# IJACSA

WHERE WISDOM SHARES

## 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 January 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

Rajalakshmi Engineering College; Matrix Vision GmbH

- **Bilian Song**
  LinkedIn

- **Brahim Raouyane**
  INPT

- **Brij Gupta**
  University of New Brunswick

- **Constantin Filote**
  Stefan cel Mare University of Suceava

- **Constantin Popescu**
  Department of Mathematics and Computer Science, University of Oradea

- **Chandrashekhar Meshram**
  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

- **Prof. D. S. R. Murthy**
  Sreeneedhi

- **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

- **Dr. Gunaseelan Devraj**
  Jazan University, Kingdom of Saudi Arabia

- **Dr. Bright Keswani**
  Associate Professor and Head, Department of Computer Applications, Suresh Gyan Vihar University, Jaipur (Rajasthan) INDIA

- **Dr VAKA MOHAN**
  TRR COLLEGE OF ENGINEERING

- **Dr. Faris Al-Salem**
  GCET

- **Dr. S Kumar**
  Anna University

- **Dr. Sanskruti Patel**
  Charotar Univeristy of Science & Technology, Changa, Gujarat, India

- **DR.RAVISANKAR HARI**
  SENIOR SCIENTIST, CTRI, RAJAHMUNDRY

- **Dragana Becejski-Vujaklija**
  University of Belgrade, Faculty of organizational sciences

- **Driss EL OUADGHIRI**

- **Dr.VUDA SREENIVASARAO**
  School of Computing and Electrical Engineering,BAHIR DAR UNIVERSITY, BAHIR DAR,ETHIOPA.

- **Duck Hee Lee**
  Medical Engineering R&D Center/Asan Institute for Life Sciences/Asan Medical Center

- **Driss EL OUADGHIRI**

- **Dr. Omaima Al-Allaf**
  Asesstant Professor

- **Elena Camossi**
  Joint Research Centre

- **Eui Lee**

- **Firkhan Ali Hamid Ali**
  UTHM

- **Fokrul Alom Mazarbhuiya**
  King Khalid University

- **Frank Ibikunle**
  Covenant University

- **Fu-Chien Kao**
  Da-Y eh University

- **G. Sreedhar**
  Rashtriya Sanskrit University

- **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)

- **Gufran Ahmad Ansari**
  Qassim University

- **Giri Babu**
  Indian Space Research Organisation

- **Giacomo Veneri**
  University of Siena

- **Gerard Dumancas**
  Oklahoma Medical Research Foundation
- **Georgios Galatas**
- **George Mastorakis**
  Technological Educational Institute of Crete
- **Gavril Grebenisan**
  University of Oradea
- **Hadj Tadjine**
  IAV GmbH
- **Hamid Mukhtar**
  National University of Sciences and Technology
- **Hanumanthappa J**
  UNIVERSITY OF MYSORE
- **Hamid Alinejad-Rokny**
  University of Newcastle
- **Harco Leslie Hendric Spits Warnars**
  Budi LUhur University
- **Harish Garg**
  Thapar University Patiala
- **Hardeep**
  Ferozaepur College of Engineering & Technology, India
- **Hamez I. 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
- **Jasvir Singh**
  Communication Signal Processing Research Lab
- **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
- **Kanak Saxena**
  S.A.TECHNOLOGICAL INSTITUTE
- **Ka Lok Man**
  Xi'an Jiaotong-Liverpool University (XJTLU)
- **Kushal Doshi**
  IEEE Gujarat Section
- **Kashif Nisar**
  Universiti Utara Malaysia
- **Kavya Naveen**
- **Kayhan Zrar Ghafoor**
  University Technology Malaysia
- **Kitimaporn Choochote**
  Prince of Songkla University, Phuket Campus
- **Kohei Arai**
  Saga University
- **Kunal Patel**
  Ingenuity Systems, USA
- **Krasimir Yordzhev**
  South-West University, Faculty of Mathematics and Natural Sciences, Blagoevgrad, Bulgaria
- **Labib Francis Gergis**
  Misr Academy for Engineering and Technology
- **Lai Khin Wee**
  Biomedical Engineering Department, University Malaya
- **Latha Parthiban**
  SSN College of Engineering, Kalavakkam
- **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
- **Mohamed Najeh Lakhoua**
  ESTI, University of Carthage
- **Mohammad Alomari**
  Applied Science University
- **Mohammad Kaiser**
  Institute of Information Technology
- **Mohammed Al-Shabi**
  Assisstant 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
- **Monji Kherallah**
  University of Sfax
- **Mostafa Ezziyyani**
  FSTT
- **Mueen Uddin**
  Universiti Teknologi Malaysia UTM
- **Mona Elshinawy**
  Howard University
- **N Ch.Sriman Narayana Iyengar**
  VIT University
- **Natarajan Subramanyam**
  PES Institute of Technology
- **Neeraj Bhargava**
  MDS University
- **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
- **Na Na**
  NA
- **Om Sangwan**
- **Oliviu Matel**
  Technical University of Cluj-Napoca
- **Osama Omer**
  Aswan University
- **Ousmane Thiare**
  Associate Professor University Gaston Berger of Saint-Louis SENEGAL
- **Pankaj Gupta**
  Microsoft Corporation
- **Paresh V Virparia**
  Sardar Patel University
- **Dr. Poonam Garg**
  Institute of Management Technology, Ghaziabad
- **Prabhat K Mahanti**
  UNIVERSITY OF NEW BRUNSWICK
- **Qufeng Qiao**
  University of Virginia
- **Rachid Saadane**
  EE departement EHTP

- **raed Kanaan**
  Amman Arab University
- **Raghuraj Singh**
- **Raj Gaurang Tiwari**
  AZAD Institute of Engineering and Technology
- **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
- **Samarjeet Borah**
  Dept. of CSE, Sikkim Manipal University
  University College of Applied Sciences UCAS-Palestine
- **Santosh Kumar**
  Graphic Era University, India
- **Sasan Adibi**
  Research In Motion (RIM)
- **Saurabh Pal**
  VBS Purvanchal University, Jaunpur
- **Saurabh Dutta**
  Dr. B. C. Roy Engineering College, Durgapur
- **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
- **Shikha Bagui**
  University of West Florida
- **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

(vii)

# CONTENTS

# A Comparative Study of Meta-heuristic Algorithms for Solving Quadratic Assignment Problem

Gamal Abd El-Nasser A. Said
Computer Science Department
Faculty of Computer & Information
Sciences, Ain Shams University
Cairo, Egypt

Abeer M. Mahmoud
Computer Science Department
Faculty of Computer & Information
Sciences, Ain Shams University
Cairo, Egypt

El-Sayed M. El-Horbaty
Computer Science Department
Faculty of Computer & Information
Sciences, Ain Shams University
Cairo, Egypt

*Abstract*—**Quadratic Assignment Problem (QAP) is an NP-hard combinatorial optimization problem, therefore, solving the QAP requires applying one or more of the meta-heuristic algorithms. This paper presents a comparative study between Meta-heuristic algorithms: Genetic Algorithm, Tabu Search, and Simulated annealing for solving a real-life (QAP) and analyze their performance in terms of both runtime efficiency and solution quality. The results show that Genetic Algorithm has a better solution quality while Tabu Search has a faster execution time in comparison with other Meta-heuristic algorithms for solving QAP.**

*Keywords*—*Quadratic Assignment Problem (QAP); Genetic Algorithm (GA); Tabu Search (TS); Simulated Annealing (SA); Performance Analysis*

## I. INTRODUCTION

Optimization problems arise in various disciplines such as engineering design, manufacturing system, economics etc. thus in view of the practical utility of optimization problems there is a need for efficient and robust computational algorithms which can solve optimization problems arising in different fields. Several NP-hard combinatorial optimization problems, such as the traveling salesman problem, and yard management of container terminals can be modeled as QAPs..

Optimization is a process that finds a best, or optimal, solution for a problem. An optimization problem is defined as: Finding values of the variables that minimize or maximize the objective function while satisfying the constraints. The Optimization problems are centered on three factors: (1) an objective function which is to be minimized or maximized. (2) A set of unknowns or variables that affect the objective function. (3) A set of constraints that allow the unknowns to take on certain values but exclude others.

In most optimization problems there is more than one local solution. Therefore, it becomes very important to choose a good optimization method that will not be greedy and look only in the neighborhood of the best solution; because this will mislead the search process and leave it stuck at a local solution. However, the optimization algorithm should have a mechanism to balance between local and global search. There are multiple methods used to solve optimization problems of both the mathematical and combinatorial types. In fact, if the optimization problem is difficult or if the search space is large, it will become difficult to solve the optimization problem by using conventional mathematics.

For this reason, many meta-heuristic optimization methods have been developed to solve such difficult optimization problems [13].

Combinatorial generally means that the state space is discrete. Combinatorial optimization is widely applied in a number of areas nowadays. Combinatorial optimization problems (COP) are those problems that have a finite set of possible solutions. The best way to solve a combinatorial optimization problem is to check all the feasible solutions in the search space. However, checking all the feasible solutions is not always possible, especially when the search space is large. Thus, many Meta-heuristic algorithms have been devised and modified to solve these problems. The Meta-heuristic approaches are not guaranteed to find the optimal solution since they evaluate only a subset of the feasible solutions, but they try to explore different areas in the search space in a smart way to get a near-optimal solution in less cost and time [11].

In this paper we focus on combinatorial optimization problem, namely the Quadratic Assignment Problem. (QAP) is one of the most difficult NP-hard combinatorial optimization problems, so, to practically solve the QAP one has to apply Meta-heuristic algorithms which find very high quality solutions in short computation time[15].

The rest of this paper is organized as follows: A brief description of QAP is given in section II. Section III provides a brief overview of related work for comparison between different Meta-heuristic algorithms for solving combinatorial problems. The Meta-heuristic algorithms are described in section IV. This section is further subdivided into three subsections namely GA, TS, and SA. Section V, include the experimental results. Our conclusions and future work are given in section VI.

## II. QUADRATIC ASSIGNMENT PROBLEM

QAP is one of the most difficult NP-hard combinatorial optimization problems; there are a set of n facilities and a set of n locations. For each pair of locations, a distance is specified and for each pair of facilities a weight or flow is specified. The problem is to assign all facilities to different locations with the aim of minimizing the sum of the distances multiplied by the corresponding flows. (QAP) is formulated as follows:-

The following notation is used in formulation of QAP

n     total number of facilities and locations

$f_{ik}$    flow of material from facility I to facility k

$d_{jl}$    distance from location j to location l

The objective function minimizes the total distances and flows between facilities

$$\min \quad f(x) = \sum_{i=1}^{n}\sum_{j=1}^{n}\sum_{k=1}^{n}\sum_{l=1}^{n} f_{ik}\, d_{jl}\, x_{ij}\, x_{kl}$$

$$\text{s.t} \quad \sum_{j=1}^{n} x_{ij} = 1,$$

$$\sum_{i=1}^{n} x_{ij} = 1,$$

where

$$x_{ij} = \begin{cases} 1, & \text{if facility i is assigned to location j} \\ 0, & \text{otherwise} \end{cases}$$

The constraints ensure that each facility i is assigned to exactly one location j and each location j has exactly one facility which assigned to it [16].

### III. RELATED WORK

Many Meta-heuristic algorithms have been proposed by researchers to find optimal or near optimal solutions for the QAP such as Genetic Algorithm [1], Tabu Search [3] and Simulated Annealing [15]. Also Many researchers presented comparison study between different Meta-heuristic algorithms for solving combinatorial problems [2,5,11].

John Silberholz and Bruce Golden [2], compared Meta-heuristic algorithms in terms of both solution quality and runtime, Their conclusions show that good techniques in solution quality and runtime comparisons will ensure fair and meaningful comparisons are carried out between Meta-heuristic algorithms, producing the most meaningful and unbiased results possible.

Comparison between simulated annealing and genetic algorithm for solving the Travelling Salesman Problem was done by Adewole et al. [4], where they have compared the performance of SA and GA. Their results show that Simulated Annealing runs faster than Genetic Algorithm and runtime of Genetic Algorithm increases exponentially with number of cities. However, in terms of solution quality Genetic Algorithm is better than Simulated Annealing.

Bajeh et al. [5] compared Genetic Algorithm and Tabu Search approaches to solve scheduling problems. The results show that TS can produce better solution, with less computing time, than those produced by GA. However, GA can produce several different near optimal solutions at the same time because of its holds the whole generation of chromosomes which may not originate from the same parents.

Marvin et al. [6] compared the relative performance of Tabu Search (TS), Simulated Annealing (SA) and Genetic Algorithms (GA) on various types of FLP under time-limited, solution-limited, and unrestricted conditions. The results indicate that TS shows very good performance in most cases. The performance of SA and GA are more partial to problem type and the criterion used.

In Karimi et al. [7] Meta-heuristic methods such as SA, TS, and PSO are presented; the research is dedicated to compare the relative percentage deviation of these solution qualities from the best known solution which is introduced in QAPLIB. The results show that TS is the most excellent method in computational time.

Paul [8] compared the performance of tabu search and simulated annealing heuristics for the quadratic assignment problem. The results shows that for a number of varied problem instances, SA performs better for higher quality targets while TS performs better for lower quality targets.

This paper presents Genetic algorithm (GA), Tabu search (TS) and simulated annealing (SA) for solving real life Quadratic Assignment Problem. The analysis of the obtained results in terms of both runtime efficiency and solution quality show the performance of each algorithm and show comparison on their effectiveness in finding the optimal solution for real life QAP.

### IV. META-HEURISTIC ALGORITHMS

A Meta-heuristic is formally defined as an iterative generation process which guides a subordinate heuristic by combining intelligently different concepts for exploring and exploiting the search space, learning strategies are used to structure information in order to find efficiently near-optimal solutions. Meta-heuristic algorithms are among these approximate techniques which can be used to solve complex problems.

Most widely known Meta-heuristic algorithms are Genetic algorithm (GA), simulated annealing (SA) and Tabu search (TS). Genetic algorithm (GA) emulate the evolutionary process in nature, whereas tabu search (TS) exploits the memory structure in living beings, simulated annealing (SA) imitates the annealing process in crystalline solids [2].

#### A. Genetic algorithm

Genetic Algorithm is a Meta-heuristic algorithm that aims to find solutions to NP-hard problems. The basic idea of Genetic Algorithms is to first generate an initial population randomly which consist of individual solution to the problem called Chromosomes, and then evolve this population after a number of iterations called Generations. During each generation, each chromosome is evaluated, using some measure of fitness. To create the next generation, new chromosomes, called offspring, are formed by either merging two chromosomes from current generation using a crossover operator or modifying a chromosome using a mutation operator. A new generation is formed by selection, according to the fitness values, some of the parents and offspring, and rejecting others so as to keep the population size constant. Fitter chromosomes have higher probabilities of being selected. After several generations, the algorithms converge to the best chromosome, which hopefully represents the optimum or suboptimal solution to the problem [12]. The process of GA can be represented as follows:

*Step 1 Generate initial population.*

*Step 2 Evaluate populations.*

*Step 3 Apply Crossover to create offspring.*

*Step 4 Apply Mutation to offspring.*

*Step5 Select parents and offspring to form the new population for the next generation.*

*Step 6 If termination condition is met finish, otherwise go to Step2*

## B. Tabu Search

Tabu search is the technique that keeps track of the regions of the solution space that have already been searched in order to avoid repeating the search near these areas [8]. It starts from a random initial solution and successively moves to one of the neighbors of the current solution. The difference of tabu search from other Meta-heuristic approaches is based on the notion of tabu list, which is a special short term memory. That is composed of previously visited solutions that include prohibited moves. In fact, short term memory stores only some of the attributes of solutions instead of whole solution. So it gives no permission to revisited solutions and then avoids cycling and being stuck in local optima.

During the local search only those moves that are not tabu will be examined if the tabu move does not satisfy the predefined aspiration criteria. These aspiration criteria are used because the attributes in the tabu list may also be shared by unvisited good quality solutions. A common aspiration criterion is better fitness, i.e. the tabu status of a move in the tabu list is overridden if the move produces a better solution [2, 3]. The process of TS can be represented as follows:

*Step 1 Generate initial solution x.*

*Step 2 Initialize the Tabu List.*

*Step 3 While set of candidate solutions X″ is not complete.*

*Step 3.1 Generate candidate solution x″ from current solution x*

*Step 3.2 Add x″ to X″ only if x″ is not tabu or if at least one Aspiration Criterion is satisfied.*

*Step 4 Select the best candidate solution x* in X″.*

*Step 5 If fitness(x*) > fitness(x) then x = x*.*

*Step 6 Update Tabu List and Aspiration Criteria*

*Step 7 If termination condition met finish, otherwise go to Step 3.*

## C. Simulated annealing

Simulated Annealing is an early Meta-heuristic algorithm originating from an analogy of how an optimal atom configuration is found in statistical mechanics. It uses temperature as an explicit strategy to guide the search. In Simulated Annealing, the solution space is usually explored by taking random tries. The Simulated Annealing procedure randomly generates a large number of possible solutions, keeping both good and bad solutions.

As the simulation progresses, the requirements for replacing an existing solution or staying in the pool becomes stricter and stricter, mimicking the slow cooling of metallic annealing. Eventually, the process yields a small set of optimal solutions. Simulated Annealing advantage over other methods is its ability to obviate being trapped in local minima.

This means that the algorithm does not always reject changes that decrease the objective function or changes that increase the objective function according to its probability function: $p = e^{\Delta f/T}$ Where T is the control parameter (analogy to temperature) and $\Delta f$ is the variation in the objective function [4,15]. The process of SA can be represented as follows:

*Step 1 Compute randomly next position .*

*Step2 Determine the difference between the next position and current position, call this different delta .*

*Step3 If delta < 0, the assign the next position to the current position .*

*Step4 If delta > 0, then compute the probability of accepting the random next position .*

*Step5 If the probability is < the e^(-delta / temperature), then assign the next position to the current position .*

*Step 6 Decrease temperature by a factor of alpha .*

*Step7 Loop to step 1 until temperature is not greater than epsilon*

## V. EXPERIMENTS AND RESULTS

Experimental results were run on a Laptop with the following configurations: i3 CPU 2.4 GHZ, 4.0 GB RAM, Windows 7. This test was conducted with GA, TS and SA algorithms. Comparison of the algorithms is based on solution quality and execution time for real life QAP, in the experiment, we analyze the solution quality and run time for solving QAP using instances presented in QAPLIB site [18].

QAPLIB problems are classified to four classes. (i) Unstructured, randomly generated instances. (ii) Grid- based distance matrix instances (iii) Real-life instances (iv) Real-life like instances.

These groups are taken from the study of Ramkumar et al. (2009) [17]. Our experiment for Real-life instances class which have Best Known Quality Solution as shown in the Table1.

TABLE I.  QAP AT DIFFERENT PROBLEMS SIZE

| Problem Name | Problem Size | Best Known Quality |
|---|---|---|
| bur26h | 26 | 7098658 |
| chr12c | 12 | 11156 |
| Chr15a | 15 | 9896 |
| Esc128 | 128 | 64 |
| esc16i | 16 | 14 |
| esc32h | 32 | 438 |
| esc64a | 64 | 128 |
| had12 | 12 | 1652 |
| had14 | 14 | 2724 |
| had20 | 20 | 6922 |
| kra30b | 30 | 91420 |
| ste36a | 36 | 9526 |

The obtained Best quality solutions for each algorithm of Meta-heuristic algorithms are compared with QAPLIB Best Known Quality solutions. For each problem instance we execute a series of runs for various parameters.

Figure.1 shows an example of solving QAP namely ste36a. The figure shows the best Quality solution and the best known Quality solution of the problem ste36a by tabu search algorithm.



Fig. 1.  solution quality of instance ste36a using tabu search algorithm

The results shown in table 2 show the relative differences of the solution quality for real life quadratic assignment problems by each algorithm for different problems size.

Relative difference represents the difference between algorithm best quality solution and the best known quality solution of the problem in percent. The difference value is calculated in the following way

Relative difference= ((Best Quality − Best Known Quality) / Best Known Quality) * 100%

TABLE II.  RELATIVE DIFFERENCE OF THE SOLUTION QUALITY FOR QAP AT DIFFERENT PROBLEMS SIZE

| Problem Nam | GA Diff% | TS Diff% | SA Diff% |
|---|---|---|---|
| bur26h | 0.84% | 1.51% | 0.96% |
| chr12c | 9% | 25.91% | 16.22% |
| Chr15a | 27.21% | 77.66% | 32.95% |
| Esc128 | 91.67% | 73.26% | 60.94% |
| esc16i | 0% | 3% | 0% |
| esc32h | 7.23% | 10.13% | 3.88% |
| esc64a | 1% | 1.07% | 0% |
| had12 | 0.61% | 1.70% | 0.87% |
| had14 | 0% | 1.31% | 0% |
| had20 | 0.93% | 2.23% | 0.82% |
| kra30b | 6.09% | 14.38% | 14.44% |
| ste36a | 36.70% | 31.81% | 38.52% |

As for the execution time, in table 3 we compare the execution time by each algorithm for QAP at different problems size. Execution time in the format (minutes: seconds. tenths of seconds).

