with the more general concept of a monotonous structure adds flexibility. The work of [4] is an explicit account for minimal connectivity models, which partially helps to mitigate this issue. With the aid of game-theory, we can further generalize these previous views on perfectly private communication from a discrete yes/no-classification towards a continuous quantitative risk assessment. Details follow in Section II-E.

The physical adversary is assumed capable of capturing any set $Y \in \mathcal{A}$. Those captured nodes are entirely under the adversary's control, meaning that he is free to block, insert, modify or passively read any message passing through nodes in $Y$. Such an adversary is said to be *k-active*, if he can conquer any union of up to $k$ sets from $\mathcal{A}$. Contrary to this, a *k-passive* adversary is only allowed to extract (read)

information, but otherwise strictly follows the protocol without any active fiddling. Moreover, any adversary (regardless of active or passive) is assumed to know the entire protocol specification, message space, topology of the network, and any inputs except for Alice's secret message $m$ and the coin flips $r$ used for transmission by Alice if the protocol uses randomness (such as most cryptographic protocols do).

### C. Security Model

We will use the security model put forth in [4]: at the beginning, the adversary chooses the plain text distribution $\Pr$ and the nodes to conquer from the adversary structure $\mathcal{A}$. For the actual transmission of a secret message $m$, the sender Alice uses a randomized protocol, taking the random coins $r$ as an input that is *unknown* to the attacker. The adversary's *view* is the information acquired from eavesdropping on the protocol. It is denoted as $\mathcal{A}(m, r)$, whenever he extracts the message $m$ from the information in his possession. For $\varepsilon > 0$, we say that the transmission is $\varepsilon$-*private*, if for every two messages $m_0 \neq m_1$ and every $r$, $\sum_c |\Pr[\mathcal{A}(m_0, r) = c] - \Pr[\mathcal{A}(m_1, r)]| \leq 2\varepsilon$. The probabilities are taken over the coin flips of the honest parties, and the sum is over all possible values of the adversary's view. For $\delta > 0$, we call the protocol $\delta$-*reliable*, if with probability at least $1 - \delta$, Bob terminates the protocol with the correct result $m$. The probability is over the choices of $m$ and the coin flips of all internal transmission nodes in $V$ and the adversary. We call a protocol $(\varepsilon, \delta)$-*secure*, if it is $\varepsilon$-private and $\delta$-reliable. It is said to be *efficient*, when the round complexity and bit complexity are both polynomial in the size of the network, $\log \frac{1}{\varepsilon}$ and $\log \frac{1}{\delta}$ if $\varepsilon > 0, \delta > 0$. Any $(0, 0)$-secure protocol is called *perfectly secure*, and a communication having this performance guarantee is called *perfectly secure message transmission (PSMT)*. In this work, we will consider a slightly weaker notion, which we will call *arbitrarily secure message transmission (ASMT)*.

**Definition II.2** (arbitrarily secure message transmission)**.** *A communication protocol is called* arbitrarily secure*, if for any (small) $\varepsilon > 0, \delta > 0$, we can efficiently run it in a way that achieves efficient $(\varepsilon, \delta)$-security.*

**Remark II.3.** *Note the kind of "duality" between intractability-based and information-theoretic security: for computational (intractability-based) security, we must assume limited computational power of the adversary, while allowing the attacker to listen to all conversation over the channel. Likewise, information-theoretic security imposes no limits on the computational power, yet must assume that not the entirety of the conversation can be eavesdropped. The latter limitation will manifest itself as a polynomial bound on the cardinality of the adversary structure (permitting infinite computational power for the analysis of whatever information the attacker acquires).*

Graph connectivity has been used in [4] with the aim of judging various network types for their suitability for perfectly secure message transmission in the sense of the above security models. An interesting classification that serves as partial motivation here too has been given by [3]. Their characterization relies on a refined graph-connectivity criterion, which explicitly refers to a given adversary structure $\mathcal{A}$. More precisely, the graph $G$ is called $\mathcal{A}^{(k)}(s, t)$-*subconnected*, if for any $k$ sets $Y_1, \ldots, Y_k \in \mathcal{A}$ the deletion of the nodes in $\bigcup_{l=1}^k Y_l$ from $G$ does not disconnect $s$ and $t$ within $G$. A graph $G$ is said to be $\mathcal{A}^{(k)}$-*connected*, if it is $\mathcal{A}^{(k)}(s, t)$-connected for all pairs $s, t \in V$ where $s \neq t$. With this, we have the following security criterion, referring to perfect secure communication in the above sense.

**Theorem II.4** ([3])**.** *Perfectly secure message transmission from the sender $s$ to the receiver $t$ in the network $G$ is possible, if and only if $G$ is $\mathcal{A}^{(2)}(s, t)$-subconnected.*

So, it suffices to consider a 2-active adversary in order to decide whether or not PSMT is possible in the given graph. This approach can indeed be improved to better match a real-life setting, using the concepts of *channel-* and *network-vulnerability* [20], which we briefly recap in section II-E later. The next section is devoted to a closer look at the ideas of how to achieve perfectly secure communication within Theorem II.4 and related results (e.g. [6], [5]).

### D. Transmission Model

The general idea underlying all (secure) multipath transmissions schemes between a sender $s$ and receiver (target) $t$ is the following: the sender $s$ chooses a set $P$ of node-disjoint $s$–$t$-paths, and encodes the message $m$ into $n$ packets. Let the entirety of nodes that are used to convey $m$ be denoted as $V(P) = \bigcup_{\pi \in P} V(\pi)$. The attacker takes over a set $Y \in \mathcal{A}$ of nodes in an attempt to learn everything that flows through the nodes in $V(P) \cap Y$. The sender performs the transmission by encoding $m$ into $|P|$ pieces $c_1, \ldots, c_{|P|}$, and sending those to $t$ over their own individual paths in $P$. In the simplest case, this can be done by conventional XOR-secret-sharing, i.e. $m = c_1 \oplus c_2 \oplus \cdots \oplus c_{|P|}$, where $\oplus$ is the bitwise XOR, and all but one of the $c_i$'s are random strings. The message is protected from discovery unless the attacker intercepts all paths in $P$. Since such encoding is prone to transmission errors and blows up the overall transmission overhead, practical schemes [6], [5] employ more flexible and efficient encodings (e.g., based on polynomial secret sharing to add error correction capabilities and thus gain robustness)[1].

Perfectly secure message transmission demands some encoding and transmission paths $P$ such that *every* attack scenario $Y \in \mathcal{A}$ gives insufficient information to recover $m$. For example, the above XOR-secret-sharing over $n = |P|$ paths displays a one-round PSMT scheme against an attacker with $|Y| < n$ for every $Y \in \mathcal{A}$ (see figure 1; and note that the case $n = 2$ is essentially equivalent to symmetric encryption).

Towards the weaker goal of arbitrarily secure message transmission, we can use randomly chosen (and changing)

---

[1]Feedback rounds can as well be used to gain efficiency and security [7], however, we confine ourselves to one-round protocols here.
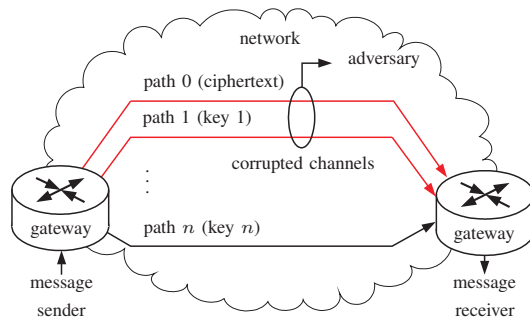
Fig. 1. Basic approach to perfectly secure message transmission

paths to deliver the packets $c_1, c_2, \ldots, c_n$, in an attempt to minimize the attacker's chances to learn enough information to discover $m$. Like for PSMT, we attempt to bypass the attacker, however unlike in PSMT, the randomly chosen paths are not fixed a-priori, thus making ASMT possible even in some cases where the attacker (e.g., thanks to a sufficient threshold) could break the respective PSMT scheme. Moreover, ASMT is doable even using (a sequence of) single-path transmissions, which cannot be used to run PSMT.

*E. Channel- and Network-Vulnerability*

Security of multipath transmission hinges on the existence of at least one path that bypasses all hostile nodes in the network. Consequently, it is the sender's (player 1) intention to optimize his path choices against an attacker (player 2) who seeks optimal spots to sniff the network traffic. This optimization can be done using game-theory.

To this end, take the collection of all existing $s$–$t$-paths, and group them together into a polynomial number of $\mathrm{poly}(|V|)$ different bundles $P_1, P_2, \ldots$ (note that the full enumeration of paths would have exponentially many entries, hence we must work with a feasibly small selection of these). Condense all these bundles in the *strategy set* $PS_1$. With this set, the game is about the sender taking his best randomized choice of a path set for communication. The *opponent strategy set* $PS_2$ is exactly the adversary structure $\mathcal{A}$. The game's payoff matrix $\mathbf{A} = (a_{ij})$ can be defined in binary terms as

$$a_{ij} = \begin{cases} 1, & \text{if the } s\text{–}t\text{-transmission remained secret;} \\ 0, & \text{otherwise,} \end{cases} \quad (1)$$

if $i \in PS_1$ is the chosen pair of paths $\pi_1, \pi_2$, and $j \in \mathcal{A}$ is the compromised set $Y \subset V$ of adversarial nodes within the network $G(V, E)$. We note $a_{ij} = 1$ if the compromised set was insufficient to extract the secret from the adversary's view (transcript). Note that this decision strongly depends on the chosen encoding of $m$, so the evaluation of equation (1) depends on the particular instantiation of the framework protocol (examples are found in [5], [6]).

The game's solution is the *saddle-point value* $v(\mathbf{A}) = \max_{\mathbf{x} \in S(PS_1)} \min_{\mathbf{y} \in S(\mathcal{A})} \mathbf{x}^T \mathbf{A} \mathbf{y}$, where $S(PS_1), S(\mathcal{A})$ denote the set of (discrete) probability distributions over the player's strategy sets. The *equilibrium* is the pair $(\mathbf{x}^*, \mathbf{y}^*) \in$

$S(PS_1) \times S(\mathcal{A})$, at which the saddle-point value $v(\mathbf{A}) = (\mathbf{x}^*)^T \mathbf{A} \mathbf{y}^*$ is attained. The definition of $v(\mathbf{A})$ directly formalizes the aforementioned competition: the sender tries to maximize his chances of keeping the message secure (maximization over all randomized choices $\mathbf{x} \in S(PS_1)$), while the attacker tries his best to discover the message (minimization of the sender's benefit over all randomized choices $\mathbf{y} \in S(\mathcal{A})$ of nodes to conquer from $\mathcal{A}$).

Such modeling might be inaccurate in a real-life scenario because assuming a zero-sum competition can be a misjudgment of the adversary's intentions. However, as eloquently noted in [21], presuming a zero-sum regime is a valid worst-case approach, since with the binary valuation as above and with $v(\mathbf{A})$ denoting the saddle-point value of the zero-sum game induced by the matrix $\mathbf{A}$, it is easy to prove that

$$\Pr[\text{successful attack}] \leq 1 - v(\mathbf{A}),$$

which holds *regardless* of how the adversary actually behaves, provided that the sender and receiver act according to their zero-sum equilibrium strategy. Notice that the matrix $\mathbf{A}$ specifically models the communication between $s$ and $t$. In [20], the upper bound $1 - v(\mathbf{A}) =: \rho(s, t)$ has been assigned the name *vulnerability*, since it measures the degree to which an attack will be successful.

Applications in which the outcome of the transmission cannot be classified in binary terms as in (1) or perhaps is even random, can arise in infrastructures that use security measures like firewalls, intrusion detection systems, etc., all of which have some positive rate of failure. A straightforward way to recover a deterministic valuation from a random outcome in a transmission scenario is taking expectations of the random outcome. This changes the game's payoff structure from a 0-1-matrix to a matrix with real values, but does no inherent change to the model nor its solution procedure. Since random or more general than binary outcomes can be treated with the very same framework, we avoid unnecessary complications here by leaving this direction aside. Respective details and examples can be found in [20], but are not needed for our upcoming considerations.

**Definition II.5.** *Let a graph $G(V, E)$, an integer $k \geq 1$ and a pair of distinct nodes $s, t \in V$ be given. Assume that an $s$–$t$-communication runs over $k$ paths in the presence of an adversary (structure) $\mathcal{A}$. The* vulnerability *of this $s$–$t$-communication is defined as $\rho(s, t) = 1 - \max_{\mathbf{x} \in S(PS_1)} \min_{\mathbf{y} \in S(\mathcal{A})} \mathbf{x}^T \mathbf{A} \mathbf{y}$, where $\mathbf{A} \in \{0, 1\}^{|PS_1| \times |\mathcal{A}|}$ models the zero-sum communication game with the payoffs as defined through* (1).

As not all nodes in a network might be actively communicating, it makes sense to restrict the attention to only a certain set of pairs $U \subseteq V \times V$ that will eventually attempt a private conversation. We call the entirety of these pairs a *communication relation*, whose vulnerability is our measure of overall security in the network $G(V, E)$.

**Definition II.6.** *For a communication relation $U \subseteq V \times V$, the network $G(V, E)$ has the* vulnerability

$$\rho(G, U) := \max_{s,t \in U} \rho(s, t). \qquad (2)$$

Convention (2) is justified by the maximum-principle that is common practice in security audits: the overall security of a system is determined by the vulnerability of its weakest component (similarly to a chain being as strong as its weakest element). In the following, we will use the following characterization of ASMT based on the vulnerability.

**Theorem II.7** ([20])**.** *Let Alice and Bob set up their game matrix with binary entries $a_{ij} \in \{0, 1\}$, where $a_{ij} = 1$ if and only if a message can securely and correctly be delivered by choosing the $i$-th pure strategy, and the adversary uses his $j$-th pure strategy for attacking. Then $\rho(\mathbf{A}) \in [0, 1]$, and*

1) *If $\rho(\mathbf{A}) < 1$, then for any $\varepsilon > 0$ there is a protocol so that Alice and Bob can communicate with an eavesdropping probability of at most $\varepsilon$ and a chance of at least $1 - \varepsilon$ to deliver the message correctly.*

2) *If $\rho(\mathbf{A}) = 1$, then the probability of the message being extracted and possibly modified by the adversary is 1.*

*F. How ASMT Relates to PSMT and Risk Management*

It is worth noting that in case of a pure binary valuation, ASMT becomes PSMT if the vulnerability is either 0 or 1, in which case the incident of zero vulnerability directly implies a certain graph connectivity. We will exploit this fact later.

Moreover, Theorem II.7 remains valid under a modified setting in which the outcome of a transmission is uncertain. More specifically, while PSMT usually presumes all-or-nothing adversarial access to a node, ASMT can be used with probabilistic security models and uncertain behavior of a node's defense (e.g., a firewalls, virus scanners, etc.). The above characterization of (im)possible ASMT still holds. As a further generalization unlike PSMT, ASMT based on games permits using different scales than zero-one, especially nominal or scales used in qualitative risk management. Since the vulnerability is the expected product of likelihood and damage in terms of the given scale, it is nothing else as a *risk*. So, the security guarantees made by ASMT are much better compatible with quantitative (and under a mapping of the vulnerability onto a nominal scale, also qualitative) risk management issues. PSMT is not explicitly designed for integration in such processes. This means that the general problems stated in the next section equivalently refer to the search for a network design that minimizes (general) risk of communication in perhaps even monetary units. Unfortunately, this particular task of risk management will be proven infeasible unless P = NP.

## III. GRAPH AUGMENTATION FOR SECRET COMMUNICATION

Theorems II.4, II.7 as well as the results of [4] and [5] indicate that – on classical grounds, i.e., in the non-quantum setting – multiple paths are inevitable for perfectly and arbitrarily secure communication. This raises the natural question of graph augmentation in order to meet these needs. Using

---

**Problem III.1** MIN-VULNERABILITY-AUGMENTATION

INSTANCE: A graph $G(V, E)$, an adversary structure $\mathcal{A} \subset 2^V$, a set of pairs $U \subseteq V \times V$ that can communicate and a set $\widetilde{E}$ of additional (candidate) edges with costs $c : \widetilde{E} \to \mathbb{Z}^+$, and a budget limit $B \in \mathbb{Z}$.

SOLUTION: An edge augmentation $E^+ \subseteq V \times V \setminus E$ within the budget limit $c(E^+) \leq B$.

MEASURE: The vulnerability $\rho(G(V, E \cup E^+), U) = \max_{(u,v) \in U} \rho(u, v)$, where $\rho(x, y)$ is the vulnerability of an $x$–$y$-communication in $G$

---

the aforementioned game-theoretic framework and Theorem II.7 in particular, the problem boils down to asking for an augmentation that yields a vulnerability $\rho(G, U) \leq \varepsilon < 1$ for a given network $G$, communication relation $U$ and risk threshold $\varepsilon$. Besides the decision-version of the problem, our main interest in the following lies in the respective *search* problem. Suppose that the network is insufficiently connected so that perfectly and arbitrarily secure transmission are both ruled out by any known conventional criterion (e.g. [3], [4], [5]). Then we seek the smallest (cheapest) edge-augmentation to $G$ that would at least give $\rho(G, U) \leq \varepsilon$, so that at least ASMT is possible, even if PSMT might still be out of reach. This is problem III.1.

Towards formulating optimization problems, we let $\widetilde{E} \subset V \times V \setminus E$ be a set of candidate edges not yet existing in the graph $G(V, E)$. Furthermore, let a function $c : \widetilde{E} \to \mathbb{Q}^+$ measure the costs for any of these edges. For reasons of tractability (theoretical as well as computational), we assume that $c(E^+)$ can be computed in poly($|E^+|$) time by a Turing-machine that leaves an encoding of $c(E^+) = \frac{a}{b} \in \mathbb{Q}^+$ on its output tape of the form $\#a\#b\#$, where $a, b$ are natural (radix-based) encodings of the integers $a$ and $b$.

The "reverse" problem III.2, which asks for the cheapest augmentation that undercuts a given vulnerability limit, is treated later.

In the following sections, we will investigate the complexity of both problems, and discover the existence of efficient exact solution algorithms as equivalent to P = NP. Both problems are known to be NP-hard [13], but even despite this fact, there is no point in looking for approximation algorithms.

Before getting to the complexity-theoretic details, let us consider the obvious variants of the above problems; why not consider vertex-augmentations or mixed (vertex- and edge-)

---

**Problem III.2** MIN-COST-SECURITY

INSTANCE: A graph $G(V, E)$, an adversary structure $\mathcal{A} \subset 2^V$, a set of pairs $U \subseteq V \times V$ that can communicate and a set $\widetilde{E}$ of additional (candidate) edges with costs $c : \widetilde{E} \to \mathbb{Z}^+$, and a vulnerability limit $\varepsilon$.

SOLUTION: An edge augmentation $E^+ \subseteq V \times V \setminus E$ that achieves the vulnerability limit $\rho(G(V, E \cup E^+), U)) \leq \varepsilon$.