TABLE III.  SOLUTION EXECUTION TIME FOR QAP AT DIFFERENT PROBLEMS SIZE

| Problem Name | Execution Time | Execution Time | Execution Time |
|---|---|---|---|
| | GA | TS | SA |
| bur26h | 00:10.4 | 00:00.4 | 00:02.7 |
| chr12c | 00:08.9 | 00:00.1 | 00:02.3 |
| Chr15a | 00:08.6 | 00:00.2 | 00:02.5 |
| Esc128 | 00:33.6 | 00:17.8 | 00:13.7 |
| esc16i | 00:10.0 | 00:00.1 | 00:03.6 |
| esc32h | 00:10.2 | 00:00.4 | 00:03.3 |
| esc64a | 00:15.4 | 00:01.6 | 00:04.9 |
| had12 | 00:08.6 | 00:00.2 | 00:02.4 |
| had14 | 00:08.7 | 00:00.2 | 00:02.4 |
| had20 | 00:09.0 | 00:00.3 | 00:02.5 |
| kra30b | 00:10.2 | 00:00.5 | 00:02.8 |
| ste36a | 00:10.8 | 00:00.6 | 00:03.0 |

Figure.2 shows the relative percentage deviation (Relative difference) of the solution quality for different problems size for GA,TS and GA algorithms, the results shows that genetic algorithm has a good solution quality more than the other Meta-heuristic algorithms for solving QAP instances.

Fig. 2. Relative percentage deviation of the solution quality for different problems size

Figure.3 shows the execution time for QAP for different problems size for GA,TS and GA algorithms, the results shows that Tabu search algorithm has a faster execution time than the other Meta-heuristic algorithms for solving Real-life QAP instances.



Fig. 3. Execution time for QAP for different problems size

## VI. CONCLUSION AND FUTURE WORK

In this paper, we applied Genetic algorithm (GA), tabu search (TS), and simulated annealing (SA) as Meta-heuristic algorithms for solving the Real life QAP. This research is dedicated to compare the relative percentage deviation of these solution qualities from the best known quality solution which is introduced in QAPLIB. The results show that GA, TS, and SA algorithms have effectively demonstrated the ability to solve QAP optimization problems. the computational results show that genetic algorithm has a better solution quality than the other Meta-heuristic algorithms for solving QAP problems. Tabu search algorithm has a faster execution time than the other Meta-heuristic algorithms for solving Real-life QAP problems.

In future research, comparisons between Meta-heuristic algorithms for more different types, different sizes of QAP instances and different algorithms can be conducted. Also apply Meta-heuristic algorithms to solve other combinatorial problems such as container terminals problems.

REFERENCES

[1] Yongzhong Wu, and Ping Ji, "Solving the Quadratic Assignment Problems by a Genetic Algorithm with a New Replacement Strategy". International Journal of Computational Intelligence Volume 4 Number 3, 2007.

[2] John Silberholz and Bruce Golden, "Comparison of Meta-heuristic " Handbook of Meta-heuristic algorithms International Series in perations Research & Management Science Volume 146, pp 625-640,2010.

[3] Elena Ikonomovska, Ivan Chorbev, Dejan Gjorgjevik and Dragan Mihajlov "The Adaptive Tabu Search and Its Application to the Quadratic Assignment Problem", Proceedings of 9th International Multiconference - Information Society 2006, pp. 26-29, Ljubljana, Slovenia, 2006.

[4] Adewole A.P., Otubamowo K. Egunjobi T.O. and Kien Ming Ng, "A Comparative Study of Simulated Annealing and Genetic Algorithm for Solving the Travelling Salesman Problem ", International Journal of Applied Information Systems (IJAIS), Volume 4– No.4, October 2012.

[5] Bajeh, A. O. and Abolarinwa, K. O. , " Optimization: A Comparative Study of Genetic and Tabu Search Algorithms", International Journal of Computer Applications (IJCA), Volume 31– No.5, October 2011.

[6] Marvin A. Arostegui Jr., Sukran N. Kadipasaoglu, and Basheer M. Khumawala, "An empirical comparison of Tabu Search, Simulated Annealing, and Genetic Algorithms for facilities location problems", International. Journal of Production Economics 103 (2006) 742–754, 2006.

[7] Mahdi Bashiri and Hossein Karimi, "Effective heuristics and meta-heuristics for the quadratic assignment problem with tuned parameters and analytical comparisons", Journal of Industrial Engineering International, 2012.

[8] Gerald Paul, "Comparative performance of tabu search and simulated annealing heuristics for the quadratic assignment problem ", Operations Research Letters 38 (2010) 577–581, 2010.

[9] Houck C.R, Joines J.A, Kay M.G, "Characterizing Search Spaces For Tabu Search", Currently under second review in European Journal of Operational Research., 2011.

[10] Clayton W. Commander, "A Survey of the Quadratic Assignment Problem, with Applications", Morehead Electronic Journal of Applicable Mathematics. Issue 4 | MATH-2005-01, 2005.

[11] Malti Baghel, Shikha Agrawa, and Sanjay Silakari Ph.D, "Survey of Meta-heuristic Algorithms for Combinatorial Optimization", International Journal of Computer Applications (0975 – 8887) Volume 58– No.19, November 2012.

[12] Frank Neumann, Carsten Witt, "Bio inspired Computation in Combinatorial Optimization", Springer Heidelberg Dordrecht London New York, 2010.

[13] Azmi Alazzam and Harold W. Lewis, "A New Optimization Algorithm For Combinatorial Problems", (IJARAI) International Journal of Advanced Research in Artificial Intelligence, Vol. 2, No.5, 2013.

[14] Mohamed Saifullah Hussin and Thomas Stutzle, "Tabu Search vs. Simulated Annealing for Solving Large Quadratic Assignment Instances", IRIDIA-Technical Report Series Technical Report No.TR/IRIDIA/2010-020 October 2010.

[15] Ghandeshtani, Mollai, Seyedkashi, and Neshati, "New Simulated Annealing Algorithm for Quadratic Assignment Problem", The Fourth International Conference on Advanced Engineering Computing and Applications in Sciences, 2010.

[16] Hossein Shahbazi, Ali Eghbali Ghahyazi, and Farhad Zeinali, "Quadratic Assignment Problem", Journal of American Science 2013;9(8), 2013.

[17] A.S. Ramkumar, S.G. Ponnambalam, N. Jawahar, "A new iterated fast local search heuristic for solving QAP formulation in facility layout design Robotics and Computer-Integrated Manufacturing 25(2009) 620–629, 2009.

[18] R. E. Bedkard, S. E. Karisch, and F. Rendl, "QAPLIB-A Quadratic assignment problem library", http://www.opt.math.tu-graz.ac.at/qaplib.

AUTHOR'S PROFILE

**Gamal Abd El-Nasser A. Said**: He received his M.Sc. (2012) ) in computer science from College of Computing & Information Technology, Arab Academy for Science and Technology and Maritime Transport (AASTMT), Egypt and B.Sc (1990) Faculty of Electronic Engineering, Menofia University, Egypt.

His work experience as a Researcher, Maritime Researches & Consultancies Center, Egypt. Computer Teacher, College of Technology Kingdom Of Saudi Arabia and Lecturer, Port Training Institute, (AASTMT), Egypt. Now he is Ph.D. student in computer science, Ain Shams University. His research areas include optimization, discrete-event simulation, and artificial intelligence.

**Dr Abeer M. Mahmoud**: She received her Ph.D. (2010) in Computer science from Niigata University, Japan, her M.Sc (2004) B.Sc. (2000) in computer science from Ain Shams University, Egypt.

Her work experience is as a lecturer assistant and assistant professor, faculty, of computer and information sciences, Ain. Shams University. Her research areas include artificial intelligence medical data mining, machine learning, and robotic simulation systems.

**Professor El-Sayed M. El-Horbaty**: He received his Ph.D. in Computer science from London University, U.K., his M.Sc. (1978) and B.Sc (1974) in Mathematics From Ain Shams University, Egypt. His work experience includes 39 years as an in Egypt (Ain Shams University), Qatar(Qatar University) and Emirates (Emirates University, Ajman University and ADU University). He Worked as Deputy Dean of the faculty of IT, Ajman University (2002-2008). He is working as a Vice Dean of the faculty of Computer & Information Sciences, Ain Shams University (2010-Now). Prof. El-Horbaty is current areas of research are parallel algorithms, combinatorial optimization, image processing. His work appeared in journals such as Parallel Computing, International journal of Computers and Applications (IJCA), Applied Mathematics and Computation, and International Review on Computers and software. Also he has been involved in more than 26 conferences.

# Generic Packing Detection using Several Complexity Analysis for Accurate Malware Detection

Dr. Mafaz Mohsin Khalil Al-Anezi

Computer Sciences
College of Computer Sciences and Mathematics,
Mosul University, Mosul, Iraq

*Abstract—* **The attackers do not want their Malicious software (or malwares) to be reviled by anti-virus analyzer. In order to conceal their malware, malware programmers are getting utilize the anti reverse engineering techniques and code changing techniques such as the packing, encoding and encryption techniques. Malware writers have learned that signature based detectors can be easily evaded by "packing" the malicious payload in layers of compression or encryption. State-of-the-art malware detectors have adopted both static and dynamic techniques to recover the payload of packed malware, but unfortunately such techniques are highly ineffective. If the malware is packed or encrypted, then it is very difficult to analyze. Therefore, to prevent the harmful effects of malware and to generate signatures for malware detection, the packed and encrypted executable codes must initially be unpacked. The first step of unpacking is to detect the packed executable files.**

**The objective is to efficiently and accurately distinguish between packed and non-packed executables, so that only executables detected as packed will be sent to an general unpacker, thus saving a significant amount of processing time. The generic method of this paper show that it achieves very high detection accuracy of packed executables with a low average processing time.**

**In this paper, a packed file detection technique based on complexity measured by several algorithms, and it has tested using a packed and unpacked dataset of file type .exe. The preliminary results are very promising where achieved high accuracy with enough performance. Where it achieved about 96% detection rate on packed files and 93% detection rate on unpacked files. The experiments also demonstrate that this generic technique can effectively prepared to detect unknown, obfuscated malware and cannot be evaded by known evade techniques.**

*Keywords—Packed Executables; Malware Detection; compression algorithms*

## I. INTRODUCTION

As a consequence of the arms race between virus writers and anti-virus vendors, sophisticated code obfuscation techniques are commonly implemented in computer viruses. Executable code polymorphism, metamorphism, packing, and encryption, have been proven very effective in evading detection by traditional signature-based anti-virus software. Traditional signature-based anti-virus software needs updating the virus database regularly, and the virus detection relying on the known virus database is a passive protection technology

without the capacity of detecting the new unknown virus, the virus deformation, and packed virus. Among these techniques, executable packing is the most common due to the availability of several open source and commercial executable packers [21][14].

According to [9][5][7], over 80% of computer viruses appear to be using packing techniques. Moreover, there is evidence that more than 50% of new viruses are simply re-packed versions of existing ones, see Fig.1. It has been reported that among 20, 000 malware samples collected in April 2008, more than 80% were packed by packers from 150 different families. This is further complicated by the ease of obtaining and modifying the source code of various packers. Currently, new packers are created from existing ones at a rate of 10 to 15 per month [7].

Although executable packing is very popular among virus writers, it is also applied for encrypting benign executables. Programmers of benign software apply packing to their applications mainly to make the resulting executables smaller in terms of bytes, and therefore faster to distribute through the network, for example. Also, packing makes reverse-engineering more difficult, thus making it harder for hackers to break the software license protections. As a matter of fact, there exist many commercial executable packing tools that have been developed mainly for protecting benign applications from software piracy. However, the percentage of packed benign executables is low (perhaps as low as 1%, although we were not able to find any study that can confirm this estimate, which is based solely on our experience) [14].



Fig. 1. Malware and packing, 80% of new malware are packed with various packers, 50% of new malware samples are simply repacked versions of existing malware [8].

An executable packing tool is a software that given a program P generates a new program P′ which embeds an encrypted version of P and a decryption routine. When P′ is executed, it will decrypt P on the fly and then run it. Assuming P contains known malicious code, signature based anti-virus would (likely) be able to detect it. However, if P has been packed the anti-virus will try to match the signature of P on P′. As the malicious code of P is encrypted in P′, no match will be found. Therefore, P will evade detection and infect the victim machine, if P′ is executed [14].

PE (Portable executable) file format is a standard Windows executable file format, which plays a very important role in the Windows operating system. PE files are widely used in Win32 executable programs including EXE, DLL, OCX, SYS, SCR and so on. PE viruses are designed in the way making use of the characteristics of PE file structure, and are portable on different hardware platforms, which is a serious security threat to the Windows operating system [21].

The very first step in the unpacking of packed file is to detect packed executable files. Recently, many researchers and analysts have focused on packed file detection techniques. In this paper, however, a new lightweight packed PE file detection technique based on the analyze the complexity of PE files by several algorithms. Packed PE files were analyzed using the proposed technique. It was found that nearly every type of packed PE file has higher complexity than it in unpacked status.

The methods always used are intelligent or it depends on database, but there are two major problems in them. Firstly, masses of malicious and benign codes as training data set are difficult to collect. Secondly, it would consume a lot of time to train the classifiers, and so the efficiency of the detection of unknown virus is dissatisfactory and difficult to use in practice.

A few generic and automatic unpacking techniques have been proposed to unpack packed binaries without specific knowledge of the packing technique used, e.g., OmniUnpack, Justin, Renovo, PolyUnpack and others [7].

The objective is to accurately distinguish between packed and non-packed executables, so that only the executables detected as packed will be sent to a computationally expensive general unpacker for hidden code extraction, before being sent to the antivirus software.

Therefore, the classification system here helps in improving virus detection while saving a significant amount of processing time. This paper do not focus on the improvements in virus detection accuracy achieved after unpacking, because this has already been studied in other researchs, for example. Instead, it focus on the accuracy and computational cost related to the classification of packed executables into the two classes packed and non-packed.

## II. RELATED WORKS

Most Packer Detection Methods can be summed up by: Signature based (Executable code signatures) and Heuristics (Entropy Checks, Import Address Table, Other Checks (not exclusive to packers)).

Coogan et al. [3] proposed an automatic static unpacking mechanism. It uses static analysis techniques to identify the unpacking code that comes with a given malware binary, then uses this code to construct a customized unpacker for that binary. This customized unpacker can then be executed or emulated to obtain the unpacked malware code.

Exeinfo PE [4] is an ongoing work for packed PE file detection and PE header information extraction. It shows the entrypoint, file offset, compiler information and the unpack information of the input file.

Renovo [6] utilizes a virtual machine. By using a virtual machine, they run a packed executable and record memory writing operations on shadow memory. When execution flow reaches one of checked bits of the shadow memory, all the checked memory bits are dumped. Shadow memory is changed to extract hidden code from packed executables with multiple hidden layers. With this mechanism, Renovo can find hidden layers as well.

OmniUnpack [8] monitors the program execution and tracks written, as well as written-then-executed, memory pages. When the program makes a potentially damaging system call, OmniUnpack invokes a malware detector on the written memory pages. If the detection result is negative, execution is resumed. If new type of malware appears, the dangerous system calls they defined on their paper could not match.

OllyDbg [12] is a debugger that emphasizes binary code analysis, which is useful when source code is not available. It traces registers, recognizes procedures, API calls, switches, tables, constants and strings, as well as locates routines from object files and libraries. According to the program's help file, version 1.10 is the final 1.x release. Version 2.0 is in development and is being written from the ground up. The software is free of cost, but the shareware license requires users to register with the author. OllyDbg is only available in 32-bit binaries. OllyDbg shows the message box that the input file is packed when the file is detected as a packed or encrypted file.

PEiD [13] is most commonly used with signature-based packers, cryptors and compilers for PE file detection. At present, it can detect more than 600 different signatures in PE files. PEiD is unique in some regard when compared to other identifiers. Its detection rates are pretty good among the current identifiers. Moreover, it has a plugin interface that supports plugins such as Generic OEP Finder and Krypto ANALyzer. Finally, it is free and easy to use.

Robert, et al.[16] present an encrypted and packed malware detection technique based on entropy analysis. In their paper, they analyzed packed PE files via the byte distribution. A set of metrics are developed that analysts can use to generalize the entropy attributes of packed or encrypted executable and thus distinguish them from native (non-packed or unencrypted) executables. As such, this methodology computes entropy at a naive model level, in which entropy is computed based only on the occurrence frequency of certain bytes of an executable without considering how these bytes were produced. Entropy analysis examines the statistical

variation in malware executables, enabling analysts to identify packed and encrypted samples quickly and efficiently.

PolyUnpack [17] performs static analysis over a packed executable to acquire a model of what its execution would look like if it did not generate and execute code at runtime. When the first instruction of a sequence not found in the static model is detected, the unknown instruction sequence is written and the execution of the packed executable is halted.

Hump-and-Dump [20] is a different approach from other research. Hump-and-Dump tries to find the OEP. Using a characteristic of unpacking, it counts the number of loops used in unpacking. When the number of loops is greater than a threshold and no more big loops are used for the period of a threshold, the address of the loop end point is the OEP.

## III. PACKER

A packer is proposed to reduce file size at first. A packed executable file is a file applied packer. This packed executable file operates functionally same as original file. Fig. 2 illustrates packing operation of packer [9].



Fig. 2.   Packer structure [9].

In packing procedure, a packer compresses or encrypts the IMAGE SECTION of input file that is the packed data, and then insert additional UNPACKING SECTION HEADER and UNPACK SECTION which can decompress or decrypt the packed data. Lastly, Packer modifies the entry point to start instruction of UNPACKING SECTION. Those are all of packing process. Therefore, packed executable file has smaller size than original file size and same functional operation as original file. Fig. 3 illustrates the procedures of execution of packed executable file.

When a packed executable file is executed, PE loader loads the packed file to virtual memory, and then the instruction of UNPACKING SECTION that is indicated by entry point is executed. Next, UNPACKING SECTION decompresses PACKED SECTION which is original section(s). Lastly, UNPACKED SECTION is executed on virtual memory. That is why operation of packed executable file is functionally same.

However, a packed executable file has a different bytes structure with original file. Namely, packed executable file has a different signature with original file. Therefore, anti-virus scanner does not consider packer that cannot detect the packed executable file by a signature of original file.

There exists various packers such as UPX, FSG, ASPack, Morphine, Exestealph, Pecompact, Yodacrypt, MEW, Packman, Upack, RLPack, Icrypt, EXE Smasher, Themida, and etc. Also, these packers have lots of versions, and manual

packers which malware makers made exists. Malware maker is able to generate variant of malware using lots of packers to evade anti-virus scanner.

For instance, there exist one malware and three packers. Malware maker generates three variant malwares using three packers. If malware maker applies packers to three variant malwares repeatedly, lots of variant malwares can be generated. In this way, malware maker makes variant malwares using various packers. As a matter of fact, 92% of malwares are packed executable in 2006. Of course, there exists that usage of packer for protection of commercial programs from malicious reverse engineering, but this normal usage is less than 2% (in fact, there is no study about normal usage of packer). Thus, anti-malware methods such as 'exepacker blacklisting are proposed, that is packed executable files are considered as malware.



Fig. 3.   Execution operation of packed executable file [9].

Among some packers, the most widely used packer is the UPX, ASPack, Themida, and so on. In next section, only will describe characteristics of UPX packer because it is used to pack the dataset in this paper.

### A. UPX

The UPX(Ultimate Packer for eXecutables) was released in March 1998. That is the first beta version. And then recently the version 3.07 was released in September 2010, the UPX is created by the Markus Oberhumer and Laszlo Molnar. And that is distributed in GPL(General Public License). The UPX offers the more high compression ratio than the Winzip or GZIP, see Fig4. And the decompress speed is faster than the others compression applications. The compression speed of the UPX is about 10MB/sec on the Pentium 133 and about 200MB/sec on the Athlon 2000. Also the UPX supports many file formats and various platforms. As explained earlier, the

UPX compression ratio is superior other applications and is a commonly used algorithm. However, the UPX is already a widely known packing algorithm so, the packed binary as the UPX is able to unpack [19][11].



Fig. 4. UPX Compression Ratio [18]

## IV. PE FILE AND PE VIRUS

PE (Portable Executable) file is an important executable file format of Windows operating system . ALL win32 executable (except VxDs and 16-bit DLLs) are PE file format. Files of 32bit DLLs, COM files, OCX controls, Control Panel Applets (CPL files) and NET executables are all PE format. The portability of PE file format means that the file format can be used on all Win32 platforms, and PE loader can recognize and use the file format in all win32 platforms. PE viruses take advantages of the PE file format to spread themselves among different Win32 platforms. The data structure of PE file in memory is consistent with that on disk. PE file uses a flat address space in which all code and data are merged into a large structure. PE loader maps the disk file to the virtual address space by the mechanism of mapping file to the memory [21]. All of the data structures of PE file are defined in WINNT.H.

### A. PE virus

PE virus is a computer virus that can infect PE format file in Windows operating system. Most of PE viruses are written with Win32 assembly language. PE virus has no data section. Variables and data are all put in code section. There are several key technologies of Win32 virus like Virus address relocation, Obtaining API Address, Searching target files, Mapping files to the memory, The general process of virus infection, Returning to the host program [21], and others.

As mentioned before, signature is a specific bytes string. But when packer technique is applied to specific file, that file will have different file structure in comparison to the original file. It means that malware makers can generate variant of malware using packers. Thus, anti-virus scanner cannot detect variant of malware by original signature. Recently, almost 92% malwares are found to be protected by packers In particular, the packing of malware is the very first problem that an analyst should address. If it is impossible to unpack a packed executable file, the analysis is impossible because the codes cannot be understood [9].

## V. FILE ANALYSIS

Security researchers need to find ways to fight malware, i.e., they need to obtain malware samples, analyze them to gain an understanding of malware tactics and weaknesses, and use that understanding to develop effective countermeasures [1].

### A. Static Analysis

Static analysis is a generic term referring to analysis methods that do not involve executing the program to be analyzed, for the sake of brevity henceforth called a specimen. Static analysis can be used to gather a variety of information about a specimen, e.g., high-level information such as its file size, a cryptographic hash, its file format, imported shared libraries. Cryptographic hashes can be used to identify a specimen. Packer signatures or its entropy may be used to determine whether it might be runtime packed.

Static analysis has several advantages over dynamic approaches. As static methods do not involve executing a potentially malicious specimen, there is a lesser risk of damaging the system that analysis is performed on. Given availability of the right tools, it is also possible to perform the analysis on a platform that differs from the platform that the specimen is designed to run on, further mitigating the risk of damaging the analysis platform (e.g., by accidentally executing it). Furthermore, static analysis typically covers the whole specimen and not just those code paths that are executed for a set of inputs, like dynamic analysis.

### B. Dynamic analysis

Dynamic analysis is a way of analyzing an unknown program by executing it and observing its behavior. When executing potentially hostile code, careful consideration must be given to securing the analysis environment, so as not to risk its destruction or even damage to other computer systems on the same network.

## VI. UNPACKING TRADITIONAL METHODS AND THEIR LIMITATIONS

### 1) Signature-based Unpacking method.

The signature-based anti-virus scanner detects the malware by signature which exists in malware as a specific bytes string, so it has low false-negative rate. If no signature is matched with the target, anti-virus scanner will classify an input file as non-malware. However, malware maker uses various evasion techniques such as control-flow obfuscation, source obfuscation, instruction virtualization, and packer which combine all evading techniques. In fact, the packer is originally proposed to reduce file size, but malware maker misuse packer to hide its malicious intention [9]. PEiD is an example of signature-based packer[2][13].

### 2) Algorithm-based Unpacking method

Use of specific unpacking routines to recover the original code (i.e., one routine per packing algorithm). Their limitations are [8]:

- Every new packer requires a dedicated unpacking algorithm.

- New packers are created from existing ones at a rate of 10-15 per month.

### 3) Generic Unpacking method

Emulation/tracing of the execution until the unpacking routine terminates (e.g., PolyUnpack and Renovo). Their limitations are [8]:

- Unpacking is slow and is not suitable for end-user environments.

- Effectiveness depends on the fidelity of the emulation environment (packers leverage anti-emulation techniques).

*4) Heuristic*

This involves searching through the code in a file to determine whether that code takes actions that appear to be actions typical of a packed file. The more packed like code that is found, the more likely that a packed is present. Heuristics approach of detection to provide protection against new and unknown packer, but it is inefficient and inaccurate where it is usually resulting in false positives. And there is a difficult to describe a heuristic which will work on all kinds of computer systems.

## VII. PROPOSED METHOD AND GOALS

The main goals are to achieve high accuracy on packed file classification with appropriate performance for practical anti-virus scanner. In addition, the proposed method is not evaded by avoidance techniques. Ultimately, the goal is that reduce the malware infection.

To achieve these goals, the arbitrator of the packed executable file classification based on complexity, not signature, entropy, or characteristics. Since packed executable file that is compressed or encrypted usually has high complexity, this is easily the judge that executable file is packed or not. A complexity concept can measure the information quantity more correctly than entropy concept.

The programs in Portable Executable (PE) 32-bit and 64-bit Microsoft Windows operating systems format is used. And in order to classify an executable program, binary static analysis is used to extract information. This information allows us to translate each executable into a sequence string of bytes. Then apply complexity measures techniques to distinguish between packed and non-packed executables.

Figure 5 shows how the classifier may be used to improve virus detection accuracy with low overhead, compared to a system where all the executables are directly sent to the general unpacker. Once a PE executable is received, the classification system performs a static analysis of the PE file in order to measure the complexity of it. After that, the complexity obtained from the PE executable is compared with a threshold. If the executable is classified as packed, it will be sent to the general unpacker for hidden code extraction, and the hidden code will then be sent to the anti-virus scanner. On the other hand, if the executable is classified as non-packed, it will be sent directly to the anti-virus scanner. It is worth noting that the PE file classifier may erroneously label a non-packed executable as packed. In this case the general unpacker will not be able to extract any hidden code from the received PE file. Nonetheless, this is not critical because if no hidden code is extracted, the AV scanner will simply scan the original non-packed code. The only cost paid in this case is the time spent by the general unpacker in trying to unpack a non-packed executable. On the other hand, the PE classifier may in some cases classify a packed executable as non-packed. In this case, the packed executable will be sent directly to the anti-virus scanner, which may fail to detect the presence of malicious code embedded in the packed executable, thus causing a false negative. However, this PE file classifier has a very high accuracy and is therefore able to limit the false negatives due to these cases.



Fig. 5. Overview of classification Method and operations of Anti-Malware Scanner

The complexity concept is proposed by to complement the entropy concept to more exact measure information quantity. The complexity C(X) of a finite string X will be defined as the length of the shortest string of X. In other words, C(X) is the length of the shortest computer program that represents X and then stops. The computer program can be programming languages or any others [9]. Complexity function is defined as

$$C(X) = \min \{X\} \qquad (1)$$

For example, the finite string X as
$$\underbrace{111111\ldots..1}_{10;000 \text{ times}}$$

then, this X can represented as follow program.
print 10,000 times a '1'

However, a serious problem of complexity concept is incomputable, since finding optimal algorithm that makes the shortest length output program from input string x is infeasible. A good news is compression algorithm as same as the complexity concept [9].

$$\text{Compress}(X) = (X`) \qquad (2)$$

where X` is the compressed string of string X. Thus, various compression algorithms are used to measure the complexity. The definition of compression algorithm is reduced input size to best smallest output size using their algorithm. Therefore, the complexity can measured for the input file using compression algorithm for classification.

Almost of packed executable file is compressed or encrypted, so to classify packed executable file is done when file has high complexity. But the difference value between file length's before and after compress is lower in the case of packed file than in the unpacked file. So if the complexity lower than Th value it will be packed else, unpacked. So setting Th and choice the compression algorithm are important for accuracy and performance.

Three steps are considered for implementation:

*1) Scan the sequence string bytes of input file for unnecessary bytes and cancel them.*

*2) Compress the bytes string throughout the step1 by several compression algorithms and entropy.*

*3) Measure the Complexity of string throughout the step 2. By the follow operation.*

- C = Length of X / Length of Compress(X) in the case of compression algorithms are used, where X is input string,

- C= 8- entropy, in the case of entropy is used,

- C <= Th : packed executable file, else: unpacked executable file.

- And, it is packed executable file, if the decision of at least 5 compression algorithms and entropy is packed, else unpacked executable file.

## VIII. THE USED COMPRESSION ALGORITHMS

The entropy and compression algorithms used to measure the complexity in this paper can be summarized as following:

### A. Entropy Analysis

In information theory, entropy is a measure of uncertainty in a series of an information unit. Information is compressed by following a logical sequence. First, some repeated patterns are found in the information, and then the redundancies of the patterns are used to reduce the size of the information. That is, the number of patterns of the information is reduced by compression and a series of bits becomes more unpredictable, which is equivalent to uncertainty. Therefore, the measured entropy of compressed information is higher than of the original information. Shannon's formula is devised to measure information entropy, as follows [5]:

$$\text{H}(x) = -\sum_{i=1}^{n} \text{p}(i) \cdot \log_b \text{p}(i) \qquad (3)$$

where H(x) is the measured entropy value and p(i) is the probability of an ith unit of information in event x's series of n symbols. The base number of the logarithm can be any real number greater than 1. However, 2, 10, and Euler's number e are chosen in general. We choose b=2 so this formula generates entropy scores as real numbers; when there are 256 possibilities, they are bounded within the range of 0 to 8.

### B. LZO

Lempel–Ziv–Oberhumer (LZO) is a lossless data compression algorithm that is focused on decompression speed. It is a portable lossless data compression library written in ANSI C. It offers pretty fast compression and very fast decompression. Decompression requires no memory.

LZO is a data compression library which is suitable for data de-/compression in real-time. This means it favours speed over compression ratio [10].

It is a block compression algorithm—it compresses and decompresses a block of data. Block size must be the same for compression and decompression. The LZO library implements a number of algorithms with the following characteristics:

- Decompression is simple and *very* fast.

- Requires no memory for decompression.

- Compression is pretty fast.

- Requires 64 KiB of memory for compression.

- Includes compression levels for generating pre-compressed data which achieve a quite competitive compression ratio.

- Algorithm is thread safe.

### C. Deflate

Deflate is a data compression algorithm that uses a combination of the LZ77 algorithm and Huffman coding. Deflate is widely thought to be implementable in a manner not covered by patents. This has led to its widespread use, for example in gzip compressed files, PNG image files and the ZIP file format for which Katz originally designed it [23][22].

Compression is achieved through two steps:

- The matching and replacement of duplicate strings with pointers.

- Replacing symbols with new, weighted symbols based on frequency of use.

### D. LZW

LZW compression is named after its developers, A. Lempel and J. Ziv, with later modifications by Terry A. Welch. It is the foremost technique for general purpose data compression due to its simplicity and versatility. Typically, you can expect LZW to compress text, executable code, and similar data files to about one-half their original size. LZW also performs well when presented with extremely redundant data files, such as tabulated numbers, computer source code,

and acquired signals. Compression ratios of 5:1 are common for these cases. LZW is the basis of several personal computer utilities that claim to "double the capacity of your hard drive". LZW compression is always used in GIF image files, and offered as an option in TIFF and PostScript [19].

### E. Gzip

gzip is based on the DEFLATE algorithm, which is a combination of LZ77 and Huffman coding. DEFLATE was intended as a replacement for LZW and other patent-encumbered data compression algorithms which, at the time, limited the usability of compress and other popular archivers. "gzip" is often also used to refer to the gzip file format[23].

Although its file format also allows for multiple such streams to be concatenated (zipped files are simply decompressed concatenated as if they were originally one file), gzip is normally used to compress just single files. Compressed archives are typically created by assembling collections of files into a single tar archive, and then compressing that archive with gzip. The final .tar.gz or .tgz file is usually called a tarball.

gzip is not to be confused with the ZIP archive format, which also uses DEFLATE. The ZIP format can hold collections of files without an external archiver, but is less compact than compressed tarballs holding the same data, because it compresses files individually and cannot take advantage of redundancy between files (solid compression) [23].

### F. QuickLZ

QuickLZ is the world's fastest compression library, reaching 308 Mbyte/s per core. It can be used under a commercial license if such has been acquired or under GPL 1, 2 or 3 where anything released into public must be open source [15]. It characterize by:

- Simple to use and easy to integrate. Get done in minutes and continue developing!

- Streaming mode for optimal compression ratio of small packets down to 200 - 300 bytes in size.

- Auto-detection and fast treatment of incompressible data.

### IX. EXPERIMENTAL RESULTS

The dataset used in this paper consists of 250 benign unpacked programs that were randomly gathered from the system files of windows XP operating system, then these files are packed using UPX. Each set of unpacked .exe files, and packed .exe files are enter alone in the classifier and the last decision is not depend on one the other.

Table 1 and figure 6 show the higher detection rate (True Positive TP = 0.96) of unpacked files is for the Totality Algs, this mean that the False Positive is (FP = 0.04). While the lower detection rate (TP = 0.83) of unpacked files is for the Qlz, this mean that the higher False Positive is (FP = 0.17).

Table 2 and figure 7 show the higher detection rate (True Negative TN = 0.97) of unpacked files is for the Entropy, this mean that the lower False Negative is (FN = 0.03). While the lower detection rate (TN = 0.9) of unpacked files is for the Qlz, this mean that the higher False Negative is (FN = 0.1).

TABLE I.        250 Unpack .exe file

| Algorithm | Unpack | Pack | Detection Rate |
|---|---|---|---|
| Entropy | 57 | 193 | 0.228 |
| LZO | 233 | 17 | 0.932 |
| QLZ | 207 | 43 | 0.828 |
| Gzip | 235 | 15 | 0.94 |
| Deflate | 235 | 15 | 0.94 |
| LZW | 239 | 11 | 0.956 |
| Totality Algs | 240 | 10 | 0.965 |



Fig. 6.   250 UnPack .exe file

TABLE II.        250 Pack .exe file

| Algorithm | Unpack | Pack | Detection Rate |
|---|---|---|---|
| Entropy | 7 | 244 | 0.976 |
| LZO | 25 | 225 | 0.9 |
| QLZ | 15 | 235 | 0.94 |
| Gzip | 16 | 234 | 0.936 |
| Deflate | 17 | 233 | 0.932 |
| LZW | 9 | 241 | 0.964 |
| Totality Algs | 18 | 232 | 0.928 |

Fig. 7.    250 Pack .exe file

## X.    Conclusion

The main goal in this paper, is to classify a packed and unpacked executable file in simple manner and achieve high accuracy, non-evade technique, and efficiency that can apply practical anti-virus scanner. These goals ultimately contribute to the anti-virus scanner that reduces malware infection with a little overhead.

Complexity is measured using entropy and five known compression algorithms. And the advantage of using complexity analysis is that it offers a convenient and quick technique for analyzing a sample at the binary level and identifying suspicious PE file (packed and encrypted Executables). This Generic unpacking has low-overhead by using existing hardware mechanisms, and it is characterized by fast, detect unknown packers, and resilient to anti-debugging.

For future works to enhance the detection rates use an artificial technique to segment the PE file and eliminate the less important segments, and further to add this packing/unpacking detection step to unpacking system.

### References

[1]    Bohne L., 2008, "Pandora's Bochs: Automatic Unpacking of Malware", Diploma thesis, Laboratory for Dependable Distributed Systems University of Mannheim.

[2]    Choi Y., Kim I., Oh J., and Ryou J., 2009, "Encoded Executable File Detection Technique via Executable File Header Analysis", International Journal of Hybrid Information Technology, Vol.2, No.2, pp. 25-36.

[3]    Coogan K., Debray S., Kaochar T., and Townsend G., "Automatic static unpacking of malware binaries". In WCRE '09: Proceedings of the 2009 16th Working Conference on Reverse Engineering, pages 167–176, Washington, DC, USA, 2009. IEEE Computer Society, Iraq Virtual Science Library.

[4]    Kang M. G., Poosankam P., and Yin H.. "Renovo: a hidden code extractor for packed executables". In WORM '07: Proceedings of the 2007 ACM workshop on Recurring malcode, pages 46–53, New York, NY, USA, 2007. ACM.

[5]    L. Limin, M. Jiang, W. Zhi, G. Debin, and J. Chunfu, "Denial-of-Service Attacks on Host-Based Generic Unpackers", 2010.

[6]    Martignoni L., Christodorescu M., and Jha S., "OmniUnpack Fast, Generic, and Safe Unpacking of Malware", ACSAC 2007, Iraq Virtual Science Library.

[7]    Noh H., 2009, "Complexity-based Packed Executable Classification with High Accuracy", Master Thesis, School of Engineering, Information and Communications University, Korea.

[8]    Oberhumer M., http://www.oberhumer.com/opensource/lzo, Version: 2.06, Date: 12 Aug 2011.

[9]    Oberhumer M., Molnar L. & Reiser J., 1996-2010, "The Ultimate Packer for eXecutables", UPX, http://upx.sourceforge.net

[10]    OllyDbg homepage, http://www.ollydbg.de/

[11]    PEiD homepage, http://www.peid.info/

[12]    Perdisci R., Lanzi A., and Lee W., 2008, "Classification of Packed Executables for Accurate Computer Virus Detection", Elsevier, Iraq Virtual Science Library.

[13]    QuickLZ 1.5.x, http://www.quicklz.com/index.php, 2013

[14]    Robert, Lyda, et al, "Using Entropy Analysis to Find Encrypted and Packed Malware", IEEE Security and Privacy, Apr. 2007, Iraq Virtual Science Library.

[15]    Royal P., Halpin M., Dagon D., Edmonds R., and Lee W, "Polyunpack: Automating the hidden-code extraction of unpack-executing malware". In ACSAC'06: Proceedings of the 22nd Annual Computer Security Applications Conference, pages 289–300, Washington, DC, USA, 2006. IEEE Computer Society, Iraq Virtual Science Library.

[16]    Shin D., Im C., Jeong H., Kim S., and Won D., 2011, "The new signature generation method based on an unpacking algorithm and procedure for a packer detection", International Journal of Advanced Science and Technology Vol. 27, February, pp 59-78.

[17]    Steven W. Smith , "The Scientist and Engineer's Guide to Digital Signal Processing", copyright ©1997-1998.

[18]    Sun L., Ebringer T., and Boztas S., "Hump-and-dump: efficient generic unpacking using an ordered address execution histogram". 2nd Int'l CARO Workshop, May 2008.

[19]    Tian Z., Sun X. and Yang H., 2011, "A Scheme of PE Virus Detection Using Fragile Software Watermarking Technique", International Journal of Digital Content Technology and its Applications. Volume 5, Number 2, February, pp. 158-164.

[20]    Wagner C., "Data Compression- DEFLATE Algorithm", Spring Semester 2011.

[21]    Wikipedia, the free encyclopedia, http://en.wikipedia.org/wiki/, 2013.

# Laguerre Kernels –Based SVM for Image Classification

Ashraf Afifi

Computer Engineering Department
Faculty of Computers and Information Technology
Taif University
Taif, KSA

*Abstract*—Support vector machines (SVMs) have been promising methods for classification and regression analysis because of their solid mathematical foundations which convey several salient properties that other methods hardly provide. However the performance of SVMs is very sensitive to how the kernel function is selected, the challenge is to choose the kernel function for accurate data classification. In this paper, we introduce a set of new kernel functions derived from the generalized Laguerre polynomials. The proposed kernels could improve the classification accuracy of SVMs for both linear and nonlinear data sets. The proposed kernel functions satisfy Mercer's condition and orthogonally properties which are important and useful in some applications when the support vector number is needed as in feature selection. The performance of the generalized Laguerre kernels is evaluated in comparison with the existing kernels. It was found that the choice of the kernel function, and the values of the parameters for that kernel are critical for a given amount of data. The proposed kernels give good classification accuracy in nearly all the data sets, especially those of high dimensions.