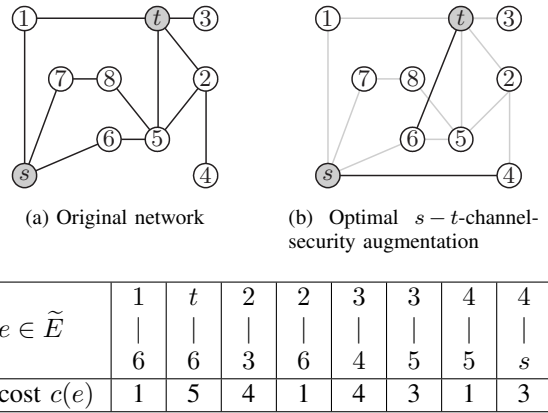MEASURE: The total cost $c(E^+)$ of the augmentation $E^+$.

---

augmentations? It is easy to see that adding only vertices does no change to the vulnerability, since the nodes are all isolated. Adding vertices and edges is equivalent to adding the vertices in first place (leaving the problem's solution unchanged), and afterwards consider a pure edge-augmentation only. So, edge augmentations cover both of these cases.

*Example*

Problem MIN-VULNERABILITY-AUGMENTATION and MIN-COST-SECURITY both admit representations as mixed-integer programming problems [22]. Therefore, solutions for small networks might be feasible in a practical setting. Moreover, the representation of either problem is trivially converted into a representation of the other, so that linear programming software (e.g. Cplex or lp_solve) can be applied to both. For example, consider the network shown in Figure 2a, being the yet unaugmented graph. We solve the respective instance of MIN-VULNERABILITY-AUGMENTATION for an adversary structure $\mathcal{A} = \{U \subset V : |U| = 3\}$ and two-path transmission from $s$ to $t$, where the encoding of the message $m$ is by a $(2, 2)$-XOR-secret sharing of the form $m = r_1 \oplus r_2$, where $r_1$ is random and $r_2 = m \oplus r_1$ (one-time pad symmetric encryption under key $r_1$). Consequently, the transmission is perfectly private unless both, $r_1$ and $r_2$ are intercepted by the attacker. Finally, let the budget limit be $B = 18$ and take the set $\widetilde{E}$ of candidate edges along with edge weights as given by Figure 2c.

Observe that $Y_{\text{cut}} = \{1, 8, 6\} \in \mathcal{A}$ so that no communication from $s$–$t$ is possible without traversing a node in $Y_{\text{cut}}$ in the *unaugmented network* shown in Figure 2a (another cut would be $\{1, 5\}$). Consequently, a fraction of $v = 0$ messages can be delivered secretly and hence the vulnerability is $\rho = 1 - v = 1$ for the unaugmented network. Contrary to this, the fully augmented network including all edges in $\widetilde{E}$ permits 141 different $s$–$t$-paths, from which we can form a set $PS_1$ having 295 pairs of node-disjoint paths. The adversary has – in either case – $|PS_2| = |\mathcal{A}| = \binom{8}{3} = 56$ possible attack strategies (where attacks on $s$ or $t$ are excluded for obvious reasons). Setting up the full game matrix results in a $(295 \times 56)$-tableau, from which we can iteratively and alternatingly delete rows and columns whose payoff is uniformly worse than for another column (in game-theory terminology, we delete the *dominated strategies*). This reduction leaves us with a $6 \times 4$ payoff matrix **A**, shown in Figure 3b, along with the remaining strategies for both players, listed in Figure 3a. All other existing strategies are either redundant (i.e., yield duplicate rows or columns in the matrix) or give less or equal revenue than another strategy (i.e., are dominated). Solving the linear program (in polynomial time [23]) gives $v(\mathbf{A}) = 0.5$ at the full cost of $c(\widetilde{E}) = 22$. Our goal is finding the *minimal* augmentation obeying the cost limit of 18.

Figure 2b displays the solution $E^+ = \{t\text{–}6, 4\text{–}s\}$ for MIN-VULNERABILITY-AUGMENTATION, having $\rho = 0.5$ as the maximal attack probability, as opposed to $\rho = 1$ in the unaugmented graph. Seeking the minimal cost augmentation



(a) Original network  (b) Optimal $s - t$-channel-security augmentation

| $e \in \widetilde{E}$ | 1<br>\|<br>6 | $t$<br>\|<br>6 | 2<br>\|<br>3 | 2<br>\|<br>6 | 3<br>\|<br>4 | 3<br>\|<br>5 | 4<br>\|<br>5 | 4<br>\|<br>$s$ |
|---|---|---|---|---|---|---|---|---|
| cost $c(e)$ | 1 | 5 | 4 | 1 | 4 | 3 | 1 | 3 |

(c) Edge augmentation set $\widetilde{E}$

Fig. 2.  Example graph augmentation

| | | $PS_1$<br>(pairs of paths) | | $\mathcal{A} = PS_2$<br>(compromised) |
|---|---|---|---|---|
| strategy number | 1 | $s$–4–2–$t$ | $s$–1–$t$ | $1, 5, 6$ |
| | 2 | $s$–6–$t$ | $s$–1–$t$ | $1, 4, 6$ |
| | 3 | $s$–6–$t$ | $s$–4–2–$t$ | $1, 4, 5$ |
| | 4 | $s$–7–8–5–$t$ | $s$–1–$t$ | $4, 5, 6$ |
| | 5 | $s$–7–8–5–$t$ | $s$–4–2–$t$ | |
| | 6 | $s$–7–8–5–$t$ | $s$–6–$t$ | |

(a) Strategy sets

attacker

| | | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| sender | 1 | 1 | 0 | 0 | 1 |
| | 2 | 0 | 0 | 1 | 1 |
| | 3 | 1 | 0 | 1 | 0 |
| | 4 | 0 | 1 | 0 | 1 |
| | 5 | 1 | 1 | 0 | 0 |
| | 6 | 0 | 1 | 1 | 0 |

(b) Payoff matrix **A**

Fig. 3.  Game-theoretic model for our example

to achieve at least $\rho = 0.5$, i.e. solving MIN-COST-SECURITY with $\varepsilon = 0.5$ gives the same solution shown in Figure 2b, coming at price $c(E^+) = 8$, and proving that the previous solution $E^+$ is as well the cheapest for this security demand.

Unfortunately, any heuristic approximation to the general problem (i.e. not all equal edge costs) is doomed to unbounded relative errors, unless $\mathrm{P} = \mathrm{NP}$, as we prove in the sequel.

## IV. COMPLEXITY OF GRAPH AUGMENTATION FOR ASMT

To answer the question whether or not it is feasible to create suitable networks for multipath transmission efficiently, we will use some complexity classes for search problems, besides the decision-problem classes P, NP, and the set NPC of problems that are complete for NP. The class FP is the set of all binary relations $P(x, y)$ such that there is an algorithm A that runs in time $\text{poly}(|x|)$ and outputs some $y$ such that $P(x, y)$ holds. The class $\mathrm{FP}^{\mathrm{NP}}$ is defined in exactly the same way, except that A is allowed to make queries to an NP-oracle, where a call to the oracle takes only one step.

An *instance* of an optimization problem is denoted by $I$. By $A(I)$, we denote the result of the algorithm $A$ when applied to the instance $I$ of the (general) optimization problem (e.g., MAX-CLIQUE). For many computationally hard problems efficient approximations are known (one example is MAX-CUT, for which an astonishingly good approximation has been found by [24]). An excellent account is given in [25], from which we will repeatedly draw in the following. Here we give our definitions only for minimization problems.

**Definition IV.1.** *Given an instance $I$ of a minimization (optimization) problem, an algorithm $A$ is called an* approximation *algorithm, if its output $A(I)$ is a feasible (not necessarily optimal) solution of $I$. Given $r \geq 1$, we call $A$ an $r$-*approximation algorithm, if*

$$opt(I) \leq A(I) \leq r \cdot opt(I), \tag{3}$$

*where $opt(I)$ denotes the optimal (minimal) value of the optimization problem $I$.*

The class APX includes all optimization problems for which a polynomial-time $r$-approximation algorithm exists. Strictly speaking, one would need to define APX in terms of the class NPO, which is roughly the set of all "NP-optimization problems". Since we will not need these classes any further, we refer the reader to [25] for details, and refrain from granting APX a full-fledged formal definition (which would unnecessarily complicate things here).

The next section contains a number of technical results needed to establish the main contributions in Section IV. First, we are concerned with the computational feasibility of evaluating the vulnerability of a given network.

*A. Computing Vulnerabilities*

**Lemma IV.2.** *Let $G(V, E)$ be a graph modeling a communication network, and let $\mathcal{A}$ be an adversary structure of size $|\mathcal{A}| = poly(|V|)$. Then it takes only polynomial time to decide whether or not ASMT is possible over $G$ and if so, the respective channel- and network-vulnerabilities can be computed in polynomial time.*

*Proof:* Take any two arbitrary fixed and distinct vertices $s, t \in V$. Observe that, if there is a set $Y$ such that any $s-t$-path $\pi$ intersects $Y$, i.e. $V(\pi) \cap Y \neq \emptyset$, then attacking $Y$ is a classical person-in-the-middle attack, which without pre-shared secrets between $s$ and $t$, trivially rules out any private conversation between $s$ and $t$ (simply because $t$ and the adversary have exactly the same information, so $t$ cannot do anything to decrypt that the adversary could not do equally well). So, fix any ordering of $\mathcal{A} = \{Y_1, \dots, Y_n\}$ and let us iterate over all elements in $\mathcal{A}$ (note that $|\mathcal{A}| = poly(|V|)$ and hence feasibly small to iterate over it). We will construct a game-matrix modeling a single-path transmission from $s$ to $t$ that attempts to circumvent the adversary as good as possible. Moreover, observe that we cannot rely on any encryption between $s$ and $t$, since no (shared) keys are available (public-key cryptography is ruled out by our demand for perfect secrecy).

Each set $Y_j \in \mathcal{A}$ makes yet another attack strategy, so the game-matrix $A$ will have exactly $n = |\mathcal{A}| = poly(|V|)$ columns. We will iterate through $\mathcal{A}$ and look for a path that lets us securely communicate if the nodes in $Y_j$ are compromised. Technically, we will choose a set of $n$ transmission strategies such that the diagonal of the payoff matrix is composed of all 1's, which will ensure a positive saddle-point value and finally enable ASMT by Theorem II.7.

So let $Y_j \in \mathcal{A}$ be given, and look for an $s$-$t$-path that explicitly avoids using any node $v \in Y_j$. This is easily accomplished in polynomial time by running a shortest-path algorithm on a transformed version of $G$. The required transformation is known from the computation of maximal flows with vertex capacities and can identically be re-used to find paths that avoid certain nodes within a graph. We refer the reader to [26] for a concise representation of this trick (where it has been used for a quite different purpose though). Depending on the outcome of the shortest-path algorithm, distinguish two cases:

Case 1: There is no $s$-$t$-path without using nodes in $Y_j$. Then attacking $Y_j$ will intercept any communication from $s$ to $t$, and hence no private channel can be set up. In that case, ASMT is ruled out for obvious reasons. Moreover, the vulnerability of the network and the $s$-$t$-channel are both 1.

Case 2: There is a path $\pi_j$ such that $V(\pi_j) \cap Y_j = \emptyset$. Then, private transmission over $\pi_j$ is possible, and we can assert that $a_{jj} = 1$ in the game-matrix $\mathbf{A}$, since player 1 wins the scenario in which he uses $\pi_j$ for transmission and $Y_j$ is attacked.

In this way, we obtain a path $\pi_j$ that avoids $Y_j$ for all $j = 1, 2, \dots, |\mathcal{A}|$, so that at least on the diagonal of the final game-matrix, we have all 1's. Computing the value of this special matrix game (i.e. a *diagonal game*) is easy, since it is known from game-theory (see [27]) that a diagonal matrix has the saddle-point value $v(\text{diag}(1, \dots, 1)) = \frac{1}{n}$. So, even if player 1 would lose the private transmission game in all other scenarios except for the diagonal of the game-matrix, we get $v(\mathbf{A}) > 0$. Now, regardless of what the off-diagonal entries in the actual game-matrix $\mathbf{A}$ actually do, we surely have $\mathbf{A} \geq \text{diag}(1, \dots, 1)$, where the inequality holds per component. This inequality is preserved if we take averages on either side, giving $\mathbf{x}^T \mathbf{A} \mathbf{y} \geq \mathbf{x}^T \text{diag}(1, \dots, 1) \mathbf{y} > 0$ for all discrete probability distributions $\mathbf{x}, \mathbf{y}$. Hence, ASMT is possible by Theorem II.7.

To compute the exact value of $v(\mathbf{A})$, i.e. the $s$-$t$-channel vulnerability, observe that the matrix $\mathbf{A}$ has exactly $n^2 = |\mathcal{A}|^2$ entries. Computing the off-diagonal elements $a_{ij}$ (with $i \neq j$) is easy because row $i$ corresponds to a path $\pi_i$, column $j$ corresponds to a compromised set $Y_j$, and the entry $a_{ij}$ is found as

$$a_{ij} = \begin{cases} 1, & \text{if } V(\pi_i) \cap Y_j = \emptyset \\ 0, & \text{otherwise.} \end{cases}$$

The saddle-point value of the full game-matrix $\mathbf{A}$ can then be computed in polynomial time by solving a linear optimization program [23]. The overall network vulnerability

can as well be computed in polynomial time, since there are no more than $O(|V|^2)$ $s$–$t$-pairs to look at. ∎

As a simple corollary, the following statement assures that the vulnerability of any augmented graph and given communication relation can be computed in polynomial time.

**Corollary IV.3.** *Let a graph $G(V, E)$ and an adversary structure $\mathcal{A}$ over $V$ be given. Then, for any augmentation $E' \subseteq V \times V$, and any set $U \subseteq V \times V$, the network vulnerability $\rho(G(V, E \cup E'), U)$ of the augmented graph can be calculated in polynomial time.*

The proof is immediate from the proof of Lemma IV.2, when one considers the obvious generalization of the above arguments to transmissions using more than one path and perhaps a different encoding. In any such setup, the same trick as above can be invoked provided that the payoffs can be computed in polynomial time, which is trivially possible in the settings that we consider.

Theorem II.4 classifies perfectly secure transmission in terms of network connectivity. Towards studying the hardness of graph augmentation for security, we relate the problem to graph augmentation for biconnectivity, which is known to be NP-complete in certain variants [19]. If we use two-path transmission and a special adversary structure, we can establish a useful relation between biconnectivity and network vulnerability.

**Lemma IV.4.** *Let a graph $G(V, E)$ be given. Put $n = |V|$ and define an adversary structure as*

$$\mathcal{A} = \{\{1\}, \{2\}, \dots, \{n\}\}. \tag{4}$$

*Then the following two statements hold for the vulnerability of $G$ w.r.t. $\mathcal{A}$ and any sender-receiver pair $s, t \in V$ that performs two-path transmission:*

1)  *$\rho \in \{0, 1\}$, and*
2)  *$G$ is biconnected if and only if $\rho = 0$.*

*Proof:* By theorem II.1, we know that $G$ is biconnected if and only if there are two node-disjoint paths between any two vertices in $G$, i.e. two disjoint channels exist for any pair in $V \times V$. Since the adversary can attack at most one node at a time, $\mathcal{A}$ cannot disconnect any pair that actually has two channels between them. Since the vulnerability is $\rho = \max_{(u,v) \in V \times V} \rho(u, v)$, and the adversary structure is such that $\rho(G, U) \in \{0, 1\}$, we conclude that $\rho = 0$ if and only if the adversary can mount a person-in-the-middle attack between at least one pair in $V \times V$. Otherwise, there is at least one pair such that all paths between them run through a node in $\mathcal{A}$, and the graph has vulnerability $\rho = 1$ and is not biconnected. ∎

### B. On the Existence of Approximations Towards ASMT

Having prepared the groundwork, we are ready to present our main findings. Our first result rules out the existence of efficient approximations for either problem if $P \neq NP$.

**Theorem IV.5.** *Unless $P = NP$, there is no $r$-approximation algorithm for* MIN-VULNERABILITY-AUGMENTATION.

One could equivalently state that MIN-VULNERABILITY-AUGMENTATION $\in$ APX implies $P = NP$. However, as Theorem IV.7 will later show, there is no point in looking for an approximation algorithm at all, since the existence would imply that there is as well a polynomial-time exact solution algorithm for the problem!

*Proof of Theorem IV.5:* Suppose there were an $r$-approximation algorithm $A$ for MIN-VULNERABILITY-AUGMENTATION, and let an instance of the BICONNECTIVITY-AUGMENTATION problem be given, which is known to be NP-complete [19]. This instance is made up by a graph $G(V, E)$, a weight function $w(u, v) \in \mathbb{Z}^+$ for each unordered pair $\{u, v\}$ of nodes in $V$, and a positive integer $B$. The question is to decide whether there is a set $E'$ of unordered pairs of vertices from $V$ such that $\sum_{e \in E'} w(e) \leq B$ such that the graph $G(V, E \cup E')$ is biconnected, i.e. cannot be disconnected by deleting a single vertex [19].

We can easily (almost directly) cast this problem into an instance $I$ of MIN-VULNERABILITY-AUGMENTATION as follows: set the network to be $G$, and use the adversary structure (4). Moreover, define $U := V \times V$, and set the additional edge weights to $w(e)$ as given by the instance of BICONNECTIVITY-AUGMENTATION for all $\tilde{E} := (V \times V) \setminus E$. The budget limit is also taken from the given instance of BICONNECTIVITY-AUGMENTATION. Lemma IV.4 characterizes biconnectivity in terms of the adversary structure $\mathcal{A}$ and its implied vulnerability. So if we solve the MIN-VULNERABILITY-AUGMENTATION problem under the given budget constraints, Lemma IV.4 implies that $G$ can be biconnected within the budget limit if and only if the optimum vulnerability is $\rho^* = 0$. Now, since we have an $r$-approximation algorithm, we conclude that

1)  In case that $A(I) = 0$, (3) implies $\rho^* = 0$ since $0 \leq \rho^* \leq A(G)$, and hence there is a feasible edge-augmentation to biconnect $G$.
2)  Otherwise, if $A(I) > 0$, then again by (3), $0 < A(I) \leq r \cdot \rho^*$, so $\rho^* \neq 0$. Lemma IV.4(1) implies that $\rho^* = 1$, which means that there is at least one pair that can be disconnected by removing a single node, and $G$ cannot be biconnected within the budget limit. ∎

An analogous result holds for MIN-COST-SECURITY too.

**Theorem IV.6.** *Unless $P = NP$, there is no $r$-approximation algorithm for solving* MIN-COST-SECURITY.

As before, one can equivalently state this by saying that MIN-COST-SECURITY $\in$ APX implies $P = NP$. Hence, by the same token as above, looking for approximations to this problem is useless.

*Proof of Theorem IV.6:* Assume an $r$-approximation algorithm $A$ for MIN-COST-SECURITY to be available, and let an instance of a HAMILTONIAN-CIRCUIT problem be given, which is a graph $G(V, E)$ and the question of whether it has a spanning circle. The reduction will be in two steps. We

start by reducing the HAMILTONIAN-CIRCUIT to an instance of the BICONNECTIVITY-AUGMENTATION problem, by modifying the construction of [28]. Consider the biconnectivity augmentation problem on the set $V$, where the edge weights are set to

$$w(u,v) = \begin{cases} 1, & \text{if } (u,v) \in E; \\ 1+rn, & \text{if } (u,v) \notin E, \end{cases}$$

and the budget limit is $n = |V|$. [28, Theorem 4] states that $G$ has a Hamiltonian circuit if and only if there is an edge augmentation of cost less than or equal to $|V|$. Now, suppose that we apply an $r$-approximation algorithm for MINIMUM-COST-SECURITY to exactly this instance, with the adversary structure being (4) again. So the condition $\rho(G,U) \le \frac{1}{2}$ enforces the approximation algorithm to look at only biconnected extensions of the network, by Lemma IV.4.