*Keywords—Laguerre polynomials; kernel functions; functional analysis; SVMs; classification problem*

## I.  INTRODUCTION

Improving efficacy of classifiers have been an extensive research area in machine learning over the past two decades, which led to state-of-the-art classifiers like support vector machines, neural networks and many more. Support vector machine (SVM) is a robust classification tool, effectively over comes many traditional classification problems like local optimum and curse of dimensionality[1].Support vector machines (SVMs) algorithm [2-3] has been shown to be one of the most effective machine learning algorithms. It gives very good results in terms of accuracy when the data are linearly or non-linearly separable. When the data are linearly separable, the SVMs result is a separating hyperplane, which maximizes the margin of separation between classes, measured along a line perpendicular to the hyperplane. If data are not linearly separable, the algorithm works by mapping the data to a higher dimensional feature space (where the data becomes separable) using an appropriate kernel function and a maximum margin separating hyperplane is found in this space. Thus the weight vector that defines the maximal margin hyperplane is a sufficient statistic for the SVMs algorithm (it contains all the information needed for constructing the separating hyperplane). Since this weight vector can be expressed as a weighted sum of a subset of training instances, called support vectors, it follows that the support vectors and the associated weights also constitute sufficient statistics for learning SVMs from centralized data.

One issue for improving the accuracy of SVMs is finding an appropriate kernel for the given data to improve the accuracy of SVMs. Most research relies on a priori knowledge to select the correct kernel, and then tweaks the kernel parameters via machine learning or trial-and-error. While there exist rules-of-thumb for choosing appropriate kernel functions and parameters, this limits the usefulness of SVMs to expert users, especially since different functions and parameters can have widely varying performance. Williamson et al.[4] published a method for the use of entropy numbers in choosing an appropriate kernel function. It was an attempt to explain kernel function choice by more analytical means rather than previous ad-hoc or empirical methods. The entropy numbers associated with mapping operators for Mercer kernels is discussed. In [5], it was stated that previous work on invariance transformations was mostly appropriate only for linear SVM classifiers. For non-linear SVM classifiers, an analytical method of utilizing kernel principal component analysis (PCA) map for incorporating invariance transformations was presented in[6].

Tsang et al.[7] discussed a way to take advantage of the approximations inherent in kernel classifiers, by using the Minimum Enclosing Ball algorithm as an alternative means of speeding up training. Training time had previously been reduced mostly by modifying the training set in some way. Their final classifiers, which they called the Core Vector Machine, converged in linear time with space requirements independent of the number of data points. Zanaty and Aljahdali [8] investigated the performance of different kernels when they are applied to different data sets. Zanaty et al. [9-10] combined GF and RBF functions in one kernel called "universal kernel" to take advantage of their respective strengths. The universal kernels constructed the most established kernels such as radial bases, gauss, and polynomial functions by optimizing the parameters using the training data. SedatOzer et al., [11] introduced a set of new kernel functions derived from the generalized Chebyshev polynomials, where the generalized Chebyshev kernel approaches the minimum support vector number and maximum classification performance. Zhi-Bin Pan et al. [12] introduced support vector machine based on orthogonal Legendre polynomials, to reduce

the redundancy in feature space due to the orthogonality of Legendre polynomials, which may enable the SVM to construct the separating hyperplane with less support vectors. These kernels satisfy Mercer's condition and converge faster than the existing kernels.

Completely achieving a SVM with high accuracy classification therefore, requires specifying high quality kernel function. In this paper, a new set of Laguerre functions is introduced that could improve the classification accuracy of SVMs. A class of Laguerre kernel functions on the basis of the properties of the common kernels is proposed, which can find numerous applications in practice. The proposed set of kernel functions provides competitive performance when compared to all other common kernel functions on average for the simulation datasets. The results indicate that they can be used as a good alternative to other common kernel functions for SVM classification in order to obtain better accuracy.

The rest of this paper is organized as follows: In section 2, SVM classifiers are discussed. The kernel functions are discussed in section 3. The generalized Legendre kernels are discussed in section 4. Section 5 presents the functional analysis of the proposed Laguerre kernels. Experimental and comparative results are given in section 6. Finally, section 7 shows the conclusion.

## II. SVM CLASSIFIER

SVMs [14] are a relatively new approach for creating classifiers that have become increasingly popular in the machine learning community. They present several advantages over other methods like neural networks in areas like training speed, convergence, complexity control of the classifier, as well as a stronger mathematical background based on optimization and statistical learning theory. In the novel learning paradigm embodied in support vector machines "learning" (selection, identification, estimation, training or tuning), the parameters are not predefined and their number depends on the training data used [14-15]. The support vector machines combine two main ideas. The first one is concept of an optimum linear margin classifier, which constructs a separating hyperplane that maximizes distances to the training point. The second one is concept of a kernel. In its simplest form, the kernel is a function which calculates the dot product of two training vectors. Kernels calculate this dot product in feature space, often without explicitly calculating the feature vectors, operating directly on the input vectors instead. When we use feature transformation, which reformulates input vector into new features, the dot product is calculated in feature space, even if the new feature space has higher dimensionality. So the linear classifier is unaffected. Margin maximization provides a useful trade off with classification accuracy, whichcan easily lead to overfitting of the training data. Consider aninput space $X$ with input vectors $x$, a target space $Y$ = {1,-1} and a training set$T_r$= *{*(x₁, $y_1$),...,(x$_N$, $y_N$)*}* with $x_i ∈ X$ and $y_i ∈ Y$. In SVM classification, separation of the two classes $Y$ = {1,-1} is done by means of the maximum margin hyperplane, i.e. the hyperplane that maximizes the distance to the closest data points and guarantees the best generalization on new, unseen examples. Let us consider two hyperplanes:

$$\langle w, x_i \rangle + b \geq 1 \qquad if\,(y_i = 1) \tag{1}$$

$$\langle w, x_i \rangle + b \leq 1 \qquad if\,(y_i = 1) \tag{2}$$

The distance from the hyperplane to a point $x_i$ can be written:

$$d(w,b;x_i) = \frac{\left|\langle w, x_i \rangle + b\right|}{\|w\|}$$

Consequently the margin between two hyperplanes can be written as:

$$\min_{x_i;y_i} d(w,b;x_i) + \min_{x_i;y_i} = -1d(w,b;x_i)$$

To maximize this margin we have to minimize ||w||. This comes down to solving a quadratic optimization problem with linear constraints. Notice however that we assumed that the data in $T_r$ are perfectly linear separable. In practice however this will often not be the case.

Therefore we employ the so called soft-margin method in contrast to the hard-margin method. Omitting further details we can rewrite the soft-margin optimization problem by stating the hyperplane in its dual form, i.e. find the Lagrange multipliers $\alpha_i \geq 0$ ($i = 1,...,N$) that *Maximize* :

$$L(\alpha_1,...........,\alpha_N) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i,j=1}^{N}\alpha_i\alpha_j y_i y_j \langle x_i, x_j \rangle \tag{3}$$

subject to $\sum_{i=1}^{N}\alpha_i y_i = 0 \qquad 0 \leq \alpha_i \leq C$

Considering the dual problem above, we can now write the maximum margin hyperplane as a linear combination of support vectors. By definition, the vectors $x_i$ corresponding with non-zero $\alpha_i$ are called the support vectors and this set consists of those data points that lie closest to the hyperplane and thus are the most difficult to classify. In order to classify a new point $x_{new}$, one has to determine the sign of

$$\sum_{x_i \in SV} \alpha_i y_i \langle x_i, x_{new} \rangle + b$$

$$(4)$$

If this sign is positive $x_{new}$ belongs to class 1, if negative to class -1, if zero x$_{new}$ lies on the decision boundary. Note that we have restricted the summation to the set of support vectors because the other $\alpha_i$are zero anyway.

## III. KERNEL FUNCTIONS

Support vector machine is one of kernel-based learning algorithms that consist of a learning algorithm and the kernel function [16, 17]. The kernel function creates the hypothesis space where the learning process searches for. The kernel can be considered as a similarity measure between two inputs which corresponds to their inner product in some feature space into which the original inputs are mapped. This is very useful, for instance, when the concept to be learned depends nonlinearly on the data, but the learning algorithm is able to learn only linear dependencies.

Since support vector machines are linear classifiers, it is necessary to map the input vectors with a nonlinear mapping in order to learn non-linear relations. The resulting vectors are usually called features. Formally, let $k$ denote the input space, which can be any set, and $F$ denote the feature vector space. For any mapping:

$$\phi : k \rightarrow F$$

The inner product of the mapped inputs is called a kernel function:

$$k(x, z) = \langle \phi(x), \phi(z) \rangle$$

A necessary condition for this is that $k(x, z)$ is symmetric and finitely positive semi definite [18-19]. There are many different types of kernels that can be found in the literature [18-19].

## IV. PROPOSED KERNELS

A critical step in support vector machine classification is choosing a suitable kernel of SVMs for a particular application, i.e. various applications need different kernels to get reliable classification results. It is well known that the two typical kernel functions often used in SVMs are the radial basis function kernel and polynomial kernel. More recent kernels are presented in [9-12,20-23] to handle high dimension data sets and are computationally efficient when handling non-separable data with multi attributes. However, it is difficult to find kernels that are able to achieve high classification accuracy for a diversity of data sets. In order to construct kernel functions from existing ones or by using some other simpler kernel functions as building blocks, the closure properties of kernel functions are essential [16-18].

For given non-separable data, in order to be linearly separable, a suitable kernel has to be chosen. Classical kernels, such as Gauss RBF and POLY functions, can be used to transfer non-separable data to separable, but their performance in terms of accuracy is dependent on the given data sets. The following POLY function performs well [20] with nearly all data sets, except high dimension ones:

$$k(x, z) = \left( x^T z + 1 \right)^d$$

where $d$ is the polynomial degree.

The same performance [20] is obtained with the Gauss RBF of the following form:

$$k(x, z) = \exp\left( -\gamma \| x - z \|^2 \right)$$

where $\gamma$ is appositive parameter controlling the radius. Zanaty et al in [9] presented the polynomial Radial basis function (RBPF) as:

$$PRBF = \left( (1 + \exp(\theta)) / V \right)^d$$

Where $\theta = |x - z|, V = p * d$

where $p$ is a parameter. Zanaty et al in [10] presented Support vector machines (SVMs) with universal kernels, called Gaussian radial *bas*is polynomials function (GRPF) given by:

$$GRBF = \left( \frac{b + e^{-\gamma^a \| x - x' \|^r}}{V} \right), V = PD, a = b-1, r >= 1.$$

where $a, b$, $\gamma = \dfrac{1}{2\sigma^2}$, and $d$ are the kernel parameters for the Gaussian, polynomial and universal kernels, respectively. $\beta$ and $\sigma$ are the scaling parameters for the polynomial kernel and determines the width of the Gaussian kernel respectively. Kernel functions should be applied onto input vectors directly instead of applying them onto each element and combining the results by a product, since the kernel functions are supposed to provide a measure of the correlation of two input vectors in a higher dimensional space.

### A. Laguerre polynomials

The Laguerre polynomials are defined by the equation:

$$(1-t)^{-1} \exp\left( -\frac{xt}{1-t} \right) = \sum_{n=0}^{\infty} \frac{L_n(x)}{n!} t^n \tag{5}$$

The exponential function can be expanded to give:

$$\sum_{n=0}^{\infty} \frac{L_n(x)}{n!} t^n = \frac{1}{1-t} \sum_{\gamma=o}^{\infty} \frac{(-1)^\gamma x^\gamma t^\gamma}{r!(1-t)^\gamma} = \sum_{\gamma=0}^{\infty} \frac{(-1)^\gamma x^\gamma t^\gamma}{r!} (1-t)^{-(\gamma+1)} \tag{6}$$

Recall the binomial expansion:

$$(1-t)^{-m} = 1 + (-m)(-t) + \frac{(-m)(-m-1)}{2!}(-t)^2 + \dots = \sum_{s=0}^{\infty} \frac{m(m+1)\dots(m+\varepsilon-1)}{\varepsilon!} t^s \tag{7}$$

using the notation:

$$(\alpha)_\gamma = \alpha\,(\alpha+1)(\alpha+2)\dots(\alpha+r-1) \tag{8}$$

Equation (6) may therefore be written as

$$\sum_{n=0}^{\infty} \frac{L_n(x)}{n!} t^n = \sum_{\gamma=0}^{\infty} \sum_{s=0}^{\infty} \frac{(-1)^\gamma (r+1)_s x^\gamma t^{\gamma+s}}{r!\ s!} \tag{9}$$

Equating powers of $t^n$, we get:

$$L_n(x) = n! \sum_{\gamma=0}^{n} \frac{(-1)^\gamma (r+1)_{n-\gamma}}{r!(n-r)!} x^\gamma \tag{10}$$

The summation is taken from $r = 0$ to $r = n$, in view of the factor $(n-r)!$ in the denominator. Note that:

$$(r+1)_{n-\gamma} = (r+1)(r+2)\dots(r+1+n-r-1) = \frac{n!}{r!} \tag{11}$$

$$(-n)_\gamma = (-n)(-n+1)\dots(-n+r-1) = (-1)^\gamma n(n-1)\dots(n-r+1) = (-1)^\gamma \frac{n!}{(n-r)!} \tag{12}$$

$$(1)_\gamma = (1)(2)\dots(1+r-1) = r! \tag{13}$$

Equation (5) then gives

$$L_n(x) = n! \sum_{\gamma=0}^{n} \frac{(-n)_\gamma}{(1)_\gamma} \cdot \frac{x^\gamma}{r!}$$

(14)

*Or* $\quad L_n(x) = n! F_1(-n ; 1 ; x)$

(15)

Note that the series in Equation (14) terminates after *n* terms, i.e. $L_n(x)$ is a polynomial of degree *n* (see Fig.(1)).



L(x)

$L_2(x)$
$L_3(x)$
$L_4(x)$
$L_5(x)$
$L_1(x)$

Fig. 1.  The Laguerre function for the first five polynomials.

B.1 *Rodrigues' Formula for the Laguerre polynomial*
Using Equation (11), Equation (11) may be written as

$$L_n(x) = e^x \sum_{\gamma=0}^{n} \frac{n!}{r!(n-r)!}(-1)^\gamma e^{-x} \frac{n!}{[n-(n-r)]!} x^{n-(n-\gamma)}$$

(16)

Recall the Leibniz formula for the $m^{th}$ derivative of a product:

$$\frac{d^m}{dx^m}(AB) = \sum_{\gamma=0}^{m} \frac{m!}{r!(m-r)!}\left(\frac{d^\gamma}{dx^\gamma}A\right)\left(\frac{d^{m-\gamma}}{dx^{m-\gamma}}B\right)$$

(17)

and that

$$\frac{d^p}{dx^p} x^n = \frac{n!}{(n-p)!} x^{n-p} \qquad (n >= p)$$

(18)

Equation (13) may therefore be written as

$$L_n(x) = e^x \sum_{\gamma=0}^{n} \frac{n!}{r!(n-r)!} \frac{d^\gamma}{dx^\gamma} e^{-x} \frac{d^{n-\gamma}}{dx^{n-\gamma}} x^n$$

(19)

*or* $\quad L_n(x) = e^x \frac{d^n}{dx^n}(e^{-x}x^n)$

(20)

Laguerre polynomials of low order can be evaluated by using the Rodrigues' formula (21) :

$L_0(x) = 1$
$L_1(x) = 1 - x$
$L_2(x) = 2 - 4x + x^2$
$L_3(x) = 6 - 18x + 9x^2 - x^3$

(21)

B.  *Recurrence Relations*
We write the defining equation (5) in the form

$$\exp\left[x\left(1 - \frac{1}{1-t}\right)\right] = (1-t)\sum_{n=0}^{\infty} \frac{L_n(x)}{n!}t^n$$

(22)

Differentiating both sides with respect to *t*, we get

$$-\frac{x}{(1-t)^2}\exp\left[x\left(1 - \frac{1}{1-t}\right)\right] = (1-t)\sum_{n=0}^{\infty} \frac{L_n(x)}{(n-1)!}t^{n-1} - \sum_{n=0}^{\infty} \frac{L_n(x)}{(n-1)!}t^n$$

(23)

Multiplying through by $(1-t)$, we get

$$x\sum_{n=0}^{\infty} \frac{L_n(x)}{n!}t^n + (1-t)^2\sum_{n=0}^{\infty} \frac{L_n(x)}{(n-1)!}t^{n-1} - (1-t)\sum_{n=0}^{\infty} \frac{L_n(x)}{n!}t^n = 0$$

(24)

and equating coefficients of $t^n$, we obtain

$$\frac{xL_n(x)}{n!} + \frac{L_{n+1}(x)}{n!} - 2\frac{L_n(x)}{(n-1)!} + \frac{L_{n-1}(x)}{(n-2)!} - \frac{L_n(x)}{n!} + \frac{L_{n-1}(x)}{(n-1)!} = 0$$

(25)

and hence the recurrence relation

$$L_{n+1}(x) + (x - 2n - 1)L_n(x) + n^2 L_{n-1}(x) = 0$$

(26)

If, on the other hand, we differentiate Equation (23) with respect to *x*, we get

$$-\frac{t}{1-t}\exp\left[x\left(1 - \frac{1}{1-t}\right)\right] = (1-t)\sum_{n=0}^{\infty} \frac{L'_n(x)}{n!}t^n$$

(27)

*or* $\quad t\sum_{n=0}^{\infty} \frac{L_n(x)}{n!}t^n + (1-t)\sum_{n=0}^{\infty} \frac{L'_n(x)}{n!}t^n = 0$

(28)

Equating coefficients of $t^n$ yields the identity

$$\frac{L'_n(x)}{n!} - \frac{L'_{n-1}(x)}{(n-1)!} + \frac{L_{n-1}(x)}{(n-1)!} + \frac{L_{n-1}(x)}{(n-1)!} = 0$$

(29)

and hence the recurrence relation

$$L'_n(x) - nL'_{n-1}(x) + nL_{n-1}(x) = 0$$

(30)

C.  *Orthogonality of the Laguerre Polynomials*
Laguerre's differential equation can be cast into self-adjoint form by first writing it as

$$x\frac{d^2y}{dx^2} + \frac{(1-x)}{x}\frac{dy}{dx} + \frac{n}{x}y = 0$$

(31)

and multiplying throughout by the "integrating factor"

$$\exp\left(\int\frac{1-x}{x}dx\right) = \exp(\ln x - x) = xe^{-x}$$

This gives

$$xe^{-x}\frac{d^2y}{dx^2} + (1-x)e^{-x}\frac{dy}{dx} + ne^{-x}y = 0$$

(32)

*or*  $\dfrac{d}{dx}\left(xe^{-x}\dfrac{dy_n}{dx}\right)+ne^{-x}y_n=0$

(33)

where we have included the subscript $n$ in order to associate the solution $y_n$ with the eigenvalue $n$. Replacing $n$ by $m$, we have:

$$\frac{d}{dx}\left(xe^{-x}\frac{dy_n}{dx}\right)+me^{-x}y_m=0$$

(34)

We now adopt the standard procedure: multiply Equation (39) by $y_m$ and Equation (40) by $y_n$, and subtract. This gives:

$$(m-n)e^{-x}y_m(x)y_n(x)=y_n(x)\frac{d}{dx}[xe^{-x}y_m'(x)]-y_m(x)\frac{d}{dx}[xe^{-x}y_n'(x)]$$

(35)

Integrating both sides from $0$ *to* $\infty$ , and using the rule for the derivative of a product, we get

$$(m-n)\int_0^\infty e^{-x}y_m(x)y_n(x)dx=\int_0^\infty \frac{d}{dx}[xe^{-x}\{y_n(x)y_m'(x)-y_m(x)y_n'(x)\}]dx=0$$

(36)

It follows that if $m\neq n$ the integral on the Left must be zero, i.e. $\int_0^\infty e^{-x}y_m(x)y_n(x)dx=0, \qquad if \ \ m\neq n$

(37)

We have shown that the Laguerre polynomial $L_n(x)$ is a solution of Laguerre's equation. We may therefore substitute $y_n(x)=L_n(x)$ in Equation (37), to give:

$$\int_0^\infty e^{-x}L_m(x)L_n(x)dx=0, \qquad if \ \ m\neq n$$

(38)

The case $m=n$ can be examined by noting that:

$$L_n(x)=e^x\frac{d^n}{dx^n}(x^ne^{-x})=n!\sum_{\gamma=0}^n\frac{(-n)_\gamma}{r!r!}x^\gamma$$

(39)

Thus we have:

$$\int_0^\infty e^{-x}L_n(x)L_n(x)dx=n!\sum_{\gamma=0}^n\frac{(-n)_\gamma}{r!r!}\int_0^\infty e^{-x}x^\gamma e^x\frac{d^n}{dx^n}(x^ne^{-x})dx$$

(40)

The Right Hand Side may be evaluated by integrating by parts $n$ times. The procedure is as follows:

$$\int_0^\infty x^\gamma\frac{d^n}{dx^n}(x^ne^{-\gamma})dx=\int_0^\infty x^\gamma d[\frac{d^{n-1}}{dx^{n-1}}(x^ne^{-x})]=[x^\gamma\frac{d^{n-1}}{dx^{n-1}}(x^ne^{-x})]_0^\infty-\int_0^\infty \frac{dx^\gamma}{dx}[\frac{d^{n-1}}{dx^{n-1}}(x^ne^{-x})]dx$$