If $G$ admits a Hamiltonian cycle, then the edge augmentation has cost $\le n$ and our approximation algorithm returns at most $\mathsf{A}(I) \le rn$. On the other hand, if $G$ does not admit a Hamiltonian cycle, then the costs come back $> n$ and at least one edge with cost $1+rn$ must have been used (since $G$ is not Hamiltonian). The minimal costs found by the approximation algorithm for MINIMUM-COST-SECURITY must therefore be at least $\mathsf{A}(I) \ge (n-1) + (1+rn) = (r+1)n > rn$. ∎

Knowing that neither of the problems stated in section III admit a polynomial time $r$-approximation, it is interesting to notice that they indeed admit an exact solution using polynomially many queries to an NP-oracle. The proof is based on a discretization of the optimization measure function, which uses Farey-sequences, and found in [14].

**Theorem IV.7.** MIN-VULNERABILITY-AUGMENTATION $\in$ FP$^{\text{NP}}$

As before, the same result (yet with a different proof) holds for MIN-COST-SECURITY. This as well admits an exact solution in polynomially many steps and calls to an NP-oracle. The proof as well employs Farey-sequences and bisective searching to discretize and narrow down the search space. A different version of this result also appears in [14], however, the proof given here is new and much simpler.

**Theorem IV.8.** MIN-COST-SECURITY $\in$ FP$^{\text{NP}}$

*Proof:* Let $n$ be the size of the given instance of MIN-COST-SECURITY. By definition, the measure function $c : V \times V \to \mathbb{Q}^+$ can be computed in polynomial time, i.e. there is a Turing-machine taking at most $p(n)$ steps to leave an encoding of $c(E) = \frac{a}{b}$ on the tape. This encoding takes the form $\#a\#b\#$, where $a$ and $b$ are nonnegative integers with radix encodings. Since this is printed within $p(n)$ steps, it follows that $a, b \le 2^{q(n)}$, for some polynomial $q$ (in fact, the polynomial $q$ is proportional to the polynomial $p$, with a constant that depends on the radix for the encoding of $a, b$). Consider the normalized costs

$$0 \le \frac{a}{2^{q(n)}b} \le 1. \tag{5}$$

Since $2^{q(n)}b \le 2^{2q(n)}$, we conclude that expression (5),

as having a bounded denominator, must lie within a Farey-sequence of order $2^{2q(n)}$. Using Theorem 28 in [29], we can lower-bound the distance between any two different such fractions as $\left|\frac{a}{2^{q(n)}b} - \frac{a'}{2^{q(n)}b'}\right| \ge \frac{1}{2^{4q(n)}}$. We multiply the last inequality by $2^{q(n)}$ to obtain

$$\left|\frac{a}{b} - \frac{a'}{b'}\right| \ge 2^{-3q(n)} = 2^{-O(p(n))} \tag{6}$$

Since $a, b \le 2^{q(n)}$, we can bound the measure value as $|c(E)| \le 2^{O(p(n))}$. Now, we can continue as in the proof of Theorem IV.8 by running a bisective search over the interval $[0, 2^{O(p(n))}]$, which terminates as soon as the search space has shrunk below the size of $2^{-O(p(n))}$. To this end, we introduce problem IV.1 for the decision version of MIN-COST-SECURITY in the analogous way as before.

---
**Problem IV.1** CHEAP-SECURITY

INSTANCE: the same as for MIN-COST-SECURITY, with an additional cost threshold $C$.

QUESTION: Is there an edge augmentation $E^+$ achieving a desired maximal vulnerability $\rho(G(V, E \cup E^+), U) \le \varepsilon$ such that the cost for $E^+$ are limited as $c(E^+) \le C$?

---

A nondeterministic Turing-machine can easily guess a solution $E^+$ and verify it in polynomial time, since by Lemma IV.2, the vulnerability threshold can be checked efficiently, and by definition of CHEAP-SECURITY, the measure can as well be calculated within $p(n)$ steps. It follows that CHEAP-SECURITY $\in$ NP.

For the bisective search, we make a call to a CHEAP-SECURITY-oracle (i.e. an NP-oracle) in order to decide the direction where to continue our search. The number of steps until we may terminate is, by (6), no more than $O(p(n)^2)$, since by then, the search space has been narrowed down to contain at most one element. This element is obtained by a final (nondeterministic) guess and returned as the result. ∎

Finally, we can state the following relation between our graph augmentation problems towards perfectly private transmissions and the P-vs-NP-question:

**Corollary IV.9.** *The following statements are equivalent:*
1) MIN-VULNERABILITY-AUGMENTATION *can be solved in polynomial time (i.e., the problem is in* FP*)*
2) MIN-COST-SECURITY *can be solved in polynomial time (i.e., the problem is in* FP*)*
3) P = NP.

*Proof:* Observe that FP = FP$^{\text{P}}$ obviously and that FP$^{\text{P}}$ = FP$^{\text{NP}}$ if P = NP. Together with Theorem IV.7, this implies

MIN-VULNERABILITY-AUGMENTATION $\in$ FP.

The claim for MIN-COST-SECURITY follows from Theorem IV.8. On the other hand, if either problem admits a polynomial time solution, then this is trivially an 0-approximation too, so that P = NP by Theorems IV.5 or IV.6. ∎

## V. Discussion and Conclusions

We stress that our treatment is entirely classical, in the sense of leaving aside arbitrarily long distance secure communication via quantum repeaters [10], [30]. Until these techniques have reached a level of maturity to see a wide range roll-out, security is necessarily somewhat tied to computational intractability. However, our treatment may be extended towards further security goals failure resilience (availability) or authenticity. Both are relevant in the quantum setting with and without quantum repeaters. By a trivial change to the modeling, similar equivalences between P = NP and reputation-based authentication [31] or network path redundancy may be derived. One aspect of future considerations will thus be looking for siblings of corollary IV.9 and its related approximation problems for reliable and authentic communication. Alas, the infeasibility of graph augmentation for perfectly private transmissions is strong, since it implies that every heuristic approach to the graph augmentation problem will inevitably perform arbitrarily bad in infinitely many cases. Hence, looking for good approximations for perfect security graph augmentations is (unconditionally) pointless.

As prefigured in remark II.3, we have demonstrated that information-theoretic security and computational security both strongly relate to computational infeasibility, only in quite different ways. The situation in which we would – in the perfect security paradigm – permit the adversary an unlimited number of compromised nodes is trivial, as there is no way of perfectly secure communication without pre-shared secrets, assuming the adversary to keep the transmission network fully under his control.

The final conclusion is nevertheless a positive one: either $P \neq NP$, then strong encryptions like McElice encryption [32] or related will continue to provide a good protection against eavesdropping. Otherwise, if P = NP, then we can feasibly construct networks that permit communication in arbitrarily strong privacy. So, no matter how $P \overset{?}{=} NP$ is ultimately settled, confidentiality remains an achievable goal.

## References

[1] C. Elliott, "The DARPA quantum network," 2007, arXiv:quant-ph/0412029v1.

[2] A. Poppe, M. Peev, and O. Maurhart, "Outline of the SECOQC Quantum-Key-Distribution network in Vienna," *Int. J. of Quantum Information*, vol. 6, no. 2, pp. 209–218, 2008.

[3] M. Ashwin Kumar, P. R. Goundan, K. Srinathan, and C. Pandu Rangan, "On perfectly secure communication over arbitrary networks," in *PODC '02: Proc. of the twenty-first annual symposium on Principles of distributed computing*. New York, NY, USA: ACM, 2002, pp. 193–202.

[4] M. Franklin and R. Wright, "Secure communication in minimal connectivity models," *J. of Cryptology*, vol. 13, no. 1, pp. 9–30, 2000.

[5] Y. Wang and Y. Desmedt, "Perfectly secure message transmission revisited," *IEEE Transactions on Information Theory*, vol. 54, no. 6, pp. 2582–2595, 2008.

[6] M. Fitzi, M. K. Franklin, J. Garay, and S. H. Vardhan, "Towards optimal and efficient perfectly secure message transmission," in *4th Theory of Cryptography Conf. (TCC)*, ser. LNCS LNCS 4392, S. Vadhan, Ed. Springer, 2007, pp. 311–322.

[7] S. Agarwal, R. Cramer, and R. de Haan, "Asymptotically optimal two-round perfectly secure message transmission." in *CRYPTO*, 2006, pp. 394–408. [Online]. Available: http://www.iacr.org/cryptodb/archive/2006/CRYPTO/1885/1885.pdf

[8] P. Shor, "Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer," *SIAM J. on Computing*, vol. 26, pp. 1484–1509, 1997.

[9] L. Adleman, "Molecular computation of solutions to combinatorial problem," *Science*, vol. 266, pp. 1021–1024, Nov 1994.

[10] W. Dür, H.-J. Briegel, J. I. Cirac, and P. Zoller, "Quantum repeaters based on entanglement purification," *Phys. Rev. A*, vol. 59, no. 1, pp. 169–181, Jan 1999.

[11] R. Alleaume, F. Roueff, E. Diamanti, and N. Lütkenhaus, "Topological optimization of quantum key distribution networks," *New J. of Physics*, vol. 11, p. 075002, 2009.

[12] S. Rass and P. Schartner, "Game-theoretic security analysis of quantum networks," in *Proc. of the Third Int. Conf. on Quantum, Nano and Micro Technologies*. IEEE Computer Society, February 2009, pp. 20–25.

[13] S. Rass, A. Wiegele, and P. Schartner, "Building a quantum network: How to optimize security and expenses," *Springer J. of Network and Systems Management*, vol. 18, no. 3, pp. 283–299, 2010, (published online: 23 March 2010).

[14] S. Rass and P. Schartner, "The NP-complete face of information-theoretic security," *Computer Technology and Application*, vol. 2, no. 11, pp. 893–905, 2011, david Publishing Company, ISSN 1934-7332.

[15] Y. Liang, H. V. Poor, and S. Shamai, *Information-Theoretic Security*. now Publishers Inc., 2010.

[16] O. Goldreich and S. Goldwasser, "On the possibility of basing cryptography on the assumption that P = NP," Cryptology ePrint Archive, Tech. Rep. 005, 1998/005.

[17] R. Impagliazzo and S. Rudich, "Limits on the provable consequences of one-way permutations," in *Proc. of the twenty-first annual ACM symposium on Theory of computing*, ser. STOC '89. New York, NY, USA: ACM, 1989, pp. 44–61. [Online]. Available: http://doi.acm.org/10.1145/73007.73012

[18] G. Chartrand and P. Zhang, *Introduction to Graph Theory*, ser. Higher education. Boston: McGraw-Hill, 2005.

[19] M. R. Garey and D. S. Johnson, *Computers and intractability*. New York: Freeman, 1979.

[20] S. Rass and P. Schartner, "A unified framework for the analysis of availability, reliability and security, with applications to quantum networks," *IEEE Transactions on Systems, Man, and Cybernetics – Part C: Applications and Reviews*, vol. 40, no. 5, pp. 107–119, 2010.

[21] T. Alpcan and T. Başar, *Network Security: A Decision and Game Theoretic Approach*. Cambridge University Press, 2010.

[22] S. Rass, "Information-theoretic security as an optimization problem," *J. of Next Generation Information Technology*, vol. 2, no. 3, pp. 72–83, August, 31st 2011.

[23] L. Khachian, "A polynomial algorithm in linear programming," *Soviet Math. Dokl.*, vol. 20, 1979.

[24] M. X. Goemans and D. P. Williamson, "Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming," *J. of the ACM*, vol. 42, no. 6, pp. 1115–1145, Nov. 1995.

[25] G. Ausiello, P. Crescenzi, G. Gambosi, V. Kann, A. Marchetti-Spaccamela, and M. Protasi, *Complexity and Approximation – Combinatorial Optimization Problems and Their Approximability Properties*. Springer, 1999.

[26] A. Abbas, "A hybrid protocol for identification of a maximal set of node disjoint paths," *Int. Arab J. Of Information Technology (IAJIT)*, vol. 6, no. 4, pp. 344–358, 2009.

[27] R. Gibbons, *A Primer in Game Theory*. Pearson Education Ltd., 1992.

[28] K. P. Eswaran and R. E. Tarjan, "Augmentation problems," *SIAM J. on Computing*, vol. 5, no. 4, pp. 653–665, 1976.

[29] G. Hardy and E. Wright, *An introduction to the theory of numbers*, 5th ed. Oxford Science Publications, 1984.

[30] H.-K. Lo and H. F. Chau, "Unconditional security of quantum key distribution over arbitrarily long distances," *Science*, vol. 283, pp. 2050–2056, 1999, arXiv:quant-ph/9803006.

[31] S. Rass and P. Schartner, "Multipath authentication without shared secrets and with applications in quantum networks," in *Proc. of the Int. Conf. on Security and Management (SAM)*, vol. 1. CSREA Press, July 12–15 2010, pp. 111–115.

[32] J. Buchmann and J. Ding, Eds., *Post-Quantum Cryptography*, ser. Lecture Notes in Computer Science 5299. Springer, 2008, Proceedings of the Second International Workshop, PQCrypto 2008 Cincinnati, OH, USA.

# Energy Saving EDF Scheduling for Wireless Sensors on Variable Voltage Processors

Hussein EL Ghor

Lebanese University - IUT Saida,
SNCS Research Center, UT, Saudi Arabia
B.P. 813 Saida, Lebanon
Email: hussein.ghor@ul.edu.lb

El-Hadi M Aggoune

Electrical Engineering Department,
Sensor Networks and Cellular Systems (SNCS) Research Center,
University of Tabuk, 71491 Tabuk, Saudi Arabia
Email: haggoune.sncs@ut.edu.sa

*Abstract*—**Advances in micro technology has led to the development of miniaturized sensor nodes with wireless communication to perform several real-time computations. These systems are deployed wherever it is not possible to maintain a wired network infrastructure and to recharge/replace batteries and the goal is then to prolong as much as possible the lifetime of the system. In our work, we aim to modify the Earliest Deadline First (EDF) scheduling algorithm to minimize the energy consumption using the Dynamic Voltage and Frequency Selection. To this end, we propose an Energy Saving EDF (ES-EDF) algorithm that is capable of stretching the worst case execution time of tasks as much as possible without violating deadlines. We prove that ES-EDF is optimal in minimizing processor energy consumption and maximum lateness for which an upper bound on the processor energy saving is derived. In order to demonstrate the benefits of our algorithm, we evaluate it by means of simulation. Experimental results show that ES-EDF outperforms EDF and Enhanced EDF (E-EDF) algorithms in terms of both percentage of feasible task sets and energy savings.**

## I. INTRODUCTION

Wireless sensors have gained wide interest as a new generation of embedded systems with a broad range of real-time applications. Examples include agriculture, environment, health care, urban development, habitat monitoring, medical care, military applications [2], fire monitoring [3], volcano monitoring [4] and highway traffic coordination. Many wireless sensors are powered by batteries with limited capacity and in many scenarios it is impossible to replace them after deployment, therefore a fundamental objective is to optimize the sensor life time.

The problem of reducing energy consumption imposes additional challenges on the design of many real-time embedded systems. Such systems are characterized by a time varying processor utilization. Simply adapting the operating voltage and frequency of the processor results in improving energy efficiency and therefore battery life of wireless sensors. Dynamic voltage and Frequency scaling (DVFS) is the most well-known technique that trades off the performance for energy consumption by lowering the operating voltage/frequency [5].

In our work, we deal with dynamic scheduling for uniprocessor systems that support periodic tasks. EDF has been shown to be an optimal dynamic scheduling algorithm in the sense that if a set of tasks can be scheduled by any algorithm, then it can be scheduled by EDF [6]. EDF algorithm is typically preemptive, in the sense that, a newly arrived task

may preempt the running task if its absolute deadline is shorter. This dynamic priority assignment allows EDF to exploit the full processor, reaching up to 100% of the available processing time [7].

Already having the EDF scheduler, it was only necessary to find a way to reduce energy consumption of tasks so as to prolong as much as possible the lifetime of the system. Dynamic voltage and frequency scaling (DVFS) is the most efficient technique for reducing CPU energy. It is feasible to run the processor at the weakest frequency while still admiring the deadlines of tasks. In other words, when the frequency is reduced, the processor can operate at a lower supply voltage and so reducing the energy consumption. However, when reducing the processor speed, tasks must take more time to complete their execution. Therefore, it is important to identify the slack time under which we can safely slow down the processor without missing any deadline.

In this paper, we present an approach to find the least-energy voltage schedule for executing real-time tasks on a DVFS processor according to a dynamic priority, preemptive policy, denoted by Energy Saving EDF (ES-EDF). For the minimization of energy consumption, we use DVFS technique that reduces the processor energy by slowing down the DVS processor and stretching the task execution time. We propose a slack-based method for stretching tasks as much as possible while still guaranteeing deadlines. Off-line computing by how long the tasks should be stretched is possible thanks to EDL properties [8].

The rest of the paper is organized as follows. The paper begins with a summary of the related work. Section III defines the system model and terminology used throughout this paper. The necessary background is presented in section IV. In section V, we propose the Energy Saving EDF (ES-EDF) algorithm. In section VI, we present the feasibility analysis for our proposed algorithm. An upper bound on energy savings is derived in section VII. The simulation results for performance evaluation are presented in section VIII. The paper is concluded in section IX.

## II. RELATED WORK

The majority of real-time schedulers are on-line and based on the concept of priority. If the priority is fixed at the initialization for all tasks, the algorithm is called fixed priority algorithm. Rate monotonic scheduling (RM) [9] and deadline

monotonic scheduling (DM) [11] are examples of such algorithms. If it evolves over time, the algorithm is said to be driven by a dynamic priority. The most known algorithm among such scheduling approaches is the Earliest Deadline First (EDF) algorithm [7]. The study reported in this work deals with dynamic priority scheduling, preemptive and without resource and precedence constraints.

Recently, researchers have started exploring energy-efficient scheduling for real-time embedded systems such as wireless sensors. Algorithms proposed in literature have either dynamic or fixed priority, also they can be preemptive or non-preemptive. Although DVS is one of the most important techniques, still some of them consider non-DVS techniques especially when the study of the energy consumption is for processor and devices.

Among the earliest works, Yao et al. [13] proposed an optimal offline-scheduling algorithm for independent set of jobs to minimize energy consumption. The same problem was targeted in [14], but for dependent tasks. Authors have proposed a scheduling algorithm using a variable voltage processor core.

Shin and Choi [22] proposed a power reduction technique for a processor by exploiting the slack times inherent in the system and those arising from variations of execution times of task instances. In this technique, the processor can either be shut down if there is no current active job or adopt the speed such that the current active job finishes at its deadline or the release time of the next job. Later in [24], the same authors have proposed a method that combines an off-line component with an on-line one. By applying this method, they first determine the lowest possible maximum processor speed where the task set is feasible while all deadlines are met. These tasks apply the WCET at all times and consequently some idle time will be obtained. An on-line component is then introduced that can dynamically reduce the processor speed according to the status of task set in order to exploit execution time variations and idle intervals, the only situation a task is stretched is when it is the only one running and has enough time until the next task arrives.

With the aim to find the least energy schedule for executing real-time tasks, authors in [15] proposed an optimal fixed priority policy. Later in [16], Kwon et al. presented important results for task scheduling over a fixed number of voltage levels.

Based on exploiting slack times, Aydin et al. [17] have addressed the problem of minimizing energy by proposing a dynamic speculative scheduling algorithm. Later in [18], authors have considered that task deadlines are different from the task period. Under this assumption, authors have addressed the problem of computing task slowdown factors. Recent works have proposed further dynamic voltage scaling techniques to enhance the energy gains at run-time [19] [20] [21].

For scheduling periodic real-time tasks on a variable speed processor with realistic discrete speeds, Mejia et al. [25] have proposed a heuristic algorithm that finds near-optimal solutions at low cost. This method produces a 2-approximate solution to the optimization problem. Later in [26], authors proposed a polynomial time $(1 + \epsilon)$-approximation algorithm for the scheduling of periodic real-time tasks, where $\epsilon$ is the tolerable error margin given by users ($0 < \epsilon < 1$).

Recently, reliability became as important as energy efficiency especially in real-time embedded systems like satellite and monitoring systems. Following this idea, several scheduling policies have been proposed for various task models. Zhu et al. [28], [27] proposed a reliability-aware power management (RAPM) algorithm for periodic real-time tasks that can study the negative effects of voltage scaling on system reliability. This work was later extended in [29], authors improved the quality of assurance for all tasks by managing only a subset of jobs from each task.

In a preemptive scheme certain low priority tasks may be suspended if higher priority tasks need to be executed. This will lead to a more flexible scheme but with a certain time overhead. Jejurikar et al. [30] focused on the system level power management via the computation of static slow-down factors under synchronization constraints where tasks are scheduled based on a preemptive scheduling policy. For a similar task model, authors in [31] proposed the concept of frequency locking and extended the Priority Ceiling Protocol (PCP) by locking the processor frequency in a restricted way, so that the cost in frequency switching is better managed. The major inconvenient is that frequency switching is shown to be not found. This work was later extended in [32] by avoiding voltage emergency. Authors explored one of the pioneering real-time task synchronization with the minimization of energy consumption and voltage emergency prevention.

## III. System Model and Terminology

### A. Task Set

The real-time system considered in this work consists of two major units: Real-time Operating System and the Storage unit. The considered RTOS is equipped with a DVFS-enabled processor. The variable speed processor is assumed to be working with $N$ discrete frequencies ranging from $f_{min} = f_1 \leq f_2 \leq \cdots \leq f_n = f_{max}$. The power consumption of the tasks running in the processor and frequency levels are in a way coupled together. When we change the speed of a processor, its operating frequency is changed and hence the power consumption of tasks is proportionately changed the voltage to a value which is supported at that operating frequency. We denote by $P_n$ and $V_n$ respectively the power consumption and voltage level correspondent to clock frequency $f_n$. We consider that $P_n$ is the overall power consumption of the RTOS. This means that $P_n$ is a combination of both dynamic power consumption and leakage power consumption. We also ignore the time and energy overhead incurred in changing the frequency and voltage of the processor.

We use the term slowdown factor $S_n$ as the ratio of the scheduled speed to the maximum processor speed. $S_n$ ranges from $S_{min}$ to 1:

$$S_n = \frac{f_n}{f_{max}} \qquad (1)$$

We consider in our work that each job has different power dissipation that varies according to its frequencies. Consequently, a task will have maximum power dissipation at its maximum frequency and this power consumption decreases as

the frequency decreases. Consequently, the power dissipation of a task must be defined as function of the task index and its corresponding slowdown factor $P_i(\tau_i, S_i)$.

Any application executed on this RTOS is normally composed of multiple tasks with different levels of priority. We consider here independent and preemptive periodic tasks. A task set $\Gamma$ of $n$ tasks is denoted as follows: $\Gamma = \{\tau_i \mid 1 \leq i \leq n\}$. $\tau_i$ is characterized by four-tuple $(r_i, C_i, D_i, T_i)$ where $r_i$, $C_i$, $D_i$ and $T_i$ indicate the release time, the worst case execution time (WCET), the relative deadline and the period respectively. Release time $r_i$ of task $\tau_i$ is equal to $kT_i$, $k = 0, 1, 2, \cdots$. We assume that $0 \leq C_i \leq D_i \leq T_i$ for each $1 \leq i \leq n$.

When a task $\tau_i$ is stretched by a slowdown factor $S_i$, then its actual execution time $(C_i(a))$ at frequency $f_i$ will be $C_i/S_i$. When the processor is running at it maximum frequency, then $C_i(a) = C_i$. The energy dissipation $(E_i)$ of a task $\tau_i$ is computed as:

$$E_i = P_i(\tau_i, S_i) \times (C_i/S_i) \qquad (2)$$

### B. Energy Storage

Our RTOS relies on an ideal energy storage unit, battery for example, that has a nominal capacity, namely E, corresponding to a maximum energy (expressed in Joules or Watts-hour). The energy level has to remain between two boundaries $E_{min}$ and $E_{max}$ with $E = E_{max} - E_{min}$.

We denote by $E(t)$, the energy stored in the battery at time $t$. At any time, the stored energy is no more than the storage capacity, that is

$$E(t) \leq E \quad \forall t \qquad (3)$$

## IV. Necessary Background

The main objective behind a scheduling algorithm is to determine, for a given set of jobs, the order in which tasks are to be executed [23]. In real-time systems, the main goal of the scheduling algorithm is to complete the execution of all jobs while guaranteeing their deadlines. Before we get into the details of the scheduler implementation it is important to understand some of the more important real-time scheduling approaches, namely fixed-priority algorithms, including rate monotonic [9] and deadline monotonic [11], and dynamic-priority algorithms, including the earliest deadline first (EDF) algorithm [10]. EDF schedules at each instant of time $t$, the ready task (i.e. the task that may be processed and is not yet completed) whose deadline is closest to $t$. EDF is an optimal scheduling algorithm on preemptive uniprocessors. The EDF algorithm can achieve an utilization of 100% of the available processing time. The processor utilization of a system is computed as follows:

$$U_p = \sum_{i=1}^{n} \frac{C_i}{T_i} \qquad (4)$$

A periodic task set with deadlines equal to periods is schedulable by EDF if and only if the total processor utilization $U_p$ is less than or equal to one [9].

### A. Static EDS Scheduling

The implementation of EDF consists in ordering tasks according to their priority and executing them as soon as possible with no inserted idle time. Such implementation is known as Earliest Deadline as Soon as possible (EDS) [12]. For a given periodic task set, the EDS schedule can be pre-computed and memorized in order to reduce scheduling overheads at run time.

### B. Static EDL Scheduling

EDL algorithm is based on the notion of delaying the execution of jobs by as much as possible without causing their deadlines to be missed. Although the usual scheduling scheme is EDS, EDL is very often considered for processor idle time analysis. In [12], Chetto and Chetto presented a simple method to determine the location and length of idle time in any window of a sequence generated by the two different implementations of EDF and EDL.

Before the system begins to operate, static EDL schedule is computed for a given task set. More precisely, the duration and position of the idle times is determined by mapping out the EDL schedule produced from time zero up to the end of the first hyperperiod. Let $T_{LCM}$ the hyperperiod be equal to the least common multiple of the task periods where $T_{LCM} = lcm(T_1, T_2, \cdots, T_n)$. Hence, determining the EDL schedule for the interval $[0, T_{LCM}]$ is realized by means of the two following vectors [8]:

*Static deadline vector* $\mathcal{K}$: it represents the times at which idle times occur within the first hyperperiod. $\mathcal{K} = \{k_0, k_1, \cdots, k_i, k_{i+1}, \cdots, k_q\}$ where $k_0 = 0$, $x_i = T_i - D_i$ and $k_i < k_{i+1}$ for all $1 \leq i \leq n$.

Let us consider $q \leq N + 1$ where $N$ denotes the number of jobs within the first hyperperiod. Then $k_q$ is equal to

$$k_q = T_{LCM} - min\{x_i \mid 1 \leq i \leq n\} \qquad (5)$$

*Static idle time vector* $\mathcal{D}$: it represents the lengths of the idle times which start at time instants given by $\mathcal{K}$.

$\mathcal{D} = (\Delta_0, \Delta_1, \cdots, \Delta_i, \Delta_{i+1}, \cdots, \Delta_q)$. $\Delta_i$ corresponds to the length of the idle time that starts at time $k_i$.

Vector $\mathcal{D}$ is defined by a recurrent formula as follows:

$$\Delta_q = min\{x_i \mid 1 \leq i \leq n\} \qquad (6)$$

$$\Delta_i = max(0, F_i), \ where \ i = q - 1 \ down \ to \ 0 \qquad (7)$$

$$F_i = (T_{LCM} - k_i) - \sum_{j=1}^{n} \lceil \frac{T_{LCM} - x_j - k_i}{T_j} \rceil C_j - \sum_{k=i+1}^{q} \Delta_k \qquad (8)$$

Where $\lceil y \rceil$ is the least integer greater than or equal to $y$.

Under energy constraints, executing the jobs as late as possible within the time interval $[0, T_{LCM}[$ consists in first ordering the jobs according to the EDF rule and second stretching the execution time $C_i$ by its corresponding $\Delta_i$.

*Illustrative Example:* Consider a periodic task set $\Gamma$ that is composed of three tasks, $\Gamma = \{\tau_i \mid 1 \leq i \leq n\}$ and $\tau_i = (C_i, D_i, T_i)$. Let $\tau_1 = (1, 3, 5)$, $\tau_2 = (2, 7, 10)$ and $\tau_3 = (3, 12, 20)$.

From formulae 6, 7 and 8, we have $\mathcal{K} = (0, 3, 7, 8, 12, 13, 17, 18)$ and $\mathcal{D} = (2, 2, 0, 1, 0, 2, 0, 2)$. The EDL schedule for $\Gamma$ produced at the first hyperperiod is described in figure 1.
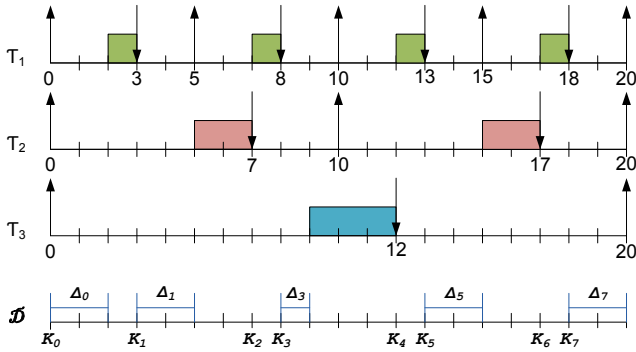


Fig. 1. Static EDL scheduling of task set $\Gamma$

## V. ENERGY SAVING - EDF (ES-EDF) SCHEDULING ALGORITHM

In this section, we describe the ES-EDF task scheduling algorithm that minimizes the CPU energy. In ES-EDF, an optimal voltage and frequency for a given task instance is computed that guarantees schedulability and minimizes energy consumption.

### A. Presentation of the Algorithm

The aim from ES-EDF is to develop dynamic task scheduling algorithm based on EDF that schedule any periodic task set feasibly, and minimize the CPU energy consumption. We use a modified Earliest Deadline First (EDF) strategy to reduce the CPU energy consumption by using the Dynamic Voltage and Frequency Selection. For this sake, We propose a slack-based method for stretching tasks as much as possible while still guaranteeing deadlines. Off-line computing by how long the tasks should be stretched is possible thanks to EDL properties.

### B. Slack Time Method

In the following analytical analysis, we determine the maximum scaling time of a DVFS system that results in minimum energy consumption for the processor and without any deadline violation. That is, we have to determine the lowest maximum processor speed to execute a real-time task set on a variable speed processor while guaranteeing the deadlines of tasks. This scaling time for each task instance is called the *optimal scaling time*.

Our approach is based on the assumption that the parameters of each task is well known off-line and that all tasks are released at time $t = 0$. This assumption is very important since otherwise we may not be able to fully utilize the benefits provided by the used variable speed processor.

As an initial schedule, we try to execute all task instances according to the earliest deadline first strategy. Let us consider that there are $M$ task instances in the ready queue. The start time and finish time of task $\tau_i$ are represented by $St_i$ and $Ft_i$ respectively.

Assume that the start time of the first task instance $\tau_1$ in the ready queue is equal to its release time.

$$St_1 = k_0 = 0 \tag{9}$$

In order to stretch the execution time $C_i$ of task $\tau_i$ as much as possible without violating the deadline $D_i$, the determination of the latest start time for every task instance requires preliminary construction of the schedule produced by the so-called Earliest Deadline as Late as possible (EDL) algorithm.

To involve an acceptable overhead at run-time, off-line computations are done by ES-EDF in order to compute efficiently the static EDL schedule without losing any time. Before the system begins to operate, we estimate the localization and the duration of the idle times produced at time $t = 0$ till the end of the hyperperiod. This means that ES-EDF computes the static deadline vector $\mathcal{K}$ and static idle time vector $\mathcal{D}$.

Thus, the start time of the remaining task instances is:

$$St_i = k_{i-1} \tag{10}$$

The total time executed by the task set within $[0, St_i]$ is denoted by $A_k$ where

$$A_k = \sum_{D_k \leq St_i} C_k \tag{11}$$

Consequently, the finish time of the remaining task instances is:

$$Ft_i = C_i + \sum_{k_j \leq St_i} \Delta_j + A_k \tag{12}$$

where $1 \leq i \leq M - 1$

To decrease the processor speed as much as possible and as long as the system will be able to meet all the deadlines, we have to compute the static idle time vector $\Delta_i$ at each task start time $St_i$. Consequently, the task's execution time will be stretched to the actual execution time $C_i(a)$ where:

$$C_i(a) = C_i + \Delta(St_i) \tag{13}$$

The slowdown factor is thus calculated thanks to the following equation:

$$S_i = \frac{C_i}{C_i(a)} \tag{14}$$

### C. ES-EDF Algorithm

The ES-EDF works as follows: First ES-EDF computes the static EDL schedule including static deadline vector $\mathcal{K}$ and static idle time vector $\mathcal{D}$. Before authorizing the execution of the task instance with highest priority, ES-EDF adds the static idle time at $St_i$ to the execution time $C_i$. Now, the execution time of the task instance will be stretched to its actual execution time $(C_i(a))$ without violating deadlines. Upon stretching

the execution time of a task instance, the energy dissipation decreases. ES-EDF can now compute the slowdown factor of the corresponding task instance and consequently the energy dissipation can be chosen.

The major components of ES-EDF are: $E(t)$, $\mathcal{K}$ and $\mathcal{D}$ where $t$ is the current time, $E(t)$ the amount of energy that is currently stored at time $t$ that means the remaining amount of energy in the energy storage at time $t$. $\mathcal{K}$ and $\mathcal{D}$ are respectively the static deadline vector and the static idle time vector. Moreover, we use the function $execute()$ to put the processor to run the ready job with the earliest deadline.

We describe in algorithm 1 the pseudo code of the ES-EDF scheduler:

---

**Algorithm 1** Energy Saving - Earliest Deadline First (ES-EDF) Algorithm

---

**Require:** A Set of $M$ periodic Tasks $\Gamma = \{\tau_i | \tau_i = (r_i, C_i, D_i, T_i, E_i)\ \ i = 1, \cdots, M\}$ According to $EDF$, current time $t$, battery with capacity ranging from $E_{max}$ to $E_{min}$, energy level of the battery $E(t)$.
**Require:** A processor working with $N$ discrete frequencies ranging from $f_1$ ($f_{min}$) to $f_N$ ($f_{max}$).
**Ensure:** $ES - EDF$ Schedule.
1: Sort task instances according to the EDF rule
2: Determine the start time of task instances
3: Compute the static EDL vectors $\mathcal{K}$ and $\mathcal{D}$
4: **for** i=1:M **do**
5:   **if** i==1 **then**
6:     $St_1 = k_0 = 0$
7:   **else**
8:     $St_i = k_{i-1}$
9:   **end if**
10: **end for**
11: **while** $E(t) > 0$ **do**
12:   Actual Execution Time $C_i(a) = C_i + \Delta(St_i)$
13:   $Ft_i = St_i + C_i + C_i(a)$
14:   Slowdown factor $S_i = C_i/C_i(a)$
15:   Update execution time
16:   Select the relative Energy Consumption ($E_i$)
17:   Calculate the remaining energy in the battery at the end of the execution.
18:   $E(Ft_i) = E(St_i) - E_i(St_i, Ft_i)$
19:   execute()
20:   Remove task $\tau_i$ from ready task list
21: **end while**

---

### D. Illustrative Example

Consider a periodic task set $\Gamma$ that is composed of three tasks, $\Gamma = \{\tau_i \mid 1 \leq i \leq n\}$ and $\tau_i = (C_i, D_i, T_i)$. Let $\tau_1 = (1, 3, 5)$, $\tau_2 = (2, 7, 10)$ and $\tau_3 = (3, 12, 20)$. We assume that the energy storage capacity is $E = 350$ energy units at $t = 0$. The processor is assumed to be working with ten discrete slowdown factors $S_i = \{1, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1\}$. The power dissipation of tasks $\tau_i$ is shown in table I:

First of all, we have to schedule $\Gamma$ according to EDF within the first hyperperiod, from 0 to 20. We verify that $\Gamma$ is not schedulable since the battery capacity is equal to zero at $t = 11$ and consequently the deadline miss rate is about 30%. In details:

At time $t = 0$, all tasks are ready. $\tau_1$ is the highest priority task and is executed until $t = 1$ where $E(1) = 320$ energy units. At time $t = 1$, $\tau_2$ is the highest priority task and is executed until $t = 3$ where $E(3) = 240$ energy units. $\tau_3$ is now the highest priority task and is executed until $t = 5$ where it is preempted by $\tau_1$. $\tau_3$ resumes its execution at time $t = 6$ and is completed at $t = 7$ where $E(7) = 30$ energy units. The remaining energy in the battery unit is sufficient only to execute $\tau_1$ until $t = 11$. Now, the battery is empty and consequently the scheduling is terminated where the deadline miss rate is about 30% (3(a)).
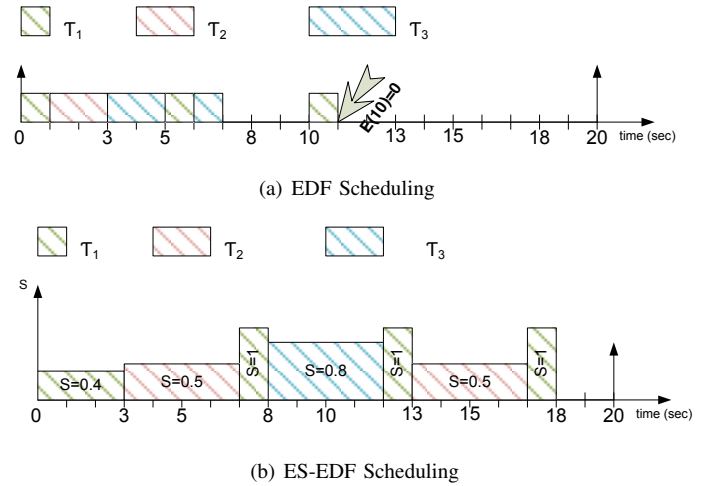


(a) EDF Scheduling

(b) ES-EDF Scheduling

Fig. 2.   Scheduling of task set $\Gamma$

To increase the efficiency of the processor and decrease the energy consumption, let us try scheduling the same task set $\Gamma$ but with ES-EDF. We find that $\Gamma$ is schedulable since all tasks are executed without violating deadlines and without getting out of energy. Let us explain how ES-EDF constructs the schedule (3(a)). Before beginning the execution of tasks, ES-EDF computes $\mathcal{K}$ and $\mathcal{D}$. $\mathcal{K} = (0, 3, 7, 8, 12, 13, 17, 18)$ and $\mathcal{D} = (2, 2, 0, 1, 0, 2, 0, 2)$.