(41)

The first term on the Right Hand Side vanishes at both limits, so we obtain:

$$\int_0^\infty x^\gamma\frac{d^n}{dx^n}(x^ne^{-x})dx=-\int_0^\infty \left(\frac{d}{dx}x^\gamma\right)[\frac{d^{n-1}}{dx^{n-1}}(x^ne^{-x})]dx=-\int_0^\infty \left(\frac{d}{dx}x^\gamma\right)d[\frac{d^{n-2}}{dx^{n-2}}(x^ne^{-x})]$$

(42)

$$\int_0^\infty x^\gamma\frac{d^n}{dx^n}(x^ne^{-x})dx=(-1)^2\int_0^\infty \left(\frac{d^2}{dx^2}x^\gamma\right)[\frac{d^{n-2}}{dx^{n-2}}(x^ne^{-x})]dx$$

(43)

A continuation of this process leads ultimately to the result:

$$\int_0^\infty x^\gamma\frac{d^n}{dx^n}(x^ne^{-x})dx=(-1)^n\int_0^\infty \left(\frac{d^n}{dx^n}x^\gamma\right)[\frac{d^{n-n}}{dx^{n-n}}(x^ne^{-x})]dx$$

(44)

Hence

$$\int_0^\infty e^{-x}L_n(x)L_n(x)dx=n!\sum_{r=0}^n\frac{(-n)_\gamma}{r!r!}(-1)^n\int_0^\infty \left(\frac{d^n}{dx^n}x^\gamma\right)(x^ne^{-x})dx$$

(45)

Recall that $\dfrac{d^n}{dx^n}x^\gamma=\begin{cases}0 & if \ \ r<n\\ n! & if \ \ r=n\end{cases}$

(46)

Thus the only term in the summation (45) which survives is the $r=n$ term; hence we obtain:

$$\int_0^\infty e^{-x}L_n(x)L_n(x)dx=n!\frac{(-n)_n(-1)^n n!}{n!n!}\int_0^\infty x^ne^{-x}dx$$

(47)

Note that $(-n)_n=(-1)^n n!$

(48)

and $\int_0^\infty x^ne^{-x}dx=n!$

(49)

So we finally obtain:

$$\int_0^\infty e^{-x}L_n(x)L_n(x)dx=(n!)^2$$

(50)

This result may be combined with the orthogonality relation (38) to give:

$$\int_0^\infty e^{-x}L_m(x)L_n(x)dx=(n!)^2\delta_{mn}$$

(51)

The weight function $e^{-x}$ may be removed by defining a new function:

$$\phi_n(x)=\frac{1}{n!}e^{-x/2}L_n(x)$$

(52)

Equation (57) may then be written as:

$$\int_0^\infty \phi_m(x)\phi_n(x)dx=\delta_{mn}$$

(53)

Functions which satisfy the relation (54) are said to be *normalized*, and we say that the $\phi_n(x)$ form an *orthonormal* set of functions.

## V. Generalized Laguerre kernels

Here, we propose a generalized way of expressing the kernel function to clarify the ambiguity on how to implement Laguerre kernels. To the best of our knowledge, there was no previous work defining the Laguerrepolynomials for vector inputs recursively. Therefore for vector inputs, we define the generalized Laguerre polynomials as:

$$
\begin{aligned}
&L_0(x) = 1 \\
&L_0(z) = 1 \\
&L_1(x) = 1 - x \\
&L_1(z) = 1 - z \\
&L_2(x) = 0.5x^2 - 2x + 1 \\
&L_2(z) = 0.5z^2 - 2z + 1 \\
&L_3(x) = -0.17x^3 + 1.5x^2 - 3x + 1 \\
&L_3(z) = -0.17z^3 + 1.5z^2 - 3z + 1 \\
&L_4(x) = 0.042x^4 - 0.667x^3 + 3x^2 - 4x + 1 \\
&L_4(z) = 0.042z^4 - 0.667z^3 + 3z^2 - 4z + 1 \\
&\ldots \\
&L_{n+1}(x) + (x - 2n - 1)L_n(x) + n^2 L_{n-1}(x) = 0
\end{aligned}
\qquad (54)
$$

Therefore, the generalized Laguerre, $L_j(x), and\ L_j(z)$, yield rowvectors, otherwise, it yields a scalar value. Thus by using generalized Laguerre polynomials, we define generalized $n^{\text{th}}$order Laguerre kernel as

$$
k(x,z) = \sum_{j=0}^{n} L_j(x).L_j^T(z)
\qquad (55)
$$

Where $x$ and $z$ are m-dimensional vectors.

TABLE I. LIST OF THE GENERATED LAGUERRE KERNEL FUNCTIONS UP TO 4TH ORDER.

| Kernel Parameter: n | Kernel function: k(x,z) |
|---|---|
| 0 | $L_0(x)*L_0(z)^T$ |
| 1 | $K_0 + L_1(x)*L_1(z)^T$ |
| 2 | $K_0 + K_1 + L_2(x)*L_2(z)^T$ |
| 3 | $K_0 + K_1 + K_2 + L_3(x)*L_3(z)^T$ |
| 4 | $K_0 + K_1 + K_2 + K_3 + L_4(x)*L_4(z)^T$ |
| 5 | $K_0 + K_1 + K_2 + K_3 + L_4(x)*L_4(z)^T + L_5(x)*L_5(z)$ |
| 6 | $K_0 + K_1 + K_2 + K_3 + L_4(x)*L_4(z)^T + L_5(x)*L_5(z) + L_6(x)*L_6(z)$ |
| 7 | $K_0 + K_1 + K_2 + K_3 + L_4(x)*L_4(z)^T + L_5(x)*L_5(z) + L_6(x)*L_6(z) + L_7(x)*L_7(z)$ |

construct the mapping $\phi$ from the eigenfunctiondecomposition of $k$ .According to Mercer's work [24],it is known that if $k$ is the symmetrical and continuous kernel of an integraloperator $O_k : L^2 \to L^2$, such that:

$$
O_{kg}(x) = \int K(x,z)g(z)dz
$$

is positive, i.e.,

$$
\int K(x,z)g(x)g(z)dxdz \geq 0 \qquad \forall g \in L^2,
$$

Then $k$ can be expanded into a uniformly convergent series

$$
K(x,z) = \sum_{i=1}^{\infty} \lambda_i \phi_i(x)\phi_i(z),
$$

With $\lambda_i \geq 0$. In this case, the mapping from input space tofeature space produced by the kernel is expressed as

$$
\phi : x \mapsto (\sqrt{\lambda_1}\phi_1(x), \sqrt{\lambda_2}\phi_2(x),...)
$$

such that $k$ acts as the given dot product, i.e.,

$$
(\phi(x),\phi(z)) = \phi^T(x)\phi(z) = K(x,z).
$$

**Theorem 1.**A nonnegative linear combination of Mercer kernels is also a Mercer kernel.

**Proof.** Let $K_i (i = 1,...,M)$ be Mercer kernels and let

$$
K(x,z) = \sum_{i=1}^{M} a_i K_i(x,z),
$$

where $a_i \geq 0$is a nonnegative constant. According toMercer's theorem, we have

$$
\int K_i(x,z)g(x)g(z)dxdz \geq 0, \forall g \in L^2, i = 1,...,M.
\qquad (56)
$$

By taking the sum of the positive combination of (56)with coefficients $a_i$ over $i$, one obtains

$$
\sum_{i=1}^{M} a_i \int K_i(x,z)g(x)g(z)dxdz \geq 0, \quad \forall g \in L^2.
\qquad (57)
$$

Therefore, one reaches

$$
\int K(x,z)g(x)g(z)dxdz \geq 0, \quad \forall g \in L^2.
$$

In particular, if $\sum_{i=1}^{M} a_i = 1, (a_i \geq 0)$, then we regard $K(x,z)$ in (57) as the convex combination of the positivedefinite kernels $K_i(x,z)$. This kind of kernel can findnumerous applications in practice.

**Theorem 2**.The product of Mercer kernels is also a Mercer kernel.

**The proof**is similar to that of the precedingtheorem.

**Theorem 3**.To be a valid SVM kernel, a kernel should satisfy the Mercer Conditions [26-27]. If the kernel does not satisfy the Mercer Conditions, SVM may not find the

optimal parameters, but rather it may find suboptimal parameters. Also if the Mercer conditions are not satisfied, then the Hessian matrix for the optimization part may not be positive definite.Therefore we examine if the generalized Leguerre kernel satisfies the Mercer conditions:

Mercer Theorem: To be a valid SVM kernel, for any finite function $g(x)$, the following integration should always be non-negative for the given kernel function $k(x,z)$ [1]:

$$\iint K(x,z)g(x)g(z)\,dxdz \geq 0 \tag{58}$$

Where

$$k(x,z) = \sum_{j=0}^{n} L_j(x)L_j^T(z) = L_0(x)L_0^T(z) + L_1(x)L_1^T(z) + \dots\dots + L_n(x)L_n^T(z) \tag{59}$$

Consider that $g(x)$ is a function where $g: R^m \rightarrow R$, then we can evaluate and verify the Mercer condition for $k(x,z)$ as follows by assuming each element is independent from others:

$$\iint k(x,z)g(x)g(z)dxdz = \iint \sum_{j=0}^{n} L_j(x)L_j^T(z)g(x)g(z)dxdz$$

$$= \sum_{j=0}^{n} L_j(x)L_j^T(z)g(x)g(z)dxdz = \sum_{j=0}^{n}\left[\int L_j(x)g(x)dx \int L_j^T(z)g(z)dz \right]$$

$$= \sum_{j=0}^{n}\left[\left(\int L_j(x)g(x)dx \right)\left(\int L_j^T(x)g(x)dx \right) \right] \geq 0 \tag{60}$$

Therefore, the kernel $k(x,z)$ is a valid kernel.

## VI. EXPERIMENTAL RESULTS

The classification experiments are conducted on different data like Cloud, Liver, Seed, Forest Fire and Yeast dataavailable at http://archive.ics.uci.edu/ml/datasets.html. These data sets have been given to the algorithm with different sizes (classes and attributes). Table II shows the classification accuracy for five different data sets using Laguerre kernel of order from 2 to 5 implementations. As shown in Figure 2, it is clear that when the order of polynomials increases, the accuracy increases for all data sets.



Fig. 2. Laguerre Kernels-Based SVM Classification Accuracy

### A. Comparative results

The performance of the proposed kernel with SVMs, in terms of classification accuracy, is evaluated by application to a variety of data sets available at:

*http://www.cs.toronto.edu.delve/data/image-set/desc.html.* Firstly, we used LIBSVM with different kernels(linear, polynomial, radial basis function [8]). The parameters used include two parameters for the RBF kernel parameter $\gamma=0.5$ and $\sigma=0.5$, $d=1$ for linear and $d=5$ for polynomial kernels. Table II lists the main characteristics of the seven datasets used in the experiments. In order to evaluate the performance of the support vector machine with different kernels, we carried out some experiments with different data sets from machine learning benchmarks domains [28].

The data has even different classes of image. They contain 210 data for training and another 2100 data for testing, Each vector has 18 elements with different minimum and maximum values. For the training, we have 30 data for the class(+1) and180 data for the class (-1) and similarly for testing. As can be seen from Table II, the generalized Laguerre kernel results show better generalization ability than the existing Gaussian, Polynomial (POLY) and Chebyshev [12] kernels.

TABLE II. CLASSIFICATION ACCURACY OF DIFFERENT DATA SETS USING LAGUERRE KERNELFUNCTION

| Sr. No. | Name of Database | No. of features | No. of Classes | Training Size | Testing Size | Classification accuracy | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | $2^{nd}$ order | $3^{rd}$ order | $4^{th}$ order | $5^{th}$ order |
| 1 | cloud | 1024 | 10 | 150 | 450 | 0.8905 | 0.9048 | 0.9190 | **0.9190** |
| 2 | Liver | 345 | 07 | 140 | 205 | 0.8905 | 0.9524 | 0.9619 | **0.9619** |
| 3 | Seeds | 210 | 07 | 105 | 210 | 0.9048 | 0.9095 | 0.9190 | **0.9333** |
| 4 | Forest fires | 517 | 13 | 195 | 210 | 0.8667 | 0.8800 | 0.8952 | **0.9143** |
| 5 | Yeast | 1484 | 08 | 150 | 500 | 0.8857 | 0.8952 | 0.9095 | **0.9238** |

For example, in Table II, the 5[th] order generalized Laguerre kernel results in classification accuracy of more than 96% for all test data sets, while the existing kernels achieve less than 96% for most test data sets. More specific, comparing the results of Tables III, the 5[th]order generalized Laguerrekernel always gives good results and may be the best at all, as shown in Figure 3.

TABLE III.    RESULTS ON IMAGE SEGMENTATION DATA WITH VARIOUSKERNELFUNCTIONS.

| Data Segmentation | Linear Kernel | Gaussian Kernel | Polynomial Kernel | Chebyshev Kernel | Laguerre Kernel |
|---|---|---|---|---|---|
| Brickface | 0.9952 | 0.9809 | 0.9904 | 0.9952 | 0.9952 |
| Sky | 1 | 0.9857 | 1 | 1 | 1 |
| Foliage | 0.9762 | 0.9571 | 0.9571 | 0.9714 | 0.9857 |
| Cement | 0.8952 | 0.9476 | 0.9714 | 0.9762 | 0.9667 |
| Window | 0.9524 | 0.9476 | 0.9667 | 0.9476 | 0.9762 |
| Path | 1 | 0.9761 | 0.9952 | 0.9762 | 0.9952 |
| Grass | 1 | 0.9809 | 1 | 0.9857 | 1 |



Fig. 3. SVM Classification Accuracy with differentKernels

## VII.    CONCLUSION

In this paper, SVMs have been improved to solve the classification problems by mapping the training data into a feature space by the aid of Laguerre kernel functions and then separating the data using a large margin hyperplane. A class of Laguerre kernel functions on the basis of the properties of the common kernels is proposed, which can find numerous applications in practice.

Experimental results illustrate the validity and effectiveness of the proposed kernel. The experimental results show that the proposed kernel function results in the best accuracy in nearly all the data sets especially in the data set with large number of attributes. The obtained results are encouraging and suggest that the proposed method is worth further consideration.

REFERENCES

[1] Vapnik V. N., "The nature of statistical learning theory", Springer-Verlag, New York, NY, USA, 1995.

[2] Kim H., Pang S., Je H., Kim D., Bang S.Y., "Constructing support vector machine ensemble", Pattern Recognition", vol.36, no.12, pp.2757–2767, 2003.

[3] Du P., Peng J., Terlaky T.. "Self-adaptive support vector machines", modeling and experiments Computational Management Science, vol. 6, no.1, pp. 41–51, 2009.

[4] Williamson, R., A. Smola, and B. Schölkopf, "Entropy numbers, operators and support vector kernels", pp. 44-127, Cambridge, MA: MIT Press, 1999.

[5] Chapelle, O. and B. Schölkopf, "Incorporating invariances in non-linear support vector machines", In T. G. Dietterich, S. Becker, and Z. Ghahramani (Eds.), Advances in Neural Information Processing Systems, Volume 14, pp. 594-609, Cambridge, MA: MIT Press, 2002.

[6] Hansheng Lei &VenuGovindaraju, "Speeding up multi-class SVM evaluation by PCA and feature selection", Center for Unified Biometrics and Sensors (CUBS), 2005.

[7] Tsang, I., J. Kwok, and P.-M. Cheung, "Core vector machines: Fast SVMs training on very large data sets", Journal of Machine Learning Research 6, pp. 271-363, 2005.

[8] E.A.Zanaty and Sultan Aljahdali "Improving the accuracy of support vector machines" accepted in 23rd International conference on computers and their application, Mexico, USA, April 9-11,2008.

[9] Zanaty E.A, Sultan Aljahdali, R.J. Cripps, "Accurate support vector machines for data classification", Int. J. Rapid Manufacturing, vol. 1, no. 2, pp. 114-127, 2009.

[10] Zanaty E.A, Ashraf Afifi, "Support vector machines (SVMs) with universal kernels ", in International journal of Artificial Intelligence, vol. 25, pp.575-589, 2011.

[11] Zhi-Bin Pan, Hong C., Xin-H., "Support vector machine with orthogonal legendre kernel" In Procceeding of the international conference on wavelet analysis and pattern recogntion, pp. 15-17, 2012.

[12] Sedat O., ChiH.C., HakanA. Cirpan, "A set of new Chebyshev kernel functions for support vector machine pattern classification", Pattern Recognition, vol. 44, pp. 1435-1447, 2011.

[13] Boser, B., Guyon I., and Vapnik V., "A training algorithm for optimal margin classifiers", In D. Haussler (Ed.), Annual ACM Workshop on COLT, Pittsburgh, PA, 144(52), ACM Press,1992.

[14] Kecman,"Support vector machines: Theory and applications", Springer, 2005.

[15] Cortes C., Vapnik V.N., "Support-vector networks", Machine Learning, vol. 20, pp. 273-297,1995.

[16] NelloCristianini and John Shawe-Taylor, "An Introduction to Support Vector Machines", Cambridge University Press, 2000.

[17] Bishop C., Pattern Recognition and Machine Learning, Springer, 2006.

[18] Shawe-Taylor J., and Cristianini N., "Kernel Methods for Pattern Analysis", Cambridge University 2004.

[19] Shawe-Taylor J., Bartlett P.L., Willianmson R.C., Anthony M., "Structural risk minimization over data-dependent hierarchies", IEEE Trans. Information Theory, vol. 44, no. 5, pp. 1926-1940, 1998.

[20] Maji S., Berg A.C., Malik J., "Classification using intersection kernel support vector machines is efficient", IEEE Conference on Computer Vision and Pattern Recognition, June pp.1–8, 2008.

[21] B. Haasdonk, "Feature space interpretation of SVMs with indefinite kernels", IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 27, no. 4, 2002.

[22] E.A.Zanaty, Ashraf Afifi, Rania El-Khateeb, "Improving the accuracy of support vector machines via a new kernel functions", International journal of intelligent Computing, sciences, pp. 55-67, 2009.

[23] Ying Tan, "A Support vector machine with a hybrid kernel and minimal Vapnik chervonenk is dimension",. IEEE Transactions on knowledge and data engineering, vol. 16, no. 4, APRIL,2004

[24] Mercer T. , "Functions of positive and negative type and their connection with the theory of integral equations", Philosophical Trans. of the Royal Soc. of London, Series A, pp. 415-446,1909.

[25] WilliamsonR.C., SmolaA.J., and Scholkopf B., "Generalization performance of regularization networks and support vector machines via entropy numbers of compact operators", Technical Report 19, Neuro COLT, 1998.

[26] Scholkopf B., A.J. Smola, "Learning with kernels: support vector machines, regularization, optimization", and Beyond, MIT Press, 2001.

[27] Vapnik V.N., "Statistical Learning Theory", Wiley-Interscience, New York, 1998.

[28] http://archive.ics.uci.edu/ml/datasets.html.

AUTHOR'S PROFILE

Ashraf Afifi is an assistance professor, faculty ofcomputers and information technology, Taif University, Saudi Arabia. He received his MSC Degree in digital communication in 1995 from zagazig University, Egypt. He completed his Ph. D. studies in 2002 from zagzig University,Egypt. His research interests aredigital communication, and image segmentation. In theseareas he has published several technical papers in refereed international journals or Conference proceedings.

# Improving Classification Accuracy of Heart Sound Signals Using Hierarchical MLP Network

Mohd Zubir Suboh, Md. Yid M.S., Muhyi Yaakob
Medical Engineering Technology Section
Universiti Kuala Lumpur
Kuala Lumpur, Malaysia

Mohd Shaiful Aziz Rashid Ali
School of Computer and Communication Engineering
Universiti Malaysia Perlis
Perlis, Malaysia

*Abstract*—**Classification of heart sound signals to normal or their classes of disease are very important in screening and diagnosis system since various applications and devices that fulfilling this purpose are rapidly design and developed these days. This paper states and alternative method in improving classification accuracy of heart sound signals. Standard and improvised Multi-Layer Perceptron (MLP) network in hierarchical form were used to obtain the best classification results. Two data sets of normal and four abnormal heart sound signals from heart valve diseases were used to train and test the MLP networks. It is found that hierarchical MLP network could significantly increase the classification accuracy to 100% compared to standard MLP network with accuracy of 85.71% only.**

*Keyword—Hierarchical MLP network; Multi-layer Peceptron Network; heart sound signals*

## I. INTRODUCTION

Heart auscultation and diagnosis are quite complicated, depending not only on the heart sound but also on other factors such as the acquisition method and patient condition [1]. In the last two decades, many research activities were conducted concerning automated and semi-automated heart sound diagnosis. The researches were concentrated at three major tasks which are segmentation of the heart sound, feature extraction and classification of heart sound signals using artificial intelligence system. The classification algorithms were mainly based on Discriminant analysis [2], k-Nearest Neighbour [3], Bayesian networks [4], Neural Networks (including radial basis function, multiplayer perceptron, self-organizing map, probabilistic neural networks) [5-7] and rule-based methods [8].