At time $0$, the residual capacity i.e. remaining energy is maximum since the storage is full. $\tau_1$ is the highest priority task where $\Delta_0 = 2$. Thus, the actual execution time for $\tau_1$ is equal to three and the slowdown factor $S_1$ is $S1 = 1/3 = 0.33$. Consequently, the the energy dissipation for $\tau_1$ is $E_1 = 12$ (see table I). Now, $\tau_1$ is executed from $t = 0$ to $t = 3$ with a slowdown factor $S_1 = 0.33$ where $E(3) = 388$ energy units.

$\tau_2$ has now the highest priority and it begins its execution at time $t = k_1 = 3$. Its execution time is stretched until $t = C_2 + \Delta_1 = 4$. The slowdown factor is then $S_2 = 0.5$ and the residual capacity at time $t = 7$ is $E(7) = 298$ energy units.

At time 7, $\Delta_2 = 0$. Consequently $\tau_1$ must be executed at full processor speed. $\tau_1$ executes completely until time 8 and consumes 30 energy units. The residual capacity then equals 268 energy units.

This procedure continues until $t = 18$ where no task requires to be processed in the interval $[18, 20]$. At the end of the hyperperiod, the battery capacity is equal to 28 energy units.

TABLE I.     ENERGY DISSIPATION OF TASKS $\tau_i$

| Energy Dissipation | $S = 1$ | $S = 0.9$ | $S = 0.8$ | $S = 0.7$ | $S = 0.6$ | $S = 0.5$ | $S = 0.4$ | $S = 0.3$ | $S = 0.2$ | $S = 0.1$ |
|---|---|---|---|---|---|---|---|---|---|---|
| Task $\tau_1$ | 30 | 27 | 24 | 21 | 18 | 15 | 12 | 9 | 6 | 3 |
| Task $\tau_2$ | 80 | 72 | 65 | 56 | 48 | 40 | 33 | 25 | 16 | 8 |
| Task $\tau_3$ | 180 | 160 | 140 | 125 | 110 | 90 | 70 | 55 | 36 | 18 |

## VI.    FEASIBILITY ANALYSIS

In contrast to EDF, ES-EDF feasibly schedules the task set in the first hyperperiod, with the same characteristics of the storage unit and the processor. It is important here to note that the processor remains busy as long as there are ready tasks in the queue. This means that the processor is busy from $[0, k_q]$. Thus, the processor utilization is equal to $1 - \Delta_q/T_{LCM} = 0.9$. This leads to deduce that ES-EDF is optimal with respect to minimizing processor energy.

THEOREM *1:* A task set $\Gamma$ that is schedulable with EDF remains schedulable with ES-EDF.

*Proof:* Let us consider the $U_p$ and $U'_p$ as processor utilization of EDF and ES-EDF respectively. If $\Gamma$ is schedulable by EDF, then $U_p \leq 1$. We have to prove that if $U_p \leq 1$, then $U'_p \leq 1$.

Suppose that $\Gamma$ is not schedulable by ES-EDF, then $U'_p > 1$. But $U'_p = \sum_{i=1}^{M} \frac{C_i(a)}{T_i} = \sum_{i=1}^{M} \frac{C_i + \Delta(St_i)}{T_i} = U_p + \sum_{i=1}^{M} \frac{\Delta(St_i)}{T_i}$. This implies that $U_p + \sum_{i=1}^{M} \sum_{j=1}^{q} \frac{\Delta_j}{T_i} > 1$. By multiplying the whole equality by $T_{LCM}$, we get $T_{LCM}U_p + \frac{T_{LCM}}{\sum_{i=1}^{M} T_i} \sum_{j=1}^{q} \Delta_j > T_{LCM}$.

But since $\Gamma$ is periodic, then $\frac{T_{LCM}}{\sum_{i=1}^{M} T_i} \sum_{j=1}^{q} \Delta_j = t_{idle}$, where $t_{idle}$ is the total idle time. By substitution, we get $T_{LCM}U_p + t_{idle} > T_{LCM}$, then $t_{idle} > T_{LCM}(1 - U_p)$

According to $EDL$, if $\Gamma$ is schedulable, then the total idle time during a whole hyperperiod $T_{LCM}$ is equal to $T_{LCM}(1 - U_p)$. Contradiction. Consequently, $U'_p$ must be less than or equal to one. ∎

COLLABORY *1:* The processor utilization $U'_p$ for ES-EDF is equal to $1 - \frac{min(x_i)}{T_{LCM}}$ where $x_i = T_i - D_i$.

*Proof:* Let $\mathcal{K}$ and $\mathcal{K}$ be respectively the static deadline vector and the static idle time vector, as defined in IV-B. According to ES-EDF, the busy period within the initial window, denoted $W(l)$, is equal to $K_q$ where $K_q$ is the latest component of $\mathcal{K}$. We note that $\Delta_q = T_{LCM} - k_q$ so that the length of busy period at $k_q$ is zero. As no task requires to be processed in such interval, the collabory is true. ∎

THEOREM *2:* Given a set of independent tasks with arbitrary computation times, deadlines and periods, ES-EDF is optimal with respect to minimizing the maximum lateness.

*Proof:* Given a set of independent tasks $\tau_i$ with arbitrary computation times, deadlines and periods. According to Horn [6], any schedule that puts the jobs in order of non-decreasing deadlines minimizes the maximum lateness. Alternatively, we concern ourselves with constructing the ES-EDF schedule that complete all tasks by their deadlines. To capture this, given the ES-EDF schedule that at every scheduling instant $k_i$ executes the task with the earliest absolute deadline among all the ready tasks and stretches its computation time $C_j$ until its deadline

$D_j$. Consequently, ES-EDF minimizes the maximum lateness. What is left in this theorem is the optimality criterion.

The maximum lateness is defined in [7] as $L_{max} = max_i(f_i - D_i)$ where $f_i$ is the finish time of the $i^{th}$ task. We can show that ES-EDF has a maximum lateness equal to zero. Let us consider $\tau_1$ to be the ready task with the highest priority at $t = k_0 = 0$. According to ES-EDF, $\tau_1$ will be executed until $t = D_1$ and hence its lateness is zero. If we repeat the above schedule until there are no more tasks in the queue, then we have constructed the ES-EDF schedule. Since this schedule has a maximum lateness $L_{max} = 0$, the ES-EDF schedule is optimal with respect to minimizing the maximum lateness. This optimum processing utilization is equal to $1 - \frac{\Delta_q}{T_{LCM}}$ where $T_{LCM}$ is the period, $\Delta_q = min\{x_j \mid 1 \leq j \leq n\}$ and $x_j = T_j - D_j$. ∎

## VII.    UPPER BOUND ON ENERGY SAVINGS

Scaling the processor frequency and voltage based on the performance requirements can lead to considerable energy savings. This is due to the quadratic relationship between voltage and dynamic power. In addition, the processor slowdown is based on slowdown factors $S_i$ in such a way that energy savings increase with decreased slowdown factor. However, the maximum energy savings possible by the any dynamic voltage and frequency selection algorithm must be bounded by the amount of energy savings with respect to the minimum slowdown factor (processing rate).

THEOREM *3:* [33] A set of periodic tasks is guaranteed to be schedulable with maximum energy savings iff the processing rate is

$$r_{min} = \sum_i \frac{C_i}{T_i} \qquad (15)$$

From this theorem, authors concluded that the maximum energy savings occurs when all tasks have the same averaged processing rate. This means that the minimum processor energy consumption occurs when all tasks are slacked by the same amount, to the maximum allowable limit such that $r_{min} = \sum_i \frac{C_i}{T_i}$.

Unfortunately, this conclusion is not always evident. To prove this, let us consider the following illustrative example: Consider a periodic task set $\Gamma$ that is composed of three tasks, $\Gamma = \{\tau_i \mid 1 \leq i \leq n\}$ and $\tau_i = (C_i, D_i, T_i)$. Let $\tau_1 = (1, 3, 5)$, $\tau_2 = (2, 7, 10)$ and $\tau_3 = (2, 12, 20)$. We assume that the energy storage capacity is $E = 350$ energy units at $t = 0$. The processor is assumed to be working with ten discrete slowdown factors $S_i = \{1, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1\}$. The power dissipation of tasks $\tau_i$ is shown in table I.

In order to demonstrate our theory, we have to schedule $\Gamma$ according to the assumption in [33] and then according to different processing rates as in ES-EDF.

Following the assumption in [33], the minimum processing rate (slowdown factor) is equal to $r_{min} = \sum_i \frac{C_i}{T_i} = 0.5$ and consequently, the computation time for each task instance is doubled (figure 3(a)). Since the same static slowdown factors are used, $\Gamma$ is scheduled till the end of the hyperperiod where the battery capacity is equal to 90 energy units.
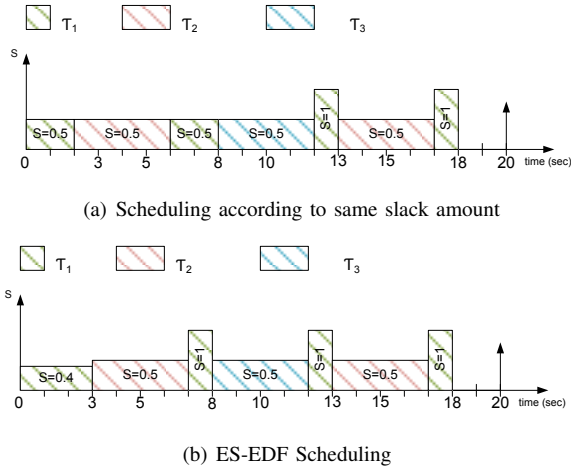


(a) Scheduling according to same slack amount



(b) ES-EDF Scheduling

Fig. 3. Scheduling of task set $\Gamma$

Now, let us try scheduling the same task set $\Gamma$ but with ES-EDF (figure 3(b)). We find that $\Gamma$ is schedulable since all tasks are executed without violating deadlines and the remaining energy in the battery at the end of the hyperperiod is equal to 108 energy units.

From this example, we can deduce that ES-EDF feasibly schedules the task set in the first hyperperiod and with more energy saving than when considering the same averaged processing rate.

THEOREM 4: A set of periodic tasks is guaranteed to be schedulable with maximum energy savings iff the optimal slowdown factor is approximated by

$$S_{opt} = (1 - \frac{x_q}{T_{LCM}})U_p \qquad (16)$$

*Proof:* It has been proved in [9] that a periodic task set is guaranteed to be schedulable by EDF *iff* $\sum_i \frac{C_i}{T_i} \leq 1$. Theorem 3 demonstrated that the maximum allowable limit, in case of equal slowdown factors, is bounded by $E_{save}(S_{min})$ where the minimum slowdown factor $S_{min} = \sum_i \frac{C_i}{T_i}$. We will now show the minimum energy consumption under different slowdown factors. Let us consider $S_i$ and $S_{opt}$ as the slowdown factor of task instance $\tau_i$ and the optimal slowdown factor respectively. From equation (14), $S_i = \frac{C_i}{C_i(a)}$. We can derive the condition that $S_{opt}$ has to satisfy in order to guarantee the schedulability of the task set, then $S_{opt} = \sum_i \frac{C_i}{C_i(a)}$. But since tasks are periodic, then we have to multiply by $\frac{T_{LCM}}{\sum_i T_i}$. Therefore, the average of slowdown factors is equal to

$$S_{opt} = \frac{1}{T} \sum_i \frac{C_i}{(C_i(a))(T_i)} \qquad (17)$$

Now, let's consider the following inequality that can be verified using the Cauchy-Schwarz inequality

$$(\sum_i \frac{C_i}{(C_i(a))(T_i)})(\sum_i C_i(a)) \geq \sum_i \frac{C_i}{T_i} \qquad (18)$$

By making some arrangement and distribution of the terms, we get

$$S_{opt} \geq (1 - \frac{x_q}{T_{LCM}})U_p \qquad (19)$$

Then, the minimum speed $S_{opt}$ that ensures feasibility is $S_{opt} = (1 - \frac{x_q}{T_{LCM}})U_p$. ∎

## VIII. PERFORMANCE EVALUATION

This section provides performance evaluation of the algorithms. Algorithms under simulations are ES-EDF, EDF and Enhanced EDF (E-EDF). E-EDF is an enhanced version of EDF in a way that tasks are slacked by the same slowdown factor $S = \sum_i \frac{C_i}{T_i}$.

### A. Experimental Setup

We implemented the proposed scheduling techniques in a discrete event simulator using C/C++. To evaluate the effectiveness of the ES-EDF algorithm, we consider a task generator of periodic tasks based on that described by Martineau in [34]. It accepts as input several parameters: the number of desired tasks $n$, the hyperperiod of task periods $T_{LCM}$ and processor utilization $U_p$. At the output, we obtain a task configuration of the scheduled task set. The execution times of tasks are randomly generated such that $\sum_i \frac{C_i}{T_i} \leq 1$. The simulator generates 30 tasks with least common multiple of the periods equal to 3360. The worst-case computation times are set according to the processor utilization $U_p$. Deadlines are less than or equal to periods and greater than or equal to the computation times ($C_i \leq D_i \leq T_i$).

To estimate the processor energy consumption, we use Intel XScale processor supporting five frequency levels [35]. The discrete frequencies, supply voltage and consumed power the processor are listed in Table II. We assume that the energy

TABLE II. XSCALE FREQUENCIES, SUPPLY VOLTAGES, AND POWER

| Frequency (MHz) | 150 | 400 | 600 | 800 | 1000 |
|---|---|---|---|---|---|
| Power (mW) | 80 | 170 | 400 | 900 | 1600 |
| Voltage (V) | 0.75 | 1.0 | 1.3 | 1.6 | 1.8 |

storage is fully charged at the beginning of the simulation. After a deadline violation is detected, the simulation terminates for ES-EDF, EDF and E-EDF.

### B. Percentage of Feasible Task Sets by Varying $U_p$

Our simulation depicts the percentage of feasible task sets by varying the processor utilization $U_p$. Here, we take interest in the percentage of task sets which are feasible with ES-EDF, EDF and E-EDF. We report the results of this simulation study where $U_p$ varies from 0.1 till 1 (figure 4).

Under low processor utilization, we observe that ES-EDF maintain 100% of feasible task sets, and exceeds that of E-EDF and EDF by about 31% and 49%. This is due to the
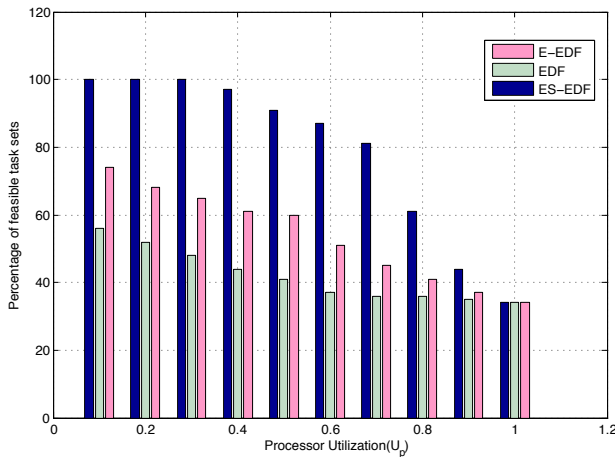
Fig. 4.   Percentage of Feasible Task Sets by Varying $U_p$

fact that the EDF algorithm runs at full processor speed and does not utilize DVS technique to save energy. On the other hand, in E-EDF, tasks are slacked with equal slowdown factor and consequently not all deadlines will be met. Hence, ES-EDF system significantly reduces the deadline miss rate due to energy shortage. In other words, the more stored energy means that more tasks are able to be finished before their deadlines.

As $U_p$ increases, the percentage of feasible task sets in ES-EDF decreases. This is because ES-EDF utilizes the slack time to slow down the execution for energy savings. Under high values of $U_p$, slack time decreases and most of tasks are just executed at the full speed, as EDF algorithm does. Hence, the ES-EDF system incurs much higher deadline miss rate but still exceeds that of E-EDF and EDF by about 34% and 46%.

It is important to note that when the processor utilization is set to one, ES-EDF, E-EDF and EDF has exactly the same percentage of feasible task sets. This is because the processor is always active and there is no processor idle time.

### C. Percentage of Feasible Task Sets by Varying Battery Capacity

This experiment depicts the percentage of feasible task sets over the energy storage capacity $E$. Here, we take interest in the percentage of task sets which are feasible with ES-EDF and not feasible with E-EDF and EDF. We report the results of this simulation study where the processor utilization $U_p$ is set to $0.4$ and $0.8$ respectively.

For each task set, we compute $E_{feas}$ as the minimum storage capacity which permits to achieve neutral operation according to ES-EDF. i.e. all tasks are executed without violating deadlines and the battery is not empty at the end of the hyperperiod. After that, we vary $E$ with $E > E_{feas}$ so as all task sets are feasible with E-EDF and EDF.

When $U_p$ is set to $0.4$ (figure 5 (a)), we observe that the battery capacity must be about 2.2 and 3.6 times bigger with E-EDF and EDF to maintain 100% feasible task sets compared to ES-EDF. This is due to the fact that EDF algorithm runs at full processor speed and does not utilize DVS technique to save energy. Thus, the SE-EDF and E-EDF systems significantly reduce the deadline miss rate due to energy shortage.

Unfortunately, E-EDF considers equal slowdown factors and this will reduce the percentage of feasible task sets.

In other words, the more stored energy means that more tasks are able to be finished before their deadlines. Hence, the ES-EDF system incurs much higher deadline miss rate when compared to E-EDF and EDF.

Figure 5 (a) shows that the deadline miss rate of E-EDF and EDF exceeds that of the proposed scheduling algorithm by about 45% and 61% when battery capacity is respectively the same. Hence, the ES-EDF algorithm is favorable even for small battery capacity.

If we increase the processor utilization to $0.8$ and run the simulation again, we find that the gain in capacity savings is decreasing (figure 5 (b)). ES-EDF obtains capacity savings of about 24% and 41% compared to E-EDF and EDF respectively. The decrease in capacity savings can be attributed to the fact that as $U_p$ increases, the slacking gets harder and the SE-EDF schedule tends to the EDF schedule with processor utilization increasingly being set to 1.

### D. Ratio of Energy Savings

We present in this section the energy gains achieved by ES-EDF when compared to E-EDF. Experiments were performed on task sets with varying processor utilization ($U_p$). Figure 6 shows the normalized energy gains for ES-EDF and EDF, that means the quantity of energy gains relative to battery capacity.
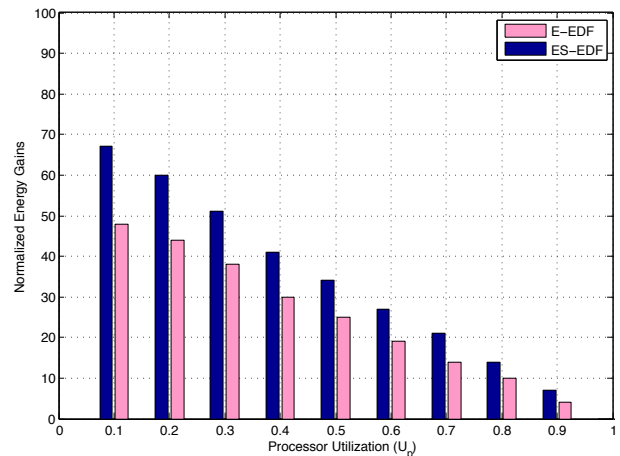


Fig. 6.   Ratio of Energy Savings of ES-EDF, E-EDF and EDF

From the beginning, we state that there are no energy gains in EDF since it operates at maximum processor frequency. As for ES-EDF and E-EDF, when the task execution time is decreased, there is more slack which results in decreasing the energy consumption by operating at a lower voltage. Furthermore, the average energy gains in ES-EDF is about 28% more than E-EDF. This is due to the fact that ES-EDF stretches tasks as much as possible while still guaranteeing deadlines and consequently the processor is always active. On the other hand, tasks in E-EDF are slacked with equal slowdown factor and consequently not all deadlines will be met.
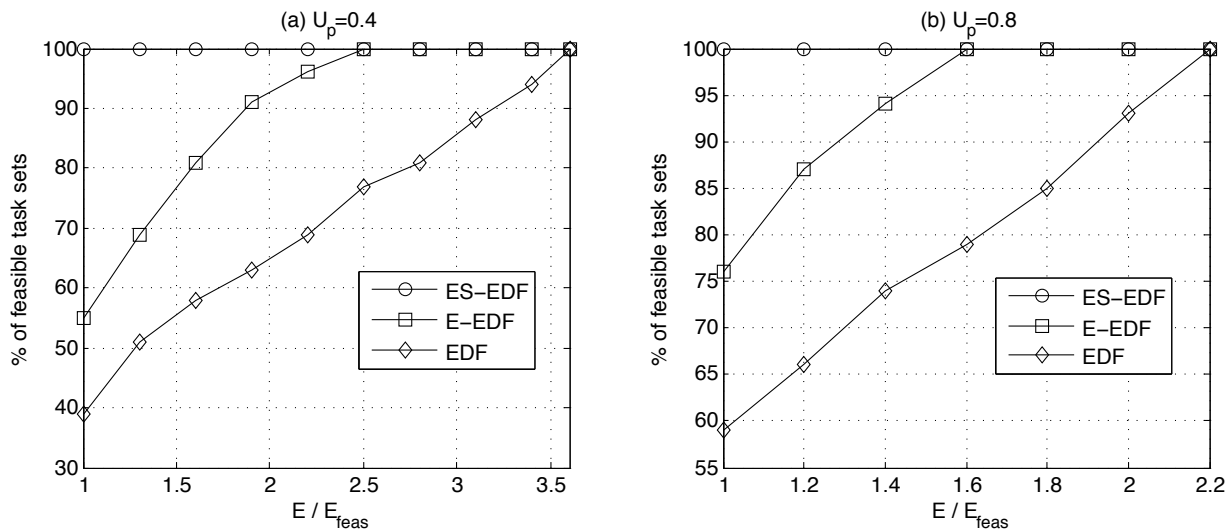
Fig. 5.    Percentage of Feasible Task Sets by Varying Battery Capacity

## IX.    CONCLUSION

In this paper, we considered the problem of developing an energy saving algorithm for periodic task sets characterized by real-time deadlines using variable voltage and frequency assignments on a monoprocessor system. To this end, we proposed and Energy Saving EDF (ES-EDF) algorithm that is proved to be optimal in minimizing the maximum lateness and the processor energy consumption. In addition, we demonstrated through an illustrative example that it is not necessary that all tasks must be slacked by the same amount so as to obtain a minimum processor energy consumption. We then determined the optimal scaling factor by which a task should be stretched to maximize energy savings while still respecting all deadline constraints. The proposed algorithm achieved an average energy savings of about 28% when compared to EDF. Further, ES-EDF yields higher percentage of feasible task sets which exceeds that of E-EDF and EDF by about 33% and 47%.

We are currently looking at extending the proposed energy saving scheduling algorithm to operate at multiprocessor systems.

### ACKNOWLEDGEMENTS

### REFERENCES

[1]

[2]    G. Simon, M. Maroti, and A. Ledeczi, *Sensor Network-Based Countersniper System*.    In Proceedings of the 2nd ACM Conference on Embedded Network Sensor Systems (SenSys04), Baltimore, Maryland, November, 2004.

[3]    C. Hartung, R. Han, C. Seielstad, and S. Holbrook, *FireWxNet: A Multi-Tiered Portable Wireless System for Monitoring Weather Conditions in Wildland Fire Environments*.    In Proceedings of the 4th International Confernce on Mobile Systems, Applications, and Services (MobiSys06), Uppsala, Sweden, June, 2006.

[4]    G. Werner-Allen, K. Lorincz, J. Johnson, J. Lees, and M. Welsh. *Fidelity and yield in a volcano monitoring sensor network*.    In Proceedings of the 7th USENIX Synposium on Operating Systems Design and Implementation (OSDI06), Seattle, Washington, November, 2006.

[5]    J. Zhuo and C. Chakrabarti, *Energy-efficient Dynamic Task Scheduling for DVS Systems*.    ACM Transactions on Embedded Computing Systems, February, 2008.

[6]    W. Horn, *Some Simple Scheduling Algorithms. Naval Research Logistics Quaterly*, 21, 1974.

[7]    ML. Dertouzos, *Control robotics: the procedural control of physical processes*.    Proceeding of International Federation of Information Processing Congress, 1974.

[8]    M. Silly-Chetto, *The EDL Server for scheduling periodic and soft aperiodic tasks with resource constraints*.    Real-Time Systems, 17(1), pp.1-25, 1999.

[9]    C-L Liu , J-W Layland, *Scheduling algorithms for multiprogramming in a hard real-time environment*.    Journal of ACM, 46-61, 1973.

[10]    Dertouzos ML, *Control robotics: the procedural control of physical processes*.    Proceeding of International Federation of Information Process Cong, 80713, 1974.

[11]    J-Y-T Leung , J. Whitehead, *On the complexity of fixed-priority scheduling of periodic real-time tasks*.    Performance Evaluation, 1982.

[12]    H. Chetto, and M. Chetto, *Some results of the earliest deadline scheduling algorithm*.    IEEE Transactions on Software Engineering, 15(10): 1261-1269, 1989.

[13]    Yao, F., Demers, A.J., Shenker, S., *A scheduling model for reduced CPU energy*.    In: Proceedings of IEEE Symposium on Foundations of Computer Science, 374382, 1995.

[14]    T. Ishihara and H. Yasuura, *Voltage scheduling problem for dynamically variable voltage processors*.    In Proceedings of the International Symposium on Low Power Electronics and Design, Monterey, CA, pp. 197-202, 1998.

[15]    Quan, G., Hu, X., *Minimum energy fixed-priority scheduling for variable voltage processors*.    In Proceedings of Design Automation and Test in Europe, 2002.

[16]    Kwon, W., Kim, T., *Optimal voltage allocation techniques for dynamically variable voltage processors*.    In Proceedings of the Design Automation Conference, 125130, 2003.

[17]    Aydin, H.,Melhem, R.,Mossé, D., Alvarez, P.M., *Determining optimal processor speeds for periodic real-time tasks with different power characteristics*.    In Proceedings of EuroMicro Conference on Real-Time Systems, 2001.

[18]    Jejurikar, R., Gupta, R., *Optimized slowdown in real-time task systems*. In Proceedings of EuroMicro Conference on Real-Time Systems, June, 2004.

[19] Aydin, H., Melhem, R., Mossé, D., Alvarez, P.M., *Dynamic and aggressive scheduling techniques for power-aware real-time systems*. In Proceedings of IEEE Real-Time Systems Symposium, 2001.

[20] Zhang, F., Chanson, S.T., *Processor voltage scheduling for real-time tasks with non-preemptible sections*. In Proceedings of IEEE Real-Time Systems Symposium, Dec., 2002.

[21] Kim,W., Kim, J., Min, S.L., *A dynamic voltage scaling algorithm for dynamic-priority hard real-time systems using slack time analysis*. In Proceedings of Design Automation and Test in Europe, March, 2002.

[22] Y. Shin and K. Choi, *Power Conscious Fixed Priority Scheduling for Hard Real-Time Systems*. In Proceedings of the 36th Design Automation Conference, DAC99, 1999.

[23] Arezou Mohammadi and Selim G. Akl, *Scheduling algorithms for real-time systems*. Technical report no. 2005-499, 2005.

[24] Y. Shin, K. Choi, and T. Sakurai, *Power optimization of real-time embedded systems on variable speed processors*. In Proceedings of the 2000 IEEE/ACM International Conference on Computer-Aided Design, pages 365-368, 2000.

[25] P. Mejia-Alvarez, E. Levner, and D. Mosse, *Adaptive scheduling server for power-aware real-time tasks*. ACM Transactions on Embedded Computing Systems, 3(2):284-306, 2004.

[26] J.-J. Chen, T.-W. Kuo, and C.-S. Shih, $1 + \epsilon$ *approximation clock rate assignment for periodic real-time tasks on a voltage-scaling processor*. In the 2nd ACM Conference on Embedded Software (EMSOFT), pages 247-250, 2005.

[27] D. Zhu and H. Aydin, *Reliability-aware energy management for periodic real-time tasks*. In Proc. of the IEEE Real-Time and Embedded Technology and Applications Symposium, 2007.

[28] D. Zhu, X. Qi, and H. Aydin, *Priority-monotonic energy management for real-time systems with reliability requirements*. In Proc. of the IEEE International Conference on Computer Design (ICCD), 2007.

[29] Dakai Zhu, Xuan Qi and Hakan Aydin, *Energy Management for Periodic Real-Time Tasks with Variable Assurance Requirements*. In Proceedings of the 2008 14th IEEE International Conference on Embedded and Real-Time Computing Systems and Applications. Pages 259-268, 2008.

[30] Ravindra Jejurikar and Rajesh Gupta, *Energy-Aware Task Scheduling With Task Synchronization for Embedded Real- Time Systems*. IEEE Transactions on Computer Aided Design of Integrated Circuits and Systems, vol. 25, no. 6, JUNE, 2006.

[31] Y.-S. Chen, C.-Y. Yang, and T.-W. Kuo, *Fl-pcp: Frequency locking for energy-efficient real-time task synchronization*. In the 13th IEEE International Conference on Embedded and Real-Time Computing Systems and Applications (RTCSA), 2007.

[32] Ya-Shu Chen and Tay-Jyi Lin, *Voltage Emergence Prevention for Energy-Efficient Real-Time Task Synchronization*. IEEE 8th International Conference on Computer and Information Technology Workshops, 2008.

[33] Amit Sinha and Anantha P. Chandrakasan, *Energy Efficient Real-Time Scheduling*. Asia-South Pacific Design Automation Conference Proceedings (ASPDAC), 2001.

[34] P. Martineau: Ordonnancement en-ligne dans les systèmes informatiques temps-réel. PhD Dissertation, University of Nantes, 1994.

[35] Intel Corp, *Intel XScale Processor Family Electrical, Mechanical, and Thermal Specification Datasheet*. Santa Clara, CA, USA, 2004.

# Mobile Receiver-Assisted Localization
# Based on Selective Coordinates in Approach to
# Estimating Proximity for Wireless Sensor Networks

Zulfazli Hussin
Graduate School of Applied Informatics
University of Hyogo
Kobe, Hyogo, Japan 650-0047
Email: zulfazlihussin AT gmail.com

Yukikazu Nakamoto
Graduate School of Applied Informatics
University of Hyogo
Kobe, Hyogo, Japan 650-0047
Email: nakamoto AT ai.u-hyogo.ac.jp

*Abstract*—Received signal strength (RSS)-based mobile localization has become popular due to its inexpensive localization solutions in large areas. Compared to various physical properties of radio signals, RSS is an attractive approach to localization because it can easily be obtained through existing wireless devices without any additional hardware. Although RSS is not considered to be a good choice for estimating physical distances, it provides some useful distance related information in adding and indicating connectivity information in neighboring nodes. RSS-based localization is generally divided into range-based and range-free. Range-based localization can achieve excellent accuracy but is too costly to apply to large-scale networks. Methods of range-free localization are regarded as cost-effective solutions for localization in sensor networks. However, the localizations are subject to the effect of radio patterns that affect variations in the radial distance estimates between nodes. It is a challenging task to select an efficient RSS value that can provide small variations in the radial distance in wireless environments. We propose a method of Mobile Localization using the Proximities of Selective coordinates (MoLPS) to localize target nodes by using information on proximities between target nodes and mobile receivers as a metric to estimate the location of target nodes. We ran a simulation experiment to assess the performance of MoLPS with 100 target nodes that were randomly deployed along a sensory field boundary. We found from the results of the simulation experiment that localization error had been reduced to below 2m in more than 80% of the target nodes.

*Keywords*—*Localization, proximity estimation, genetic algorithm, wireless sensor networks, received signal strength.*

## I. INTRODUCTION

Wireless sensor networks (WSNs) [1] are composed of many sensor nodes that have sensing and computational and wireless communication capabilities. Although WSNs have demonstrated their importance and capabilities in emergency applications, if the positions of sensor nodes are known, the use of these applications could be even more effective.

Localization is fundamentally a serious problem that deals with how to use information from sensor nodes to determine position coordinates. Locating an item is a critical process at distribution centers since poor performance results in unsatisfactory customer services (long processing and lagged delivery) and high costs. Suppose that a sensor node is attached to an item at a distribution center. Although placing an item

at a fixed location makes it easier to locate it, it is not always the most space-efficient method of storage for products that are less predictable due to uncertain demand [2]. In contrast, random-location storage uses less storage space even though it requires the use of a locator to identify the locations of items. A straightforward solution would be to equip all sensor nodes with GPS receivers that could provide them with the exact locations of items. However, this is not a cost-effective solution and it has limited applications because GPS only works in open areas with no obstructions to satellite signals.

Receiver-assisted localization has attracted a great deal of attention in estimating the positions of items that are equipped with sensor nodes. Receivers detect sensor nodes by using the radio signals received from them. There are generally two types of deployments used to detect sensor nodes. The first is to fix several receivers that cover particular regions [3]. Thus, the numbers of receivers and their distributions have a direct impact on the accuracy of localization. A large number of distributed receivers will lead to improved accuracy. However, costs will be high if they are applied to large areas. The second method is to use mobile receivers to sense locations. Since mobile receivers are portable and easy to use, they are suitable for location sensing in large areas (e.g. distribution centers).

On this basis, we reassess existing localization scheme and exlore the possibility of using selected coordinates of mobile receivers. We evaluated the concentration of center coordinates, which are computed from the selected coordinates of mobile receivers, to estimate the position of a target node without deploying fixed receivers or fixed anchors. We called the coordinates of the mobile receivers footprints. We divided the footprints into multiple sets in which each set represented the footprints that received signals in the given range of path loss values for each set. Path loss describes a signal's energy loss that varies continuously as it travels to a receiver [4]. Instead of selecting all the footprints to compute the average for each set, we selected footprints that had fewer variations in the radial distance to the true target node in each set. The center coordinates of selected footprints were individually computed and the concentration of center coordinates of selected footprints were evaluated to estimate the true position of the target node.

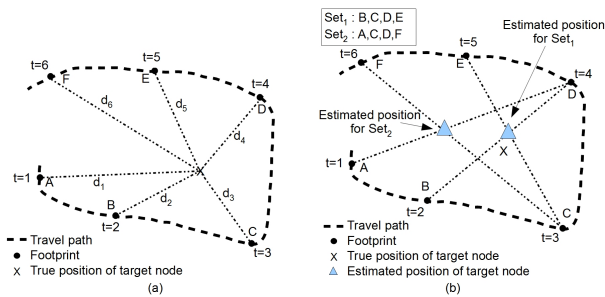The rest of the paper is organized as follows: The motiva-

Fig. 1: (a) Variations in distance using centroid and (b) comparison of localization accuracy with two different selection of footprints

TABLE I: Set of footprints.

| Set, $q$ | Path Loss [dB], $PL$ |
|---|---|
| 1 | $41.5 < PL \leq 42.5$ |
| 2 | $42.5 < PL \leq 43.5$ |
| 3 | $43.5 < PL \leq 44.5$ |
| 4 | $44.5 < PL \leq 45.5$ |
| 5 | $45.5 < PL \leq 46.5$ |
| 6 | $46.5 < PL \leq 47.5$ |
| 7 | $47.5 < PL \leq 48.5$ |
| 8 | $48.5 < PL \leq 49.5$ |
| 9 | $49.5 < PL \leq 50.5$ |
| 10 | $50.5 < PL \leq 51.5$ |

tion for our research is presented in Section 2. We investigated state-of-the-art range-based and range-free localization techniques, which are described in Section 3. Section 4 describes the problem of how to localize target nodes by using our proposed algorithm. Section 5 describes our overall algorithm. Our evaluation of performance is described in Section 6. and are followed by closing remarks in Section 7.

## II. MOTIVATION

Mobile receivers are used in mobile localization to measure the physical properties of collected radio signals from target nodes. Localization based on received signal strength (RSS) has become popular because it is an inexpensive solution to the problem of localizing target nodes. Compared to various physical properties of radio signals, such as Time of Arrival (ToA) [5], Time Difference of Arrival (TDoA) [6] or Angle of Arrival (AoA) [7], RSS is an attractive approach to localization because it can easily be obtained through existing wireless devices without the need for any additional hardware. The major challenge to accurate RSS-based positioning results from the variations in RSS that change over time and space due to dynamic and unpredictable signal propagation. Although RSS is not considered to be a good choice for estimating physical distances in many scenarios that involve unknown radio path loss factors, hardware discrepancies, and antenna orientation [8], [9], it provides useful information that is distance related in addition to indicating connectivity information between neighboring nodes. Many techniques of RSS-based localization have been proposed in the past two decades. They generally fall in two categories of range-based and range-free.

RSS readings are used in range-based localization techniques [8], [9] to directly estimate the physical location of target nodes. However, RSS measurements are easily corrupted by surrounding environments. Moreover, techniques of range-based localization require expensive and power-intensive measuring devices or synchronization that may incur cost and energy problems. In contrast, range-free approaches have been proposed as an alternative to pursue cost and energy effectiveness in WSNs [10], [11]. Range-free approaches are typically used for connectivity between nodes as a metric to estimate the position of target nodes without computing the actual distance between nodes.

The accuracy of localization in range-free approaches is subject to the effect of radio patterns that affect variations in estimates of the radial distance between nodes. Many of these techniques use an average of all anchor positions in their communication range [10] or in the same hop-count values [11] to localize the target nodes, which results in variations in the radial distance being underestimated thereby causing large localization errors. Localization errors vary between estimates of target nodes caused by variations in the radial distance that result from target nodes that have not been uniformly deployed. It is a challenging task to select efficient RSS values that can provide small variations in the radial distance from accumulated RSS values in wireless environments where complex and dynamic RSS values can affect the estimates of radial distances.