Artificial Neural Network (ANN) is one of the popular method used in classifying the heart sound signal. Sinha et al. [9] and Ari et al. [10] have done several researches and proved that ANN can classify a few type of heart valve diseases with good accuracy. ANN is a mathematical model that is inspired by the biological neural networks in human brain. In general, neural networks provide good solutions to problems with the following features.

- The problem related to noisy data. ANN has been proved to be robust for noise data applications [5-7, 11].

- A good and fast processing may be required instead the most perfect solution.

- There are no simple rules for solving the problem but only a set of sample solutions. The network can be `trained' on these so that it produces good responses to similar new cases.

Regardless the methods are, classification accuracy is the most important since poor performance given by the system or classifier to recognize the significant cardiac lesions might leads to adverse outcomes to the patient as well as unnecessary costs for inappropriate and even potentially hazardous laboratory test. Hence this study is done to provide an alternative method in improving classification accuracy of normal and abnormal heart sound signals from heart valve disease by using standard and improvised hierarchical Multi-Layer Perceptron (MLP) networks. Normal (N) and four abnormal heart sound signals of Mitral Regurgitation (MR), Mitral Stenosis (MS), Aortic Regurgitation (AR) and Aortic Stenosis (AS) from heart valve disease are used as the data in the classification process. Stenosis and Regurgitation problems are chosen in this study since they always affect the heart valves. There are cases where one or more valves affected by both problems. So there will be multiple types of heart valve disease that make the classification of the heart sound signal is very difficult [12]. That is why this study is limited with two heart valves (Mitral and Aortic valves) which having regurgitation or stenosis problems.

## II. HEART SOUND SIGNAL

The data used in this study is solely from heart sound signal, no ECG or other biomedical signals are involved. The heart sound signals are taken from three sources, which are heart sound manipulator software, recorded signals from Hospital Tuanku Fauziah (HTF) and signal that is available in the internet. The heart sound manipulator software was used to test the reliability of the method used in feature extraction. The recorded signals from HTF are not enough for the study since only 6 subjects with specific required diseases are obtained. Hence, the data that are available in the internet are collected to be used for the analysis in this study. Even so, the data collected from the internet are obtained only from trusted medical and electronic stethoscope websites. These data were verified first to ensure that no artificial signals involved.

For classification purpose, all heart sound signals from the three sources are divided into two data sets. The first data set are the simulated heart sound signals from the heart sound manipulator software. There are seven subjects (recording) taken and each subject contributes 30 samples after manual

segmentation process is done (two heart sound cycles of each sample). The second data set contains 39 samples and 646 samples. It is collected from the real recording heart sound signals from HTF and collected heart sounds signal from the trusted websites. Since the duration of each signals of the second data sets are different, it contributes different numbers of heart sound samples, which also manually segmented by two cycles of heart sound signal. The summary of subjects and samples taken from all sources is shown in Table I.

TABLE I.  GROUP OF DATA SET FOR HEART SOUND SIGNALS CLASSIFICATION

| Heart Sound Category | | First Data Set | | Second Data Set | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | Subjects | Samples | Subjects | | Samples | |
| | | | | Patient | Internet | Patient | Internet |
| Normal | | 1 | 30 | 8 | 0 | 217 | 0 |
| Abnormal | Aortic Regurgitation | 1 | 30 | 0 | 7 | 0 | 103 |
| | Aortic Stenosis | 1 | 30 | 2 | 7 | 38 | 72 |
| | Mitral Regurgitation | 1 | 30 | 3 | 4 | 59 | 51 |
| | Mitral Stenosis | 3 | 90 | 1 | 7 | 23 | 83 |
| Total | | 7 | 210 | 14 | 25 | 337 | 309 |
| | | | | 39 | | 646 | |

Segmentation on heart sound signal needs to be done to obtained uniform samples to be used for the analysis. The reason is that the collected signal from three different sources varies in term of duration, number of cycles and so on. Segmentation is done based on cycles since other features such as time and frequency are not consistent. Two cycles of the recorded heart sound signal are taken as a sample. Two cycles are taken as a sample because the two-cycle-sample offered different feature of S1, S2, systolic and diastolic components for each cycle. This will make the analysis more accurate.

## III. FEATURE EXTRACTION PROCESS

Frequency analysis method is applied in this study in order to obtain the frequency features of the heart sound samples. Cross-correlation function is used in understanding the strength of a linear relationship between two variables in this case normal and abnormal heart sound sample. Correlation analysis method is selected because it plays a major role in statistical signal processing. The cross-correlation function is used extensively in pattern recognition and signal detection [13].

In this study, cross-correlation was used to find the frequency (power spectrum) relation of a reference sample with a normal or abnormal heart sound sample. The reference sample is an average of 100 normal heart sound samples (also in power spectrum form). The power spectrum gives a plot of the portion of a signal's power (energy per unit time) falling within given frequency bins. The most common way of generating power spectrum is by using DFT function, but other techniques such as the Maximum Entropy Method can also be used (Weisstein & Eric, 2010). DFT is used in order to find the frequency components of a signal buried in a noisy time domain signal.

The procedure starts with converting the two-cycle-sample (time domain) into frequency domain using Discrete Fourier Transform (DFT). The equation to calculate the DFT is shown in (1) where $x$ is the heart sound signal in time domain while $N$ is the length of the $x$ signal. Then power spectrum of the signal is determined using complex conjugate (*conj*) function as in (2), based on the DFT signal.

$$DFT, X(k) = \frac{1}{N}\sum_{n=0}^{N-1} x(j)e^{-j\left(\frac{2\pi kn}{N}\right)} \quad (1)$$

$$Power\ Spectrum, P_X = \frac{(X) \times conj(X))}{N} \quad (2)$$

Knowing that the heart sound signals frequency components are at below 600 Hz [14,15], the first 600 components of the power spectrum are considered as the most important signals to be analyzed. The first 600 components are the energy of the first 600 Hz frequency components of the heart sound signals. The process of taking out the 600 significant values can be viewed as a filtering technique (in frequency domain) where undesired frequencies components can be simply cut-off. Power spectrum of the reference sample ($P_r$), also 600 values, is then cross-correlated with power spectrum of a test sample ($P_X$) using (3) where * denotes complex conjugation and $k$ is a variable to complete the cross-correlation [16]. As for the reference sample, it is obtained by averaging 100 sets of normal heart sound samples. Fig. 1 shows the power spectrum plot of the reference heart sound sample.

$$Cross\ Correlation, r_{P_r P_y}(k)$$
$$= \frac{1}{N+1}\sum_{n=0}^{N} P_r(n) \times P_X(n-k),$$
$$= \frac{1}{N+1}\sum_{n=0}^{N} P_X(n) \times P_r(n-k), \quad k = 0, \pm1, \pm2,.. \quad (3)$$

Fig. 1.    Power spectrum of reference heart sound sample

Cross-correlation between these signals will give 1200 points of plot pattern which depends on the type of the signal (normal and abnormal). For a normal testing sample, the cross-correlation plot should be symmetrical or almost symmetrical since it is correlated with reference sample that is also normal heart sound sample. For abnormal sample cases, even there are times when the correlation plot has the symmetric pattern, but the position of the plot are different. The plot patterns (slopes, curves and peaks with their position) were used in this study as the features to classify the samples into their category of heart sound. Fig. 2 shows examples of cross-correlation plot pattern for each category of normal and abnormal heart sound. The 1200 points for a sample is too many to be processed and thus it is averaged to 50 points only. The points cannot be reduced less than 50 since it will significantly interrupt the plot pattern. This 50 points data (a heart sound sample for classification) will be used to train and test the MLP classifiers.



Fig. 2.    Cross-correlation plot for five category of heart sound

## IV.    CLASSIFICATION PROCESS

Two sets of feed forward Multi-Layer Perceptron (MLP) network structures are used in the classification process to obtain the best classification accuracy. The first MLP network classifies 5 categories of normal and four type of heart valve disease. Fig. 3 shows the network structure. The network used 50 input neurons for 50 cross-correlation values of a sample and 5 output neurons for 5 categories of heart sound signals. Using an output neuron for a category will make the training process easier. This is because the network only has to

produce output valued 1 to the corresponding neuron and 0 to the others.



Fig. 3.    Classification of 5 categories of heart sound signal using usual MLP network.

The second MLP network that was used to classify the five categories of heart sound signal is hierarchical MLP network. Two MLP networks are used to classify normal and abnormal heart sound signal as well as classify the abnormal signal. This method reduces the complexity of training process and could increase the classification accuracy of the system. Fig. 4 shows the network structures. The first network in this structure has two output neurons for normal and abnormal category of heart sound. Abnormal category is the combination of four abnormal signals, which are AR, AS, MR and MS. The second network is used to classify the four categories of heart sound signals if the first network has abnormal output. Both networks used the same 50 values of the cross-correlation plot as the input.



Fig. 4.    Classification of 5 categories of heart sound signal using hierarchical MLP network.

As for the number of neurons in hidden layer, each network is tested with different hidden neurons from a single neuron to a maximum of 30 hidden neurons. 30 hidden neurons were set to be a maximum considering the time consumed to complete the process is acceptable. The initial weights and biases were set using random function in MATLAB software. Levenberg-Marquardt training (TRAINLM) algorithm was used for the training since it is the fastest convergent method available in the MATLAB Neural Network toolbox. Log-sigmoid transfer function (logsig) was the transfer function used in the network. The function logsig

generates outputs between 0 and 1 as the neuron's input goes from negative to positives infinity.

## V. RESULTS AND DISSCUSSION

Accuracy of classification on both networks is discussed in this section. The classification accuracy is calculated based on different random values (the initial values for weights and biases) with 30 different number of hidden neurons (from 1 to 30) to find the best classification accuracy. The proposed classifier was validated using two different data sets of heart sound signals. The first set is a simulated heart sound data and the second set is a combination of heart sound data from real patients and internet as in Table 1.0 before. The samples of both data sets are approximately divided into 60% and 40% for neural network training and testing purposes respectively. For the second data set, subjects and its samples were also divided into the ratio of 60% and 40%, which means testing is done using different samples from completely different subjects.

### A. Classification Accuracy of Standard MLP Network

Normal (N) and the other four categories of heart valve disease (AR, AS, MR, MS) sounds are classified using a standard MLP network. Training parameters for goal and gradient are set to $1 \times 10^{-24}$ and 0, respectively. Number of epochs is maximized until 1000. The training will stop if it achieved the target values for goal or gradient or else after the epoch reached 1000. Other training parameters were set to default values assigned by MATLAB. The classification accuracy is described based on 10 trials, which means 10 times of training testing using different random values, at 30 different hidden neurons. The 10 sets random values are set using *seed* 1 to 10 of the random generator provided by the neural network toolbox in MATLAB software. This random is used to set to give 10 different initial values for weight and biases in the neural network structure. The initial weight and bias values can affect the training and testing accuracy of the network. Therefore, the network is trained ten times to obtained the best trained network with the highest classification accuracy.

The best trained network for each number of hidden neuron (from 1 to 30) after 10 trials will be selected so that the results of training and testing accuracies, number of epochs as well as the MSE can be compared. The selection is made based on the performance of network training, training accuracy and finally testing accuracy. The number of input neurons is 50 while the number of output neurons is five for five categories of heart sound. The samples used for training and testing are shown in Table II.

TABLE II. NUMBER OF SAMPLES FOR THE STANDARD MLP TRAINING AND TESTING

| Heart Sound Category | First Data Set | | Second Data Set | |
|---|---|---|---|---|
| | Training Samples | Testing Samples | Training Samples | Testing Samples |
| Normal | 18 | 12 | 130 | 87 |
| Aortic Regurgitation | 18 | 12 | 63 | 40 |
| Aortic Stenosis | 54 | 36 | 65 | 45 |
| Mitral Regurgitation | 18 | 12 | 65 | 45 |
| Mitral Stenosis | 18 | 12 | 62 | 44 |
| Total | 126 | 84 | 385 | 261 |
| | 210 | | 646 | |

For the first data set, 100% classification accuracy was easily obtained at any different numbers of hidden neuron after trials or ten times of testing except for 1 and 2 hidden neurons networks. These results are shown in Table III. This table shows the average testing accuracy with the best trained network for each hidden neuron after ten times of training and testing. The highest average testing accuracy is 95.713% at 22 hidden neurons. This show that 22 hidden neurons is the best number of hidden neurons because it can correctly classify the heart sound signals even the initial values for weights and biases are different.

The classification accuracy using the same method for second data set is given in Table IV. Most of the networks with different number of hidden neurons show a good training accuracy which is over 96% accept for the first network of 1 hidden neuron. The best classification accuracy obtained is only 86.12% even after 10 trials at each hidden neuron were tested. The best accuracy obtained was at 15 hidden neurons with MSE $9.99 \times 10^{-25}$ and epoch of 164. However, the 15 hidden neurons network have too low average accuracy of 67.67% after 10 trials were made. So another network need to be chose. It is found that the network is best trained at 26 hidden neurons where the average accuracy after 10 trials is 80.08% and the testing accuracy is 85.71%. The average testing accuracy of 80.08% is selected because it is an acceptable value after 10 times of testing using 10 different initial weight and bias values. The classification accuracy of 85.71% at 26 hidden neurons is also not too different form the highest one (86.12%). The classification accuracy of second data set significantly dropped compared to the accuracy obtained by using the first data set. This is because samples from the second data set are from 39 subjects compared to the first data set, 7 subjects only. Hence the second method that uses Hierarchical MLP network is used to improve the classification accuracy of the second data set.

TABLE III. CLASSIFICATION ACCURACY OF 5 CATEGORIES OF HEART SOUND SIGNALS AFTER 10 TRIALS AT 30 DIFFERENT HIDDEN NEURONS USING THE FIRST METHOD ON FIRST DATA SET

| Hidden Neuron | Average Accuracy after 10 Trials (%) | The Best Trained Network for Each Hidden Neuron | | | |
| --- | --- | --- | --- | --- | --- |
| | | No. of Epoch | Mean Square Error (MSE) | Training Accuracy (%) | Testing Accuracy (%) |
| 1 | 33.572 | 1000 | 0.067336 | 57.14 | 53.57 |
| 2 | 54.168 | 1000 | $5.77 \times 10^{-24}$ | 100 | 97.62 |
| 3 | 88.929 | 493 | $9.99 \times 10^{-25}$ | 100 | 100 |
| 4 | 76.905 | 431 | $9.99 \times 10^{-25}$ | 100 | 100 |
| 5 | 87.617 | 378 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 6 | 88.69 | 340 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 7 | 87.142 | 294 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 8 | 82.857 | 246 | $9.96 \times 10^{-25}$ | 100 | 100 |
| 9 | 71.31 | 243 | $9.96 \times 10^{-25}$ | 100 | 100 |
| 10 | 89.762 | 334 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 11 | 82.738 | 196 | $9.98 \times 10^{-25}$ | 100 | 100 |
| 12 | 71.428 | 178 | $1 \times 10^{-24}$ | 100 | 100 |
| 13 | 84.642 | 191 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 14 | 82.619 | 259 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 15 | 80 | 213 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 16 | 91.428 | 157 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 17 | 90 | 186 | $9.95 \times 10^{-25}$ | 100 | 100 |
| 18 | 68.809 | 200 | $9.95 \times 10^{-25}$ | 100 | 100 |
| 19 | 85.594 | 205 | $9.93 \times 10^{-25}$ | 100 | 100 |
| 20 | 94.285 | 160 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 21 | 68.572 | 167 | $9.93 \times 10^{-25}$ | 100 | 100 |
| **22** | **95.713** | **230** | $\mathbf{9.91 \times 10^{-25}}$ | **100** | **100** |
| 23 | 72.856 | 152 | $9.96 \times 10^{-25}$ | 100 | 100 |
| 24 | 84.166 | 202 | $9.92 \times 10^{-25}$ | 100 | 100 |
| 25 | 92.619 | 352 | $9.91 \times 10^{-25}$ | 100 | 100 |
| 26 | 85.714 | 190 | $9.93 \times 10^{-25}$ | 100 | 100 |
| 27 | 81.429 | 234 | $9.96 \times 10^{-25}$ | 100 | 100 |
| 28 | 82.738 | 168 | $9.93 \times 10^{-25}$ | 100 | 100 |
| 29 | 88.571 | 142 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 30 | 82.857 | 141 | $9.91 \times 10^{-25}$ | 100 | 100 |

TABLE IV. CLASSIFICATION ACCURACY OF 5 CATEGORIES OF HEART SOUND SIGNALS AFTER 10 TRIALS AT 30 DIFFERENT HIDDEN NEURONS USING THE FIRST METHOD ON SECOND DATA SET

| Hidden Neuron | Average Accuracy after 10 Trials (%) | The Best Trained Network for Each Hidden Neuron | | | |
| --- | --- | --- | --- | --- | --- |
| | | No. of Epoch | Mean Square Error (MSE) | Training Accuracy (%) | Testing Accuracy (%) |
| 1 | 22.53 | 1000 | 0.0929123 | 45.39 | 48.16 |
| 2 | 28.776 | 1000 | 0.00299252 | 98.5 | 53.47 |
| 3 | 53.306 | 416 | $9.99 \times 10^{-25}$ | 100 | 77.55 |
| 4 | 50.939 | 1000 | 0.00798005 | 96.01 | 77.96 |
| 5 | 67.715 | 1000 | 0.000498753 | 99.75 | 80.82 |
| 6 | 69.307 | 1000 | 0.000498753 | 100 | 81.22 |
| 7 | 62.407 | 227 | $9.98 \times 10^{-25}$ | 100 | 76.73 |
| 8 | 64.572 | 272 | $9.97 \times 10^{-25}$ | 100 | 78.37 |
| 9 | 71.835 | 238 | $9.98 \times 10^{-025}$ | 100 | 80.41 |
| 10 | 71.429 | 201 | $9.94 \times 10^{-25}$ | 100 | 84.08 |
| 11 | 66.449 | 262 | $9.94 \times 10^{-25}$ | 100 | 82.86 |
| 12 | 50.754 | 137 | $9.96 \times 10^{-25}$ | 100 | 81.63 |
| 13 | 67.06 | 190 | $9.97 \times 10^{-25}$ | 100 | 81.63 |
| 14 | 64.407 | 164 | $9.99 \times 10^{-25}$ | 100 | 83.27 |
| **15** | **67.672** | **164** | $\mathbf{9.99 \times 10^{-25}}$ | **100** | **86.12** |
| 16 | 65.102 | 158 | $9.90 \times 10^{-25}$ | 100 | 84.08 |
| 17 | 80.898 | 184 | $9.95 \times 10^{-25}$ | 100 | 83.27 |
| 18 | 70.408 | 140 | $9.98 \times 10^{-25}$ | 100 | 84.49 |
| 19 | 67.426 | 169 | $9.89 \times 10^{-25}$ | 100 | 82.04 |

| 20 | 67.429 | 154 | $1x10^{-24}$ | 100 | 84.9 |
|----|--------|-----|--------------|-----|------|
| 21 | 70.693 | 153 | $9.90x10^{-25}$ | 100 | 82.04 |
| 22 | 63.797 | 136 | $9.90x10^{-25}$ | 100 | 82.86 |
| 23 | 67.673 | 207 | $9.91x10^{-25}$ | 100 | 83.67 |
| 24 | 65.183 | 156 | $9.94x10^{-25}$ | 100 | 84.08 |
| 25 | 70.941 | 171 | $1x10^{-24}$ | 100 | 84.49 |
| **26** | **80.08** | **196** | **$9.91x10^{-25}$** | **100** | **85.71** |
| 27 | 80.286 | 167 | $9.96x10^{-25}$ | 100 | 84.08 |
| 28 | 73.672 | 200 | $9.96x10^{-25}$ | 100 | 85.71 |
| 29 | 67.511 | 132 | $9.89x10^{-25}$ | 100 | 84.08 |
| 30 | 75.634 | 176 | $9.95x10^{-25}$ | 100 | 85.31 |

### B. Classification Accuracy of Hierarchical MLP Network (First Network)

A total of 646 of normal and abnormal heart sound samples from 39 subjects of the second data set are used in this study to obtain the classification accuracy of normal and abnormal heart sound signals. This data is exactly the same data used in the classification of the first method only the output or heart sound categories are different. The samples in each data set are divided manually about 60% for neural network training and 40% for neural network testing. The details about the samples are shown in Table V. The classification result is shown in Table VI.

From the results shown in Table 6, training accuracy of all selected networks at each number of hidden neurons are 100%

where each network were trained and achieved the goal of $1x10^{-24}$ with maximum epoch of 523. 100% of testing accuracy of normal and abnormal classification was achieved at several numbers of hidden neurons. The other testing accuracies were exceeded 94% except for 1 hidden neuron network, 87.35%. This shows that the network have successfully learnt and classified the heart sound signals very well. The best average testing accuracy obtained is at 25 hidden neurons network with the accuracy of 94.654%. However, the best testing accuracy at this network is only 98.78%. 100% of testing accuracy had become priority in selecting the best network only if the average accuracy is acceptable values, over 80%. Hence the best network is at 22 hidden neurons because it produced 100% accuracy for training and testing with an average accuracy of 83.184%.

TABLE V.    NUMBER OF SAMPLES FOR THE FIRST NETWORK TRAINING AND TESTING

| Heart Sound Category | Training Samples | Testing Samples |
|---|---|---|
| Normal | 130 | 87 |
| Abnormal | 255 | 174 |
| **Total** | **385** | **261** |
| | **646** | |

TABLE VI.    CLASSIFICATION ACCURACY OF NORMAL AND ABNORMAL HEART SOUND SIGNALS AFTER 10 TRIALS AT 30 DIFFERENT HIDDEN NEURONS USING THE FIRST NETWORK OF THE SECOND METHOD

| Hidden Neuron | Average Accuracy after 10 Trials (%) | The Best Trained Network for Each Hidden Neuron | | | |
|---|---|---|---|---|---|
| | | No. of Epoch | Mean Square Error (MSE) | Training Accuracy (%) | Testing Accuracy (%) |
| 1 | 57.837 | 376 | $1x10^{-24}$ | 100 | 87.35 |
| 2 | 69.143 | 145 | $9.99x10^{-25}$ | 100 | 94.29 |
| 3 | 70.001 | 394 | $1x10^{-24}$ | 100 | 95.1 |
| 4 | 66.083 | 232 | $9.96x10^{-25}$ | 100 | 98.78 |
| 5 | 76.122 | 175 | $9.93x10^{-25}$ | 100 | 98.78 |
| 6 | 81.185 | 146 | $9.99x10^{-25}$ | 100 | 100 |
| 7 | 81.265 | 155 | $9.97x10^{-25}$ | 100 | 97.55 |
| 8 | 82.163 | 186 | $9.92x10^{-25}$ | 100 | 94.69 |
| 9 | 71.836 | 178 | $9.94x10^{-25}$ | 100 | 100 |
| 10 | 61.185 | 183 | $9.94x10^{-25}$ | 100 | 98.78 |
| 11 | 68.816 | 227 | $1x10^{-24}$ | 100 | 97.55 |
| 12 | 61.675 | 266 | $9.96x10^{-25}$ | 100 | 94.29 |
| 13 | 59.674 | 216 | $9.96x10^{-25}$ | 100 | 97.55 |
| 14 | 76.165 | 327 | $9.90x10^{-25}$ | 100 | 98.78 |
| 15 | 67.225 | 223 | $9.98x10^{-25}$ | 100 | 97.55 |
| 16 | 73.633 | 365 | $9.92x10^{-25}$ | 100 | 97.55 |
| 17 | 80.57 | 178 | $9.93x10^{-25}$ | 100 | 99.18 |
| 18 | 79.837 | 370 | $9.93x10^{-25}$ | 100 | 96.33 |
| 19 | 71.96 | 108 | $9.89x10^{-25}$ | 100 | 98.78 |

| 20 | 82.654 | 168 | $9.95 \times 10^{-25}$ | 100 | 94.69 |
|---|---|---|---|---|---|
| 21 | 76.369 | 201 | $9.99 \times 10^{-25}$ | 100 | 100 |
| **22** | **83.184** | **162** | **$9.93 \times 10^{-25}$** | **100** | **100** |
| 23 | 80.94 | 236 | $9.93 \times 10^{-25}$ | 100 | 98.78 |
| 24 | 88.368 | 329 | $9.91 \times 10^{-25}$ | 100 | 98.78 |
| **25** | **94.654** | **523** | **$9.94 \times 10^{-25}$** | **100** | **98.78** |
| 26 | 81.633 | 165 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 27 | 81.96 | 113 | $9.95 \times 10^{-25}$ | 100 | 100 |
| 28 | 73.837 | 158 | $9.90 \times 10^{-25}$ | 100 | 98.78 |
| 29 | 74.286 | 308 | $9.98 \times 10^{-25}$ | 100 | 97.55 |
| 30 | 62.083 | 179 | $9.94 \times 10^{-25}$ | 100 | 97.96 |

## C. *Classification Accuracy of Hierarchical MLP Network (Second Network)*

The second MLP network with four output neurons is used to classify the four categories of heart valve disease. The network is trained and tested using the second data set as shown in Table VII while the accuracy of the network is shown in Table VIII. The results shows that the second network of hierarchical method can achieve 100% testing accuracy at many numbers of hidden neurons. All networks except the 1 hidden neuron network have successfully trained with training accuracy of 100%. The best average testing accuracy is 96.168% at 27 hidden neurons network. After ten times of testing, the 27 hidden neurons network have perfectly categorized 174 abnormal heart sound samples to their classes with MSE $1 \times 10^{-24}$ and 99 epochs. The combination of first and second network of hierarchical method had produced 100% accuracy, which is better than the first classification method with only 85.71% accuracy. This results show that the division of classification category using hierarchical technique of MLP network had improved the classification accuracy because the difficulties or complexity of classification had been reduced.

TABLE VII.    NUMBER OF SAMPLES FOR THE SECOND NETWORK TRAINING AND TESTING

| Heart Sound Category | Second Data Set | |
|---|---|---|
| | Training Samples | Testing Samples |
| Aortic Regurgitation | 63 | 40 |
| Aortic Stenosis | 65 | 45 |
| Mitral Regurgitation | 65 | 45 |
| Mitral Stenosis | 62 | 44 |
| **Total** | **255** | **174** |
| | **429** | |

TABLE VIII.    CLASSIFICATION ACCURACY OF 4 CATEGORIES OF HEART SOUND SIGNALS AFTER 10 TRIALS AT 30 DIFFERENT HIDDEN NEURONS USING THE SECOND NETWORK OF THE SECOND METHOD

| Hidden Neuron | Average Accuracy after 10 Trials (%) | The Best Trained Network for Each Hidden Neuron | | | |
|---|---|---|---|---|---|
| | | No. of Epoch | Mean Square Error (MSE) | Training Accuracy (%) | Testing Accuracy (%) |
| 1 | 30.633 | 1000 | 0.0811067 | 97.05 | 97.47 |
| 2 | 46.013 | 245 | $9.98 \times 10^{-25}$ | 100 | 53.80 |
| 3 | 50.063 | 304 | $9.98 \times 10^{-25}$ | 100 | 84.18 |
| 4 | 58.101 | 261 | $1 \times 10^{-24}$ | 100 | 93.04 |
| 5 | 64.557 | 248 | $9.96 \times 10^{-25}$ | 100 | 95.57 |
| 6 | 62.468 | 282 | $9.98 \times 10^{-25}$ | 100 | 91.14 |
| 7 | 61.772 | 248 | $9.98 \times 10^{-25}$ | 100 | 91.14 |
| 8 | 72.848 | 214 | $1 \times 10^{-24}$ | 100 | 86.08 |
| 9 | 71.076 | 214 | $9.96 \times 10^{-25}$ | 100 | 93.04 |
| 10 | 73.291 | 177 | $9.95 \times 10^{-25}$ | 100 | 95.57 |
| 11 | 61.139 | 138 | $9.98 \times 10^{-25}$ | 100 | 95.57 |
| 12 | 63.418 | 109 | $1 \times 10^{-24}$ | 100 | 95.57 |
| 13 | 69.873 | 163 | $9.99 \times 10^{-25}$ | 100 | 95.57 |
| 14 | 80.190 | 131 | $9.93 \times 10^{-25}$ | 100 | 95.57 |
| 15 | 77.405 | 117 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 16 | 87.342 | 150 | $9.92 \times 10^{-25}$ | 100 | 100 |
| 17 | 80.633 | 103 | $9.96 \times 10^{-25}$ | 100 | 100 |
| 18 | 73.734 | 157 | $9.97 \times 10^{-25}$ | 100 | 100 |
| 19 | 46.835 | 160 | $9.98 \times 10^{-25}$ | 100 | 100 |
| 20 | 63.038 | 106 | $9.96 \times 10^{-25}$ | 100 | 95.57 |
| 21 | 84.367 | 140 | $9.98 \times 10^{-25}$ | 100 | 100 |
| 22 | 66.835 | 99 | $9.88 \times 10^{-25}$ | 100 | 95.57 |

| 23 | 83.861 | 98 | $9.96 \times 10^{-25}$ | 100 | 100 |
|---|---|---|---|---|---|
| 24 | 68.038 | 146 | $9.91 \times 10^{-25}$ | 100 | 100 |
| 25 | 77.658 | 112 | $9.89 \times 10^{-25}$ | 100 | 97.47 |
| 26 | 82.215 | 139 | $9.90 \times 10^{-25}$ | 100 | 95.57 |
| **27** | **96.168** | **99** | **$1 \times 10^{-24}$** | **100** | **100** |
| 28 | 71.582 | 106 | $9.94 \times 10^{-25}$ | 100 | 100 |
| 29 | 73.165 | 86 | $9.84 \times 10^{-25}$ | 100 | 100 |
| 30 | 77.722 | 106 | $9.84 \times 10^{-25}$ | 100 | 100 |

## VI. CONCLUSION

This study has proved that classification of heart sound signal using standard MLP network can be increased using the hierarchical MLP network. Two data sets had been used. The first data set is the simulated data which can be easily classified by using the standard MLP network with 100% accuracy. However the standard MLP network can only classifies the heart sound signals up to 85.71% accuracy only when second data set is used (real data from patients). The accuracy is improved to 100% when hierarchical MLP network is used. The results show that the division of classification category using hierarchical technique of MLP network had improved the classification accuracy because the difficulties or complexity of classification had been reduced.

### REFERENCES

[1] A. C. Stasis, E. N. Loukis, S. A. Pavlopoulos and D. Koutsouris, "A decision tree – based method for the differential diagnosis of Aortic Stenosis from Mitral Regurgitation using heart sounds," BioMedical Engineering OnLine, 2004, 3:21.

[2] T. Leung, P. White, W. Collis, E. Brown and A. Salmon, "Analysing paediatric heart murmurs with discriminant analysis," Proceedings of the 19th Annual conference of the IEEE Engineering in Medicine and Biology Society, 1998, 1628-1631.

[3] L. G. Durand, H. Sabbah and P. Stein, "Comparison of spectral techniques for computer-assisted classification of spectra of heart sounds in patients with porcine bioprosthetic valves," Med Bio Eng Comput, 1997, 31(3), 229-36.

[4] L. G. Durand, and P. Pibarot, "Digital signal processing of the phonocardiogram: review of the most recent advancements," Critical Reviews in Biomedical Engineering, 1995, 23,3/4, 163-219.

[5] C. G. DeGroff, S. Bhatikar, J. Hertzberg, R. Shandas, L. Valdes-Cruz and R. L. Mahajan, "Artificial neural network-based method of screening heart murmurs in children," Circulation 2001, 103, 2711-2716.

[6] J. E. Hebden and J. N. Torry, "Neural network and conventional classifiers to distinguish between first and second heart sounds," IEE Colloquium (Digest), 1996, 3/1-3/6.

[7] T. Leung, P. White, W. Collis, E. Brown and A. Salmon, " Classification of heart sounds using time-frequency method and artificial neural networks," Proceedings of the 22nd Annual International Conference of the IEEE Engineering in Medicine and Biology Society, 2000, 2, 988-991.

[8] Z. Sharif, M. S. Zainal, A. Z. Sha'ameri and S. H. S. Salleh, "Analysis and classification of heart sounds and murmurs based on the instantaneous energy and frequency estimations," Proceedings IEEE 2, 2000, 130-134.

[9] R. K. Sinha, Y. Aggarwal and B. N. Das, "Backpropagation Artificial Neural Network Classifier to Detect Changes in Heart Sound due to Mitral Valve Regurgitation," J Med Syst, 31(2007), 205-209.

[10] S. Ari, P. Kumar and G. Saha, "On an algorithm for boundary estimation of commonly occurring heart valve diseases in time domain," Annual India Conference (INDICON), 2006, 1-6.

[11] Z. Sharif, M. S. Zainal, A. Z. Sha'ameri and S. H. S. Salleh, "The design of heart classification system," International Symposium on Signal Processing and its Applications (ISSPA), 2001, Kuala Lumpur, Malaysia.

[12] A. Voss, J. Herold, R. Schroeder, F. Nasticzky, A. Mix1i, P. Ullrich and T. and Huebner, "Diagnosis of aortic valve stenosis by correlation analysis of wavelet filtered heart sounds," Proceedings of the 25th Annual International Conference of the IEEE EMBS Cancun, 2003, Mexico, pp. 2873-2876.

[13] J. O. Smith III, Mathmathics of the Disrete Fourier Transform (DFT), with Audio Applications. (2nd. Edition). W3K Publishing, 2007.

[14] T. S. Yogeeswaran, D. D. N. B. Daya and K. D. I. Wasudeva, "Construction of a Low Cost Electronic Heart Sound Monitoring System. Proceedings of the Technical Sessions, 2008, 24, 72-77.

[15] T. Xin and T. Zhong, "Analysis and Decision of Heart Sounds via Arma Models," Measurement, 1987, 5(3), 102-106.

[16] T. Boss, Digital Signal and Image Processing. United States of America: John Wiley & Sons Inc., 2004.

# A Generic Framework for Automated Quality Assurance of Software Models –Implementation of an Abstract Syntax Tree

Darryl Owens and Dr Mark Anderson

Department of Computing

Edge Hill University

Ormskirk, Lancashire

*Abstract*—**Abstract Syntax Tree's (AST) are used in language tools, such as compilers, language translators and transformers as well as analysers; to remove syntax and are therefore an ideal construct for a language independent tool. AST's are also commonly used in static analysis. This increases the value of ASTs for use within a universal Quality Assurance (QA) tool. The Object Management Group (OMG) have outlined a Generic AST Meta-model (GASTM) which may be used to implement the internal representation (IR) for this tool. This paper discusses the implementation and modifications made to the previously published proposal, to use the Object Management Group developed Generic Abstract Syntax Tree Meta-model core-components as an internal representation for an automated quality assurance framework.**

*Keywords—software quality assurance; software testing; automated software engineering; programming language paradigms; language independence; abstract syntax tree; static analysis; dynamic analysis*

## I. INTRODUCTION

To ensure the reliability of output, it is imperative that Software Quality Assurance (QA) is adopted in the development and maintenance of scientific software systems [1]. The integration of such techniques can either be performed manually, which is labour intensive, or utilise automated toolkits [2] [3] which alleviate these problems. The automated toolkits are limited in the respect that they are language-, paradigm- or problem-specific. This paper proposes a framework that would address these limitations by introducing a taxonomy of generic techniques combined with a generic internal representation (IR) of languages. The framework also covers a range of different language paradigms. This paper also proposes a form of IR representation that could be used as an intermediary between QA techniques and source code.

When considering the broad range of programming paradigms the differences between the languages, such as the constructs and data types, need to be addressed. This paper focuses on addressing issues in procedural and object-oriented languages as these are the most widely adopted paradigms in the development of scientific software [4].

## II. ABSTRACT SYNTAX TREES

In order to address syntactical differences, Abstract Syntax Trees (AST) are adopted as 'a formal representation of the software syntactical structure' [5]. At a surface level, the underpinning constructs of many procedural languages appear similar, and removing syntax from these would make all ASTs analogous. However, the resulting ASTs produced following analysis of source code are based on a broad range of factors, such as the context-free grammar used to define the language syntax [6]. It is therefore highly likely that the generated ASTs for simple algorithms implemented in different programming languages can prove to be fundamentally different. These differences can become even more significant when addressing additional language features, such as data types.

## III. CURRENT APPLICATIONS OF ASTs

The primary usage of ASTs is to facilitate the implementation of compiler tools. For this purpose, ASTs are built from token streams after lexical analysis of source code [7]. However, the usage of ASTs now encompasses the implementation of many language related tools, such as interpreters, document generators and syntax-directed editors, etc. [8].

One such use that an AST can support is in duplicate code detection, whereby an AST designed to support data matching efficiency [8] requires only a pattern to be found. The significance of this is that only an initial node needs to be identified which is subsequently followed by a predetermined pattern of nodes. This technique is similar to that found in the plagiarism detection techniques, which makes use of code comparison, described by Cui et al [9]. Equally, code analysis techniques can be implemented using node counting. Removing code comments and disregarding layout metrics produces better comparison metrics from this technique compared to those collected from source code [10].

A major development of ASTs lies in language translation, which occurs by producing an AST for a specific language. The AST can then be parsed whilst introducing the syntax of the output language.

After the tree has been parsed, the resultant code should be complete and functional in the output language syntax [11] [12]. A working example of this technique is adopted by Mono as a working and functional example of the approach in action [13]. Mono is a framework that allows the use of the .NET platform upon other devices than just windows via the use of language translation.

## IV.    GENERATING ASTs

In order to generate ASTs, ANother Tool for Language Recognition (ANTLR) is a tool which uses a grammar input and can produce recognisers, compilers and translators [14]. Of significance to the project presented within this paper is that the ANTLR compilers can build ASTs from source code [14]. Utilising ANTLR, an example of the fundamental differences that can be generated in ASTs from simple algorithms is presented. In this example, a simple "Hello World" program written in Java and C#.



Fig. 1.    Java Hello World AST



Fig. 2.    C# Hello World AST

It can be seen from the ASTs depicted in Figure 1 and Figure 2 that there are fundamental differences between the

representations of a simple program in two similar object-oriented programming languages that are included in the ANTLR repository [15].

Clearly there are some similarities between the ASTs; for example the class node has the name, modifiers and body. However, the significant differences between the AST's is directly related to the grammar files.

## V.    LANGUAGE INDEPENDENCE

A key requirement for the successful implementation of the proposed analysis framework requires language independence to be implemented in order to separate any reliance on the QA procedures from the syntax and semantics of the source code programming language. The Object Management Group (OMG) has initiated a number of projects to investigate the development of generic ASTs. Broadly there have been two tiers adopted for the approach, the Abstract Syntax Tree Metamodel (ASTM) and the Knowledge Discovery Metamodel (KDM). The KDM is a standard to facilitate interoperability for exchange of data between tools that may be provided by different vendors [16]. The KDM complements the ASTM and both are designed to work together. The extent to which they do is questionable as the link between the ASTM and KDM can best be described as fuzzy [17]. However the KDM is less relevant in the development of the proposed framework, as the KDM focus on migration of software artifacts and not representation of language, so the focus is on ASTM and the Generic ASTM (GASTM) which is defined in the ASTM specification [18].

## VI.    TESTING / QUALITY ASSURANCE AND AUTOMATED APPROACHES

There is a wide range of toolkits developed for testing software [19] [20] [21] [22] [23] [24] and, broadly, these are targeted at the automation of testing to reduce workload required to test software applications. An initial survey of the available toolkits has revealed that most tools that apply QA techniques to multiple languages only do so on a small-scale. Generally this is in the range of $2 - 5$ languages [2], and also focused on languages which share a programming paradigm. It is also noted that the more generic toolkits with a broader coverage of programming paradigm instead have a restriction in terms of the areas of testing which are covered [2].

There are two types of analysis within QA; dynamic and static [25] [26]. Whilst both offer advantages, combining these techniques results in a broader impact as a result of QA [18] [19]. Static analysis facilitates an abstract view of a program and examination of source code without code execution [27], and also supports the identification of such potential issues as memory corruption errors, buffer overruns, out-of-bound array accesses, or null pointer de-references [28].

Dynamic analysis is the analysis of code as it is executing, and therefore extracting accurate values of variables under set circumstances is a key target [25]. This technique can be used to run functional, logical, interface and bottom-up tests amongst a range of supported testing [25]. The combination of static and dynamic analysis allows for a larger coverage of QA

techniques and a tool which implements both of these analysis types would be more comprehensive [29][30].

### VII. PROPOSAL FOR ABSTRACT SYNTAX TREE

A proposal for the use of the GASTM as a form of Internal Representation (IR) for automated quality assurance was theorized under the ideas that both static and dynamic analysis could be implemented upon this IR via the processes described by the flow diagrams in figure 3, 4 and 5 [31].



Fig. 3. System Data Flow

It was identified that by using the GASTM, static analysis could be possible via implementing tree walkers. In essence this would entail replacing source code analysers, and could implement automated quality assurance techniques such as metric and pattern matchers as well as allowing the GASTM to be converted to a control flow graph for data flow analysis.

Dynamic analysis however is a more complicated matter. By using a generic monitor class, nodes could be inserted into a program before conversion into a runnable language. These inserted nodes would call method in the generic monitor class allowing for information about data, properties or runtime information to be pulled out and recorded or analyses whilst the program is running.

### VIII. IMPLEMENTATION

LIQA (Language Independent Quality Assurance), is the implementation of the research discussed in this paper. Utilizing tools that have been previous developed by third parties, LIQA implements a middle layer to facilitate interaction using bespoke code and breaks down a subset of the Java language forming the GASTM representation of the source code.



Fig. 4. Static Analysis Data Flow



Fig. 5. Dynamic Analysis Data Flow

## A. LIQA Functionality

LIQA was designed for practical use within a lab setting. A GUI (Graphic User Interface) was developed to control the overall work flow as well as visual representation of the GASTM IR (Internal Representation) structure to allow quick assessment of correctness. This has been accompanied by a feature to create and load projects into the software.

## B. GASTM Representation

The diagram shown below (figure 7) is the graphical output from LIQA and is a sample of the function definition 'HelloWorld' as shown in figure 6.

```
public class HelloWorld {
    public static void print() {
        printout("Hello, World");
    }
}
```

Fig. 6.   HelloWorld Java

As was identified earlier, the IR can become very complicated from even a simple program. After the Java code has been parsed into the classes that represent the GASTM nodes, the object is then walked using a separate class to 'pretty print' the IR into XML which is then taken and placed in a SVG (Scalable Vector Graphics) format culminating in figure 7.

Several modifications have had to be made to the GASTM representation developed by Modisco [32], these modifications have been made for one of two reasons. The small change of implementing java.io.Serializable on the classes GASTMFactoryImpl, GASTMObjectImpl,

GASTMPackageImpl, GASTMSemanticObjectImpl, GASTMSourceObjectImpl and GASTMSyntaxObjectImpl, was to enable the IR to be saved as a binary file thus making it simple to implement a project based file system and allow users to save their work.

The other modifications listed below were implemented to better mirror the properties of some of the selected procedural languages. The ClassType and ClassTypeImpl were modified to include a link to the AccessKind class via the methods getAccessKind and setAccessKind. This was because in Java, C#, C++ and recent high-level languages allow programmers to assign classes with the access modifier i.e. public, private or protected. The FunctionMemberAttribute and FunctionMemberAttributeImpl have had the property IsStatic added, which is a boolean variable. The methods to modify the property setIsStatic and getISStatic have also been added. This was as Java, C#, and C++ allow programmers to assign functions with the static modifier.

## C. Tools used in development

LIQA utilizes several tools to achieve various tasks, these tasks (and therefore tools) are not all necessary however they make LIQA easier to use and simple to test for issues. The tools used as listed below:

- Modisco - GASTM Core Model [32]
- JavaCC  - Produced tokinizer for Java (Grammar from library) [33]
- XsdVi  - Used to generate a .svg file from .xsd [34]
- Batik  - Toolkit to visualize .svg file in JFrame [35]

The Modisco library has a Java representation of the GASTM core objects and therefore can be used instead of having to write the object in Java or another language. This links with the JavaCC tool which generated a Java tokenizer using a grammar located in the JavaCC library which is used by LIQA to convert the source code and generate the GASTM IR via the Modisco library. These two tools were required to make the production of LIQA a quicker and simple process. However the other tools are used to simplify the use of LIQA and simplify fault finding within the parsing and IR generating process. After the IR is generated it is then walked by LIQA and written to and .xsd file.



Fig. 7.   HelloWorld GASMT

XsdVi is then used to generate a .svg file from the .xsd, this is so a graphical representation of the IR is viewable. Following this, the Batik toolkit has been implemented into LIQA to allow the .svg file to be viewed in a java JForm utilizing the JSVGScrollpane and JSVGCanvas.

### D. Limitations

The limitations of LIQA are underpinned by a variety of contributing factors and are segmented to specific sections of the program. LIQA itself has the limitation of only being able to handle single file programs rather than full programs/projects built up over multiple files. This is due to time constraint. Although extending LIQA to handle full programs/projects would be a simple process, at this stage of the research it is not required as LIQA is only a 'proof of concept' for the larger framework.

The following limitations are to the Java language parser and are either due to time constraints for this research however will be part of future development and not being required to test the basic function of the IR, or are due to the GASTM not supporting the specific syntax. The following limitations are due to the GASTM core limitations, the program that is being analyzed cannot contain:

$$<= , >= , += , -= , /= , *=$$

The GASTM cannot handle these operators however considering the operators can be broken down i.e. 'x += 1' = 'x = x + 1' and 'x <= 2' = 'x < 3' it would be possible to integrate these into the IR. Due to time constraints, the lack of importance with these operators, as they can be replaced and effort it would take to code the conversion, they have not been included within the parser that generates the IR from the tokens.

The following limitations are due to be implemented in further developments of LIQA. However they are not necessary for testing the framework at this stage.

- Operators that are not implemented are '?' and '!'
- List types are not implemented i.e. 'List<String>'
- Re-type casting has not been implemented i.e. 'String str = (String) x;'
- The assignment of arrays via block statement has not been implemented i.e. 'int[] x = {3,2,1};'
- Inline if statements have not been implemented, if statements must have a block containment i.e. 'if (condition) statement;' is not supported and 'if (condition) {statement}' is supported.

### E. Analysis Test

After the initial implementation of the GASTM IR was finished, a proof of concept addition was made to LIQA, this was to implement a single form of static and dynamic analysis to test the proposal before a deeper analysis of quality assurance techniques takes place.

For dynamic analysis a simple profiler was developed, this required a tree walker which analyses the IR to find all the function definitions and the variable definitions in them allowing the user control over which methods and variables were monitored, after the user made their choice LIQA then inserts nodes to monitor the variable values wherever modified and counts method calls through the monitor class interface. The monitor class must be written in the original language as it may require language specific method calls itself in later development.

For static analysis a simple metric was written using a similar tree walker as the one for the profiler however this returned Logical Lines of Code value and is currently setup for further metrics to be included.

### F. Proposal Modifications

During the implementation small changes had to be made to the proposal as technological limitations arose, the only change that was significant is the formulation of a plausible and feasible way of running a GASTM IR that has been modified to perform dynamic analysis. This was achieved by running a conversion back into source code from the GASTM IR via a tree walker, this would have to be implemented with every language due to library and specific method calls that are unique to that language. A further implementation of LIQA could include a method and parameter mapping system which would also allow for language conversion.

The tree walker however provides a further form of quality assurance as some coding standards require code to be formatted is a specific way, as the tree walker generates the code a formatting system can be applied for easier maintenance.

A smaller modification is to the flow of data with regards to how static analysis is run, the DFD (figure 8) shows the initial stage of the IR being converted to a CFG, though not all static analysis techniques utilize a CFG the conversion must take place for those that do before the application of a static quality assurance techniques can be applied. The tree walker is utilized by all techniques not just the conversion to a CFG and would also be required before a technique can be applied.

### IX. CONCLUSION

This paper presented the implementation of an internal representation fit to allow both quality assurance via static and dynamic analysis and also to allow the representation of multiple languages. The major issues of this implementation are differences in languages and how to apply the analysis upon the final structure.

The ASTM is a standard prepared by OMG for language based tools and implements a set of core components that can represent a subset with many procedural and object oriented languages. It is therefore an obvious choice for this internal representation.

Fig. 8.   LIQA Data Flow

The implementation of the proposal within the tool LIQA represents a proof of concept showing that the GASTM is a suitable IR and that quality assurance techniques can be applied to this. However exactly what techniques can be applied is uncertain, it is certain that at least a simple level of static and dynamic analysis can be performed. Further work will demonstrate what techniques can be applied and these techniques will be derived from many tools currently used as industry standards.

### REFERENCES

[1]   Rosenberg, L. (2002) Software quality assurance engineering at NASA. *Aerospace Conference Proceedings, 2002. IEEE*. 5 pp. 5-2569 - 5-2575.

[2]   Owens, D. and Anderson, M. "A Generic Framework for Automated Quality Assurance of Software Models: Supporting Languages of Multiple Paradigms". In 5th International Conference on Computer Engineering and Technology (ICCET), Vancouver, Canada, 13-14 April

[3]   Collins, J., Farrimond, B., Anderson, M., Owens, D., Bayliss, D. and Gill, D. "Automated Quality Assurance Analysis: WRF – a case study". Accepted for 5th International Conference on Computer Engineering and Technology (ICCET), Vancouver, Canada, 13-14 April

[4]   Pickering, R. (2010). *Beginning F#*. Apress.

[5]   Newcomb, P. (2005, October). *Abstract Syntax Tree Metamodel Standard ASTM Tutorial 1.0*. Retrieved 2 5, 2013, from Object Managment Group: http://www.omg.org/news/meetings/workshops/ADM_2005_Proceedings_FINAL/T-3_Newcomb.pdf

[6]   Tripp, A. (2006, February 22). *Manual Tree Walking Is Better Than Tree Grammars.* Retrieved 2 5, 2013, from ANTLR v2: http://www.antlr2.org/article/1170602723163/treewalkers.html

[7]   Fischer, G., Lusiardi, J., &  Gudenberg, J. (2007, August 25-31). Abstract Syntax Trees – and their Role in Model Driven Software Development. *Software Engineering Advances, 2007. ICSEA 2007. International Conference on* , 38.

[8]   Van Den Brand, M.,  Moreau , P., &  Vinju, J. (2005). A generator of efficient strongly typed abstract syntax trees in Java. *EE Proceedings - Software Engineering 152, 2 (2005) 70--87* , 70-87.

[9]   Cui, B., Li, J., Guo, T., Wang, J.-X., & Ma, D. (2010). Code Comparison System based on Abstract Syntax Tree. *Broadband Network and Multimedia Technology (IC-BNMT), 2010 3rd IEEE International Conference on*, 668- 673 .

[10]  Fischer, G., Lusiardi, J., &  Gudenberg, J. (2007, August 25-31). Abstract Syntax Trees – and their Role in Model Driven Software Development. *Software Engineering Advances, 2007. ICSEA 2007. International Conference on* , 38.

[11]  Ichisugi, Y. (2003). *Patent No. 6516461.* United States.

[12]  ASM. (2012). *Model Driven Modernization* . Retrieved 2 5, 2013, from Automated Software Modernization: http://www.automatedsoftwaremodernization.com/component/content/article/3.html

[13]  OMG. (2012, July 19). *Catalog Of Omg Modernization Specifications.* Retrieved 2 5, 2013, from Object Management Group: http://www.omg.org/technology/documents/modernization_spec_catalog.htm

[14]  Parr, T. (n.d.). *ANTLR.* Retrieved 2 4, 2013, from ANother Tool for Language Recognition: http://www.antlr.org

[15]  *Grammar List.* (n.d.). Retrieved 2 5, 2013, from ANTLR v3: http://www.antlr3.org/grammar/list.html

[16]  Deltombe, G., & Goaer, O. L. (2012). Bridging KDM and ASTM for Model-Driven Software Modernization . *SEKE* , 517-524.

[17]  OMG. (2011, January). *OMG Architecture-driven modernization: Abstract Syntax Tree etamodel (ASTM).* Retrieved 2 5, 2013, from Object Managment Group: http://www.omg.org/spec/ASTM/1.0/PDF/

[18]  Salah, M., Mancoridis, S., Antoniol, G. & Di Penta, M. (2006) Scenario-driven dynamic analysis for comprehending large software systems. *Software Maintenance and Reengineering, 2006. CSMR 2006. Proceedings of the 10th European Conference on* pp. 80 - 90.

[19]  Fairley, R. (1978) Tutorial: Static Analysis and Dynamic Testing of Computer Software. *Computer*. 11(4) pp. 14-23.

[20]  Anywhere, A. (n.d.) *TestingAnywhere*.        HYPERLINK "http://www.automationanywhere.com/Testing/"

http://www.automationanywhere.com/Testing/  [accessed 09 November 2012].

[21]  Systems, Q. (n.d.) *Cantata - The Unit Testing Tool for C/C++*. HYPERLINK "http://www.qa-systems.com/cantata.html%20" http://www.qa-systems.com/cantata.html [accessed 09 November 2012].

[22]  Artho, C. et al. (2004) JNuke: Efficient Dynamic Analysis for Java. *Proc. CAV '04*.

[23]  MathWorks (1994) *Static Analysis with Polyspace Products*. HYPERLINK "http://www.mathworks.co.uk/products/polyspace/" http://www.mathworks.co.uk/products/polyspace/___ [accessed 09 November 2012].

[24]  SimCon (1995) *SimCon - Fortran Analysis, Engineering & Migration*. HYPERLINK "http://www.simconglobal.com/" http://www.simconglobal.com/ [accessed 09 November 2012].

[25]  Fairley, R. (1978) Tutorial: Static Analysis and Dynamic Testing of Computer Software. *Computer*. 11(4) pp. 14-23.

[26]  Austin, A. & Williams, L. (2011) One Technique is Not Enough: A Comparison of Vulnerability Discovery Techniques. *Empirical Software Engineering and Measurement (ESEM), 2011 International Symposium on* pp. 97-106.

[27]  Austin, A. & Williams, L. (2011) One Technique is Not Enough: A Comparison of Vulnerability Discovery Techniques. *Empirical Software Engineering and Measurement (ESEM), 2011 International Symposium on* pp. 97-106.

[28]  Bell, D. & Brat, P.G. (2008) Automated Software Verification & Validation: An Emerging Approach for Ground Operations. *Aerospace Conference, 2008 IEEE* pp. 1 - 8.

[29]  Harrison, K. (1999) Static Code Analysis on the C-130J Hercules Safety-Critical Software. *UK International Systems Safety Conferance*.

[30]  Wong, W.E. (2000) An Integrated Solution for Creating Dependable Software. *Computer Software and Applications Conference* pp. 269 - 270.

[31]  Owens, D. and Anderson, M. "A Generic Framework for Automated Quality Assurance of Software Models - Application of an Abstract Syntax Tree". Science and Information Conference 2013, Heathrow, London, 7-9 October 2013

[32]  Modisco. (n.d.). Modisco. Retrieved from Eclipse: http://www.eclipse.org/MoDisco/

[33]  JavaCC. (n.d.). Java Compiler Compiler tm (JavaCC tm) - The Java Parser Generator. Retrieved from JavaCC: http://javacc.java.net/

[34]  Slavětínský, V., & Kosek, J. (2013, March 22). XsdVi. Retrieved from Source Forge: http://sourceforge.net/projects/xsdvi/

[35]  Apache. (n.d.). The Apache™ Batik Project. Retrieved from The Apache™ XML Graphics Project: http://xmlgraphics.apache.org/batik/

# Computerized Kymograph for Muscle Contraction Measurement Using Ultrasonic Distance Sensor

Suhaeri[1], Vitri Tundjungsari [2]

YARSI University
Faculty of Information Technology (FTI)
Jl. Letjen Suprapto, Jakarta Pusat, Indonesia

*Abstract*—**Kymograph is a device to record the magnitude of physiological variables, such as: muscle contraction. However, we observe some lacks of the conventional kymographs, such as: result's visualisation and accuracy. Hence, we propose a computerized-kymograph which can automatically measure, record, and display graphical data of muscle contraction on the computer by using ultrasonic distance sensor. We develop hardware and software systems to support computerized kymograph and then test our device with live frog. The result shows that the device works well by displaying better visualization than the conventional kymograph.**

*Keywords—kymograph; computerized; distance sensor*

## I. Introduction

Kymograph is an aid for recording the magnitude as well as time course of a wide variety of physiological variables. The fundamental principle of converting time into distance has remained unchanged from the kymograph through polygraphs and oscilloscopes. We combine the fundamental works of conventional kymograph and a computerized device. A computerized device is defined as a machine that performs calculations and processes automatically. Hence, we define computerized kymograph as an automatic device that performs kymograph functionality, by recording and displaying the result on the computer.

We observe that the conventional kymograph at our laboratory having some lacks of visual appearance and accuracy. For example, the conventional kymograph did not demonstrate the value of muscle's contraction. Therefore we propose the computerized kymograph to perform better than the conventional one.

This paper presents an automated system to: (i) generate and display graphical data by computer, (ii) enhance the conventional kymographs result visibility for more accurate result. Section 2 discusses some related works about kymograph and its application. Section 3 explains how the conventional kymograph works and why we need to improve its performance. Section 4 describes how computerized kymograph is developed and tested. In section 4, we also explain the steps of the algorithm. Section 5 discusses the conclusion of our proposed device.

## II. Related Works

There are several works have been done related to kymograph for motion analysis and detection. Spatiotemporal method has been used in kymograph for motion analysis and detection, such as: using kymograph of video sequences to form a spatiotemporal image representation [1-4]. Kymographs have also been used to track image features over time; for example, the algorithm in [1, 3] used a variable-rate particle filter to enhance the accuracy of the extracted edges corresponding to the tip of micro-tubules from a kymograph-like image representation.

Mukherjee et al. [1] propose an automated method to profile the velocity patterns of small organelles (BDNF granules) being transported along a selected section of axon of a cultured neuron imaged by time-lapse fluorescence microscopy. The proposed method starts by generating a two-dimensional spatiotemporal map (kymograph) of the granule traffic along an axon segment, instead of directly detecting the granules as in conventional tracking. Author in [3] evaluated axonal transport by cross- and auto-correlation of kymograph columns; while author in [2] generated kymographs by hand and analysed them using the Radon transform to detect peak velocities of particles.

Our paper proposes a computerized-based kymograph in order to study the characteristics of certain physiological events such as muscle contractions, by providing more accurate recording device than a conventional kymograph. We argue that our proposed device can provide more accurate result of various physiological changes because it can directly record and display the skeletal muscle contractions (muscles can be removed from anesthetized frogs) on the computer. These muscles can be attached to recording systems and measured by ultrasonic distance sensor, stimulated by electrical shocks of varying strength, duration, and frequency. Recordings obtained from such procedures can be used to study the basic characteristics of skeletal muscle contractions.

## III. Research Background

Our research starts from Indonesian's medical faculty for physiology experiment by observing live frog's muscle contraction using conventional kymograph as a device. However, we notice some weaknesses from conventional kymograph's usage in order to observe and record muscle contraction.

Before discussing how our proposed device can overcome conventional kymograph's weaknesses, we will discuss the principal works of conventional kymograph.

## A. How Conventional Kymograph Works

To observe the phenomenon of skeletal muscle contractions, muscles can be removed from anesthetized frogs. These muscles can be attached to kymograph's recording systems and stimulated by electrical shocks of varying strength, duration, and frequency. Recordings obtained from such procedures can be used to study the basic characteristics of skeletal muscle contractions. Figure 1 and 2 show the kymograph from upper and front view. The procedures of using conventional kymograph are described as follows:

- The Kymograph, which is used to study muscle physiology, consists of a drum, which rotates at a pre-set speed, and traces are produced on the paper by means of an ink writing pointer. Thus, the drum can be rotated rapidly if rapid physiological events are being recorded or rotated slowly for events that occur more slowly.

- A *stylus* that can mark on the paper is attached to a *movable lever,* and the lever, in turn, is connected to an isolated muscle. The origin of the muscle is fixed in position by a *clamp,* and its insertion is hooked to the muscle lever.

- To study muscle physiology, a simple experiment can be carried out by using the frog gastrocnemius muscle/sciatic nerve. The principles of muscle excitation, contraction and work performed in the frog are similar to all vertebrates including man. Using frog

tissue has the practical advantage that it will function at room temperature without a blood supply; its oxygen requirements are met by diffusion from the air into the solution bathing the preparation.

- The muscle also is connected by wires to an *electronic stimulator.* The stimulator within the kymograph delivers small electric pulses to the muscle or nerve