We propose a method of Mobile Localization using Proximities of Selective coordinates (MoLPS) to localize target sensor nodes by using the connectivity between them and mobile receivers as a metric to estimate their locations to solve these stated challenges. Locations are estimated from the coordinates of footprints where the signals are collected from target sensor nodes. MoLPS assume the presence of a tentative coordinate that is arbitrarily located at a known location in the field in which it is deployed. The distance between a tentative coordinate and center coordinates of selected footprints are iteratively compared to select effective footprints to localize target nodes. We used a genetic algorithm (GA) to search the best selection of footprints that had fewer variations in the radial distance from other selected footprints to the tentative coordinates. We iteratively improved the positions of the tentative coordinate by evaluating the concentration of center coordinates of selected footprints in the vicinity of the tentative coordinate.

## III. RELATED WORKS

Theoretical or empirical models are used in range-based localization techniques to translate RSS into estimates of distance. Range-based localization can achieve better accuracy but is costly in requiring either per-node ranging hardware [12] or careful system calibration and environment profiling [13], [14], and thus it is not appropriate for large-scale sensor networks. The correlation of noise due to shadowing from obstacles in wave propagation has been exploited to estimate the locations of transmitters [15]. Cumulative errors in measurement with positioning methods have been treated as problems with localization where data sampled over time have generated points in high dimensional space [16], [17], [18]. The multi-dimensional scaling (MDS) model has been used to reduce dimensionality to estimate locations [16]. However, the

linear relationship requirement between correlation coefficients and radial distance in MDS has restricted its applications to wireless environments where RSS correlations are highly nonlinear if there is a radial distance [19] between receivers. Manifold learning (reduced nonlinear dimensionality) algorithms such as Isomap, Local Linear Embedding (LLE) and Hessian LLE have been used to centralize localization [17], [18]. The linearity between correlation measurements and radial distance is restricted in these approaches to a small area containing K nearest neighbors. However, the linearity between RSS and radial distance does not hold in Li and Liu [19], even in the immediate vicinity of operating frequencies greater than 10 MHz.

Range-free approaches localize nodes based on simple sensing, such as wireless connectivity [11], [20], [21] and anchor proximity [10], [22], [23]. Wireless connectivity information between neighboring nodes is used to estimate the location of a target node by using MDS [20]. Their major limitation is that they all rely on a large number of uniformly-distributed anchors in the networks. Embedding the combinatorial Delaunay complex in the landmark Voronoi diagram [21] has improved the localization of target nodes in various network topologies. However, using a number of landmarks to achieve precise accuracy in localization is costly.

The approximate-point-in-triangulation (APIT) algorithm [22] was proposed for area-based range-free localization, where all sensor nodes were localized by using the location information of GPS-equipped anchors. The areas occupied by sensor nodes were divided into many triangular regions between anchors in this approach by using the location information provided by GPS. This approach provided excellent accuracy when irregular radio patterns and random node placements were considered. Moreover, the large number of distributed anchors will counteract problems such as high deployment costs when applied to large areas. In Centroid [10], all possible anchors broadcast their location information to all other target nodes. The target nodes use the location information from anchors that are located in their vicinity to estimate their own location coordinates. The main difficulty with the centroid is the large number of anchors to be considered in the estimates. Moreover, if anchors are not uniformly distributed, the distance between them and target nodes varies, which deteriorates the accuracy of localization. It is necessary to take into consideration the distance between anchors and target nodes to solve this problem. The distances between anchors and target nodes are considered in the distance vector-hop (DV-hop) localization algorithm [11] and resilient Ethernet protocol (REP) [23] as a form of hop counting, which is a range-free approach that does not use RSS to compute the distance between nodes. DV-hop performs well when deployed sensor nodes have regular node density and distances between them. However, the resulting estimates may not be optimal if the radio patterns are irregular and random node deployment is used in practice.

In MoLPS, the coordinates of anchors were determined from the selected coordinates of footprints by using our proposed method. Instead of selecting all the anchors to estimate the location of target nodes, we select the anchors that had fewer variations in the radial distance which can minimize the uncertain radial distance contamination problem to the
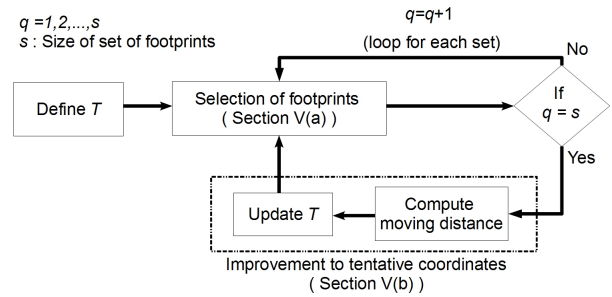


Fig. 2: Selection of footprints and improved $T$ in MoLPS.

localization of target nodes.

## IV. PROBLEM DEFINITION

We assumed that a mobile receiver would travel within a sensory boundary field that was deployed by sensor nodes that were transmitting their signals to the receiver on a periodic basis in this research. The mobile receiver is traveling at a constant speed while it is collecting signals from target nodes at $\tau$ intervals. The mobile receiver and target nodes are capable of communicating within their communication range $D$. Every time the mobile receiver receives a signal from target node $i$, it measures the RSS value of the signal and then stores it as tuple $(t, r_{i,t})$, where $t$ is the time denoted as $t = t_1, t_2, \ldots, t_\tau$ and $r_{i,t_j}$ is an RSS value denoted as $r_{i,t_j} = r_{i,t_1}, r_{i,t_2}, \ldots, r_{i,t_\tau}$. Each tuple contains a different RSS value, each of which is collected from a different position of the footprint in each $t$.

We assumed that the ranging levels of RSS received from a target node would decrease due to path loss effects as the distances between each footprint and target node increased. All footprints $P_{k_q,q} = \left( X_{k_q,q}, Y_{k_q,q} \right)$ were divided into $s$ sets according to the path loss values. Here, $P_{k_q,q}$ denotes $k_q$-th footprints in set $q$ where $k_q = 1, 2, \ldots, m_q$ and $q = 1, 2, \ldots, s$ as listed in Table I.

We took into consideration noisy environments in measuring RSS that contributed to variations in radial distances between footprints that received the same RSS values from a target node. The average of the positions from surrounding footprints are used in the centroid to estimate the locations of target nodes [10]. Accuracy with this approach greatly depends on variations in the radial distances of a target node at each footprint. The propagation of wireless signals is ideal in a noise-free environment, such that a target node can communicate with a mobile receiver from any footprint that is located within a perfect sphere centered on the target node and with a radius equal to its standard interrogation range. It is possible in this case to estimate the position of a sensor node by averaging all footprint coordinates that are located within its radius. However, it is difficult to guarantee whether the radial distance of the target node will be accurate at each footprint in practice in noisy environments. Moreover, since the position of a true target node is unknown, there is no way of selecting footprints that have fewer variations in distance to estimate the position of a target node.

It is necessary to select footprints that have fewer variations in radial distance in range-free mobile localization that relies
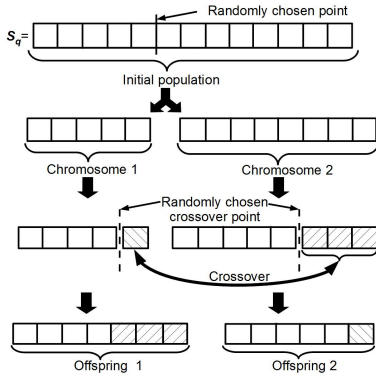
Fig. 3: The chromosome representation and generation process in GA.

**Algorithm 1** Selection of footprints with GA

1: $r \leftarrow rand()$
2: $S1 \leftarrow Chromosome1 = \{P_{1,q}, P_{2,q}, \ldots, P_{r,q}\}$
3: $S2 \leftarrow Chromosome2 = \{P_{r+1,q}, P_{r+2,q}, \ldots, P_{m_q,q}\}$
4: $F1 \leftarrow$ Distance between center of $S1$ and $T$
5: $F2 \leftarrow$ Distance between center of $S2$ and $T$
6: **while** $loop < 50$ **do**
7:    **if** $F1 < F\_Best$ **then**
8:       $S\_Best \leftarrow S1$
9:       $F\_Best \leftarrow F1$
10:    **end if**
11:    **if** $F2 < F\_Best$ **then**
12:       $S\_Best \leftarrow S2$
13:       $F\_Best \leftarrow F2$
14:    **end if**
15:    $crossover() \leftarrow$ Generation of offspring
16:    $S1 \leftarrow Offspring1$
17:    $S2 \leftarrow Offspring2$
18:    $F1 \leftarrow$ Distance between center of $S1$ and $T$ after crossover
19:    $F2 \leftarrow$ Distance between center of $S2$ and $T$ after crossover
20:    $loop \leftarrow loop + 1$
21: **end while**

on the average of the footprints as the positions of the estimates. For example, the mobile receiver in Fig.1(a) is traveling around the sensor node field and is collecting signals at each time interval and $d_1, d_2, \ldots, d_6$ denote the radial distances of a target node to all footprints. Assuming that there are two sets of footprint coordinates in Fig.1(b), localization using the information from $Set_1$ is more accurate than that from $Set_2$. The small variations in radial distance between the target node and each footprint in $Set_1$ contribute to greater accuracy. Therefore, it is important to select appropriate footprints that have fewer variations in radial distance to estimate the positions of target nodes and to obtain accurate localization.

## V. DESCRIPTION OF ALGORITHM

We iteratively improved the position of tentative coordinates $T = (X^T, Y^T)$ in close proximity to the true position of a target node to estimate the position of the target node. Here, we used $T$ to select the collection of footprints that had a center coordinate nearest to $T$. The algorithm we propose is outlined in Fig. 2. First, we arbitrarily define $T$ in the sensory boundary field without any knowledge of the position of the target node. Then, we select the collection of footprints that have the nearest center coordinates to $T$ for each set $q$. After the footprints have been selected, we determine the moving distance for $T$ by evaluating the concentration of plotted center coordinates of selected footprints in the vicinity of $T$ by iteratively improving $T$ until the number of cycles of improvements to obtain the best solution is satisfied.

### A. Selection of footprints

We selected collections of footprints that had fewer variations in radial distance to obtain accurate estimates of a target node with MoLPS. We used a genetic algorithm (GA) approach in this research to select $P_{k_q,q}$ by searching the nearest center coordinates of selected footprints to a tentative coordinate, $T$. The pseudocode for the GA in selecting the footprints is given in Algorithm 1. GA is a search algorithm that searches an optimal solution to solve a combinatorial problem, such as the NP-complete traveling salesman problem (TSP). The solution to a given problem is represented as a chromosome in GA. A population of solutions is created, and operators such as mutation and crossover are applied to derive the solutions.

The relative accuracies (fitnesses) of the solutions are then compared to find the best solution.

Assuming we have a total of $m_q$ footprints in a $q$ set, we determine all footprints in $P_{k_q,q}$ as an initial solution, $S_q = \{P_{1,q}, P_{2,q}, \ldots, P_{m_q,q}\}$. We determine the center coordinates of $S_q$ as $C_q$, which is computed as:

$$C_q = \left( \frac{\sum_{k_q=1}^{m_q} X_{k_q,q}}{m_q}, \frac{\sum_{k_q=1}^{m_q} Y_{k_q,q}}{m_q} \right) \quad (1)$$

The chromosome representation and the generation process in GA is outlined in Fig.3. We let $S_q$ represent the initial population in GA and compute the distance between $C_q$ and $T$ as the initial fitness of GA, $F_q$ as:

$$F_q = \sqrt{(X^T - X^{C_q})^2 + (Y^T - Y^{C_q})^2} \quad (2)$$

where $(X^{C_q}, Y^{C_q})$ denotes the center coordinate $C_q$ for $S_q$. We divide $S_q$ into two sets of footprints at randomly chosen points to represent two chromosomes in GA and compute the fitness of each chromosome as $F_{1,q}$ and $F_{2,q}$ using the same Eq. (2). The best solution is selected from a pair of chromosomes based on the least fitness (i.e., nearest distance) by comparing $F_{1,q}$ and $F_{2,q}$. A chromosome that has better fitness is selected as the best solution, $S_q^{best}$.

Crossover is applied to generate offspring chromosomes from dominant parent chromosomes. The crossover operator separates each chromosome into two sets at randomly chosen crossover points and exchanges separate sets to form new offspring. If crossover does not occur, the new offspring are exact copies of their parent chromosomes. Then, GA is repeated using the new offspring until the number of iterations of the computation satisfies a bound (e.g., fifty).

We assume the selected footprints at the last iteration of GA to be the best selection of footprints $S_q^{Best} = \{P_{1,q}, P_{2,q} \ldots, P_{n_q,q}\}$ where $n_q$ is the total number of selected footprints. We search $S_q^{Best}$ for each set of footprints and use the center coordinates $C_q^{Best}$ of $S_q^{Best}$, which are computed with Eq. (3) to improve the position of $T$ as will be described in the next subsection.

$$C_q^{Best} = \left(\frac{\sum_{k_q=1}^{n_q} X_{k_q,q}}{n_q}, \frac{\sum_{k_q=1}^{n_q} Y_{k_q,q}}{n_q}\right) \quad (3)$$

### B. Improvements to tentative coordinate

We compute a direction vector by evaluating the concentration of center coordinates $\{C_1^{Best}, C_2^{Best}, \ldots, C_q^{Best}, \ldots, C_s^{Best}\}$ in the vicinity of $T$ to improve the position of $T$.

As can be seen from Fig. 5, we assume a square region centered at $T$ that is divided into $3 \times 3$ frames. The square region has an $A$ length along each side that initially covers all the center coordinates. We call the square region the sequence spatial density (SSD). Each frame contains a value that indicates the number of $C_q^{Best}$ coordinates.

Let $I_{x,y}$ denote the number of $C_q^{Best}$ points in the frame $(x, y)$ where $x, y$ are the indexes of the frames in SSD shown in Fig. 5. We determine a direction vector by computing the partial derivatives of SSD as the sum of the differences between two adjacent frames in SSD as:

$$\begin{aligned}
\frac{\partial I}{\partial x} &= (I_{2,1} - I_{1,1}) + (I_{3,1} - I_{2,1}) \\
&\quad + (I_{2,2} - I_{1,2}) + (I_{3,2} - I_{2,2}) \\
&\quad + (I_{2,3} - I_{1,3}) + (I_{3,3} - I_{2,3}) \\
\frac{\partial I}{\partial y} &= (I_{1,2} - I_{1,1}) + (I_{1,3} - I_{1,2}) \\
&\quad + (I_{2,2} - I_{2,1}) + (I_{2,3} - I_{2,2}) \\
&\quad + (I_{3,2} - I_{3,1}) + (I_{3,3} - I_{3,2})
\end{aligned} \quad (4)$$

We compute the partial derivatives of SSD horizontally $\frac{\partial I}{\partial x}$ and vertically $\frac{\partial I}{\partial y}$ to compute the direction vector at each $h$-th cycle of the improvements by using the direction vector function, $\overrightarrow{DVF}$ as:

$$\overrightarrow{DVF}_h = \left(\frac{\partial I}{\partial x}, \frac{\partial I}{\partial y}\right) \quad (5)$$

.

We improve $T_h$ by using $\overrightarrow{DVF}_h$ with length $\Delta T_h$ proportional to the vector's magnitude, $|\overrightarrow{DVF}_h|$, computed as:

$$|\overrightarrow{DVF}_h| = \sqrt{\left(\frac{\partial I}{\partial x}\right)^2 + \left(\frac{\partial I}{\partial y}\right)^2} \quad (6)$$

$$\Delta T_h = \left(\frac{\frac{\partial I}{\partial x}}{|\overrightarrow{DVF}_h|} \times \nu, \frac{\frac{\partial I}{\partial y}}{|\overrightarrow{DVF}_h|} \times \nu\right) \quad (7)$$

Here, $\nu$ denotes a scale factor parameter for the unit vector computed from $\overrightarrow{DVF}_h$ to determine the length of improvement, $\Delta T_h$. We improve the position of $T_h$ with direction vector $\Delta T_h$ for the next cycle of improvement as:

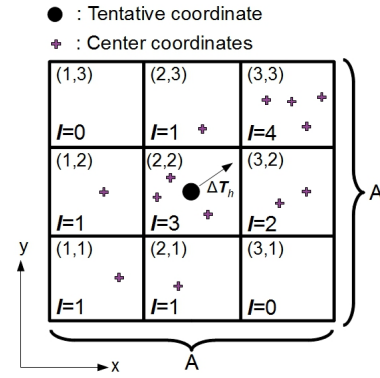$$T_{h+1} = T_h + \Delta T_h \quad (8)$$



Fig. 5: Frames of Sequence of Spatial Density (SSD)

The improvement of $T$ will also affect the concentration of center coordinates. As we can see from Fig. 4, $T$ has improved its position approaching the true positions of target nodes at 20, 40, and 60 cycles. The concentration of center coordinates has also simultaneously improved as they are computed by the average of selected footprints that are nearest to $T$ as described in Subsection V-A.

Therefore, instead of using the same values of parameters $A$ and $\nu$ in all cycles, we reduced both parameters by the fraction of $a/b$ every $m$ cycles of improvement to avoid phenomena where improvement was not taking effect because the frames were too large, as shown in Fig. 7. We called these phenomena *zero vector effects*, where the direction vector became zero as all the center coordinates were located inside the center frame of SSD. If none of the center coordinates are located in the frame other than the center frame of SSD, the direction vector will become zero as they are computed from the sum of the differences between two adjacent frames.

## VI. PERFORMANCE EVALUATION

We conducted a simulation on the proposed algorithm to evaluate the performance of MoLPS to localize sensor nodes in a noisy environment. This simulation experiment was used to demonstrate what effect a noisy environment and the shrinking size of SSD had on the localization error of sensor nodes.

The remaining part of this section presents the simulation setup, path loss model, and the results we obtained from evaluating performance.

### A. Simulation Setup

We implemented the algorithm in a custom C simulator, where we randomly deployed a set of 100 sensor nodes with one mobile receiver traveling in a 50m$\times$50m square region at a constant speed, as seen in Fig. 6. The mobile receiver and sensor nodes had the same communication range of 10m.

The mobile receiver traveled in a sensory boundary field and received signals from sensor nodes within their communication range at each time interval $t$. The positions of sensor nodes were estimated with our proposed algorithm by measuring the direction vector of the SSD from $T$. The tentative coordinate was arbitrarily deployed within the sensory

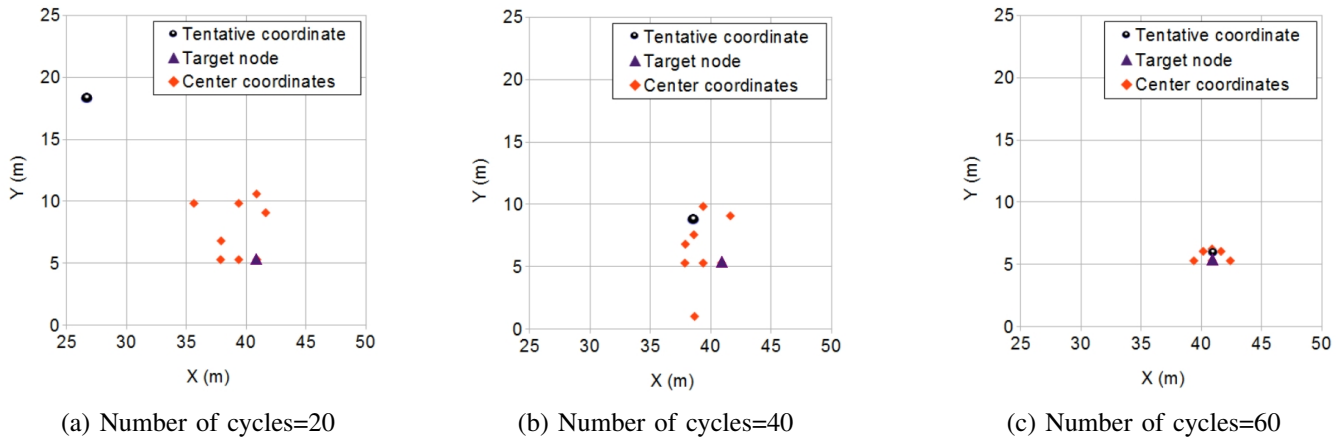(a) Number of cycles=20           (b) Number of cycles=40           (c) Number of cycles=60

Fig. 4: Improvement of $T$ approaching true position of target node

boundary field. We deployed the tentative coordinate in the first cycle of improvement at the center of the sensory boundary field as an initial coordinate of $T$. The SSD frames were used to compute the moving distance of $T$ in which the initial value of $\nu$ was 70m per number of cycles and the initial value of $A$ was 50m where all center coordinates were included in the coverage area of SSD for each tentative coordinate.

We defined the criterion in MoLPS for the error in localization as the difference between $T$ and the true position of a target node in the cycle of improvement, $h$. Localization error indicated the degree of accuracy in estimates that the algorithm could achieve.

*B. Path Loss model*

We used an extended model of log-distance path loss by combining it with the $DoI$ model [22]. The log-distance path loss model is used in many indoor and outdoor environments in which multipath propagation is presented.

The RSS reading was a value from our degree of irregularities (DoI) extended log-distance path loss in Eq. 9. There is a plot of the path loss values with our model in Fig. 8.

$$PL = \{(PL_o + 10\gamma log\frac{d}{d_o}) \times (1 \pm (rand() \times DoI))\} + S \quad (9)$$

Here, $d_o$ is the reference distance (i.e., 1 m) and $PL_o$ denotes the path loss in decibels at $d_o$, which was assumed to be 47 dB. The $d$ is the distance between sensor nodes and the mobile receiver computed from the real coordinates of the simulation system. The $\gamma$ refers to the path loss exponent, which depends on channels and the environment. According to residential indoor models [24], the path loss exponent, $\gamma$, in this model is a random variable, and requires sufficient measurements on the spot in various residential environments before effectively being applied to generic scenarios. We have used the measurements in Sohrabi et al. [25] in this paper, which denote the value of average path-loss exponents as 1.9 in an engineering building. The $S$ is log-normal shadow fading in decibels. The $S$ is usually a random variable with a Gaussian distribution with zero mean and standard deviation $\sigma$, which was assumed to be 5.7 according to Sohrabi et al. [25]. The
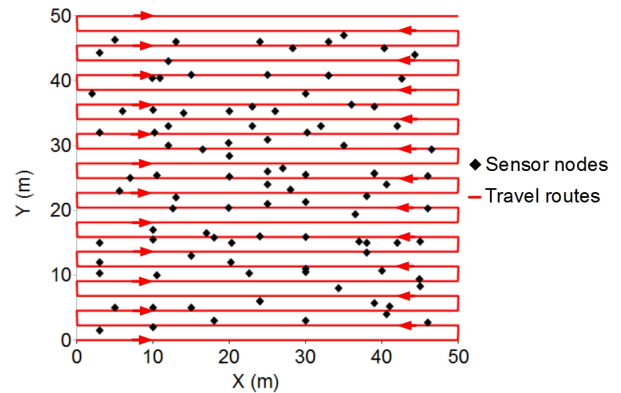


Fig. 6: Deployment of sensor nodes.

$DoI$ is the radio irregularity and $rand()$ is a random number, $\mathcal{U}(0, 1)$. We ran five simulations with different $DoI$ values in a range of $0 \sim 1.0$.

*C. Experiment Results*

We compared the cumulative distribution function (CDF) of cycles for the number of nodes that had their localization error reduced below 2m with different values of parameters $a$ and $b$ when parameters $A$ and $\nu$ were reduced by the fraction of $a/b$ through cycles of improvement. As shown in Fig. 9, less than 71% of tentative coordinates had their localization error reduced below 2m approaching the true position of target nodes when the number of cycles reached 35 where parameters $A$ and $\nu$ were reduced by the fraction of $a/b = 1/4$. However, the percentages were larger where parameters $A$ and $\nu$ were reduced by the fraction of $a/b = 3/4$, as seen in Fig. 10. In this case, the percentage of tentative coordinates that had their localization error reduced below 2m was more than 80% when the number of cycles reached 80.

The size of the square region for sequence spatial density (SSD) has an impact on computing the direction vector
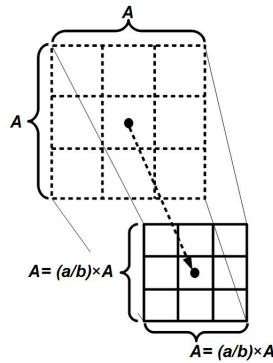
Fig. 7: Reduced $A$ length for SSD every $m$ cycles of improvement.



Fig. 8: Plot of path loss values

to improve $T$. A larger SSD will increase the number of center coordinates included in SSD that increases the total value of differences between the number of center coordinates between two adjacent frames. However, if $A$ is reduced too much between the two cycles of improvement, the number of center coordinates that are included in SSD will decrease. The decreased number of center coordinates in SSD will increase the possibility of *zero vector effects* that take place when the center coordinates are only located in the center frame of SSD. As we can see from Fig. 9, the *zero vector effects* took place when the number of cycles reached 35 and many of the center coordinates were excluded from SSD because the value of parameter $A$ was reduced too much when $a/b = 1/4$ compared to the improvement in Fig. 10 where the $T$ coordinates were continuously improved when $a/b = 3/4$ until the last number of cycles.

We also compared the required number of cycles to improve the tentative coordinates in different numbers of footprints. We fixed two values of localization error as a threshold in this evaluation scenario to assess how many cycles were needed for the tentative coordinates to improve their positions below these two thresholds (i.e., 2m and 5m). The average number of cycles for each tentative coordinate was used to represent how many cycles were required for each number of footprints. As shown in Figs. 11 and 12, the numbers of cycles were almost equal in all numbers of footprints under both conditions. However, the parameter of $a/b$ affected the number of cycles that reduced the localization error of $T$. The tentative coordinates required less than 21 cycles (2m) and 17 cycles (5m) for each threshold in which parameters $A$ and $\nu$ were reduced by the fraction of $a/b = 1/4$, as seen in Fig. 11. However, the tentative coordinates required greater numbers of cycles to reduce their localization error below 2m and 5m, as seen in Fig. 12. They needed 46 cycles for the former and 35 cycles averagely for the latter in which parameters $A$ and $\nu$ were reduced by the fraction of $a/b = 3/4$.

SSD reduced by a large fraction of $a/b$ yielded a small difference in the number of center coordinates in the frames in SSD between cycles compared to the condition in which SSD was reduced by a small fraction of $a/b$. The small fraction of $a/b$ enabled SSD to reduce its size by a larger $A$, which created large differences in the number of center coordinates in each frame as the center coordinates that were located were
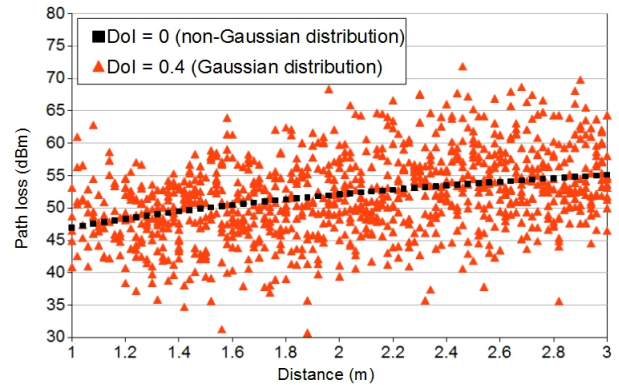
separated from one another. These will increase the value of the moving distance that improved the position of $T$ in fewer numbers of cycles.

We also compared what impact DoI had on localization error in the estimates of target node positions in various locations. The average localization error for our proposed algorithm was not entirely different for all DoI values, as seen in Fig. 13. MoLPS use range-free approaches that only use RSS to detect the proximity of footprints. Therefore, the irregularities in RSS did not have a huge impact on localization error in any target nodes on average as they did not directly use the RSS values as a metric to estimate the position of target nodes. The mean number of cycles under both conditions where the threshold of localization error was set to 2m and 5m corresponded to about 45 and 34.

## VII. CONCLUSION

We proposed range-free mobile localization based on the proximities of selected footprints in noisy environments. We used GA to iteratively search the best selection of footprints that had the nearest center coordinates to $T$. We improved the positions of tentative coordinates by measuring the direction vector from the concentration of center coordinates in the vicinity of $T$. The footprints in our proposed algorithm were divided into sets by using ranging levels to decrease variations in the radial distance between footprints in a set. We evaluated our method based on a variety of metrics that proved that it was resistant to the number of footprints used in calculations and high DoI environments at a given number of cycles while providing low localization error.

The ability to localize sensor nodes without any reference nodes in noisy environments for mobile localization can improve the localization environment in large areas. However, determining the estimates of target node positions still remains unsolved as we determined the positions of estimates by continuously improving the tentative coordinates approaching the true target nodes until the number of cycles to improve tentative coordinates satisfied a boundary (i.e., 100 cycles).

Determining suitable values for parameters $A$ and $\nu$ were major causes of difficulties in our investigations. We plan to
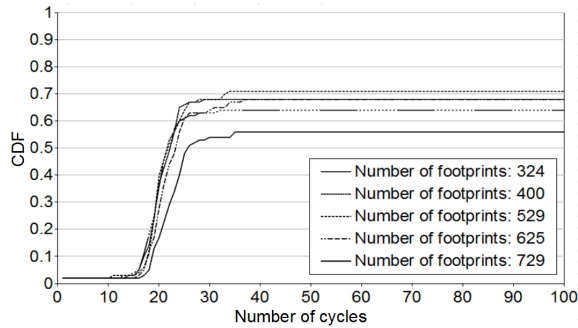
Fig. 9: Cumulative distribution of number of cycles for number of nodes that had their localization error reduced below $2m$ when parameters $A$ and $\nu$ were reduced by fraction of $a/b = 1/4$.
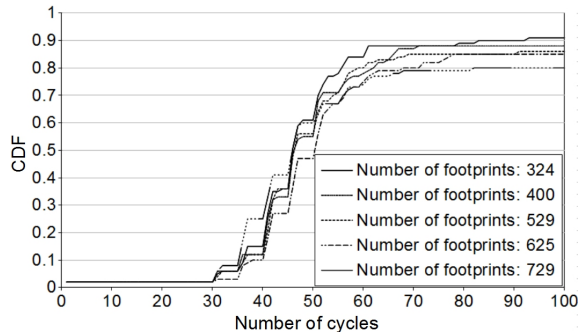


Fig. 11: Effect of localization error on number of cycles when parameters $A$ and $\nu$ are reduced by fraction of $a/b = 1/4$.



Fig. 10: Cumulative distribution of number of cycles for number of nodes that had localization error reduced below $2m$ when parameters $A$ and $\nu$ were reduced by fraction of $a/b = 3/4$.



Fig. 12: Effect of localization error on number of cycles when parameters $A$ and $\nu$ are reduced by fraction of $a/b = 3/4$.

design a method of determining receiver mobility to obtain accurate estimates by using localization based on proximity techniques. We also plan to apply our method to a real environment by running empirical experiments that focus on accurate proximity-based estimates of positions for mobile localization in the future.



Fig. 13: Effect of DoI on number of cycles.

REFERENCES

[1] F. Akyildiz, W. Su, Y. Sandarasurbramaniam, and E. Cayirci, "Wireless sensor networks: A survey," *Computer Networks Journal*, vol. 38, no. 4, pp. 393–422, 2002.

[2] M. Ang, Y. F. Lim, and M. Sim, "Robust storage assigment in unit-load warehouses," *Journal Management Science*, vol. 58, no. 11, pp. 2114–2130, 2012.

[3] J. Maneesilp, C. Wang, H. Wu, and N. F. Tzeng, "RFID Support for Accurate 3-Dimensional Localization," in *IEEE Transactions on Computers*, 2012.

[4] G. Zhou, T. He, S. Krishnamurthy, and J. A. Stankovic, "Models and solution for radio irregularity in wireless sensor networks," *ACM Trans. on Sensor Networks*, vol. 2, no. 2, pp. 221–262, 2006.

[5] K. Fukuda, and E. Okamoto, "Performance improvement of TOA Localization Using IMR-Based NLOS Detection in Sensor Networks," in *International Conference on Information Networking*, 2012, pp. 13–18.
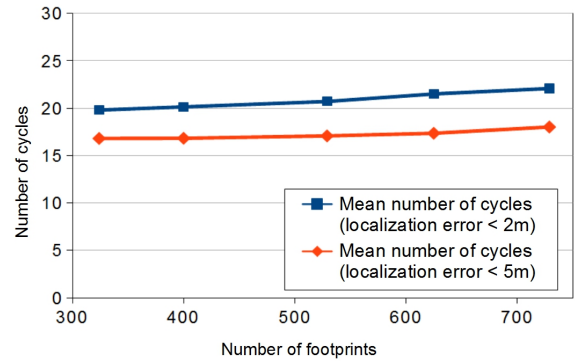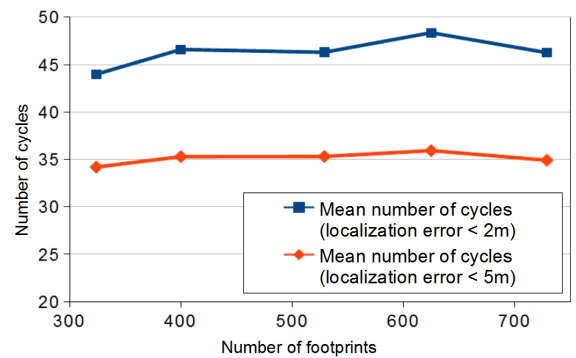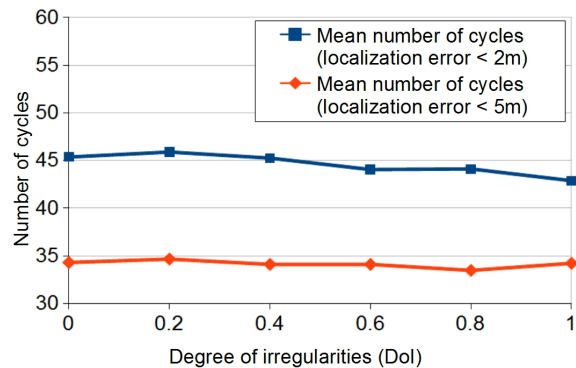
[6] B. Xu, Guodong, R. Yu, and Z. Yang, "High-accuracy tdoa-based localization without time synchronization," *IEEE Trans. on Parallel and Distribution Systems*, 2012.

[7] D. Niculescu, and B. Nath, "Ad Hoc Positioning System (APS) Using AOA," in *22nd Annual Joint Conference of the IEEE Computer and Communications. IEEE Societies*, 2003, pp. 1734–1743.

[8] K.Whitehouse, C. Karlof, and D. Culler, "A practical evaluation of radio signal strength for ranging-based localization," *ACM SIGMOBILE Mobile Computing and Communications Review*, vol. 11, no. 1, pp. 41–52, 2007.

[9]   K. Srinivasan, P. Dutta, A. Tavakoli, and P. Levis, "Understanding the causes of packet delivery success and failure in dense wireless sensor networks," in *Proc. of the 4th International Conference on Embedded Networked Sensor Systems*, 2006, pp. 419–420.

[10]  N. Bulusu, J. Heidemann, and D. Estrin, "Gps-less low cost outdoor localization for very small devices," *IEEE Personal Communication Magazine*, vol. 7, no. 5, pp. 28–34, 2000.

[11]  D. Niculescu, and B.Nath, "Dv based positioning in ad hoc networks," *Journal of Telecommunication Systems*, vol. 22, no. 1-4, pp. 267–280, 2003.

[12]  N. B. Priyantha, A. Chakraborty, and H. Balakrishnan, "The Cricket location-support system," in *Proc. of the 6th Annual International Conference on Mobile Computing and Networking*, 2000, pp. 32–43.

[13]  P. Bahl, and N. Padmanabhan, "RADAR: An In-Building RF-based User Location and Tracking System," in *19th Annual Joint Conference of the IEEE Computer and Communications. IEEE Societies*, 2000, pp. 775–784.

[14]  A. Savvides, C. C. Han and M. B. Strivastava, "Dynamic Fine-grained Localization in Ad-hoc Networks of Sensors," in *Proc. of the 7th annual international conference on Mobile computing and networking*, 2001, pp. 166–179.

[15]  P. Agrawal, and N. Patwari, "Correlated link shadow fading in multi-hop wireless networks," *IEEE Transactions on Wireless Communications*, vol. 8, no. 8, pp. 4024–4036, 2009.

[16]  J. Xiang, and Z. Hongyuan, "Sensor positioning in wireless ad-hoc sensor networks using multidimensional scaling," in *23rd Annual Joint Conference of the IEEE Computer and Communications. IEEE Societies*, 2004, pp. 2652–2661.

[17]  N. Patwari, and A. O. Hero, "Manifold learning algorithms for localization in wireless sensor networks," in *IEEE International Conference on Acoustics, Speech, and Signal Processing, 2004. Proceedings. (ICASSP '04)*, 2004, pp. 857–860.

[18]  W. Chengqun, C. Jiming, S. Youxian, and S. Xuemin, "Wireless Sensor Networks Localization with Isomap," in *IEEE International Conference on Communications, 2009. ICC '09*, 2009, pp. 1–5.

[19]  M. R. Basheer, and S. Jagannathan, "Localization of objects using stochastic tunneling," in *IEEE Wireless Communications and Networking Conference (WCNC)*, 2011, pp. 587–592.

[20]  Y. Shang, W. Ruml, Y. Zhang, and M. P. J. Fromherz, "Localization from mere connectivity," in *Proc. of the 4th ACM international symposium on Mobile ad hoc networking and computing*, 2003, pp. 201–212.

[21]  S. Lederer, Y. Wang. and J. Gao, "Connectivity-based localization of large-scale sensor networks with complex shape," *ACM Transactions on Sensor Networks (TOSN)*, vol. 5, no. 4, 2009.

[22]  T. He, C. Huang, B. M. Blum, J. A. Stankovic, and T. Abdelzaher, "Range-Free Localization Schemes for Large Scale Sensor Networks," in *Proc. of MobiCom*, 2003, pp. 81–95.

[23]  M. Li, and Y. Liu., "Rendered Path: Range-Free Localization in Anistropic Sensor Networks With Holes," in *Proc. of MobiCom*, 2007, pp. 51–62.

[24]  S. S. Ghassemzadeh, R. Jana, C. W. Rice, W. Turin, V. Tarokh, "Measurement and modeling of an ultra-wide bandwidth indoor channel," *Transactions of Communication*, vol. 52, no. 10, pp. 1786–1796, 2004.

[25]  K. Sohrabi, B. Manriquez, and G. Z. Pottie, "Near Ground Wideband Channel Measurement in 800-1000MHz," in *Proc. of IEEE Vehicular Technology Conference*, 1999, pp. 1222–1226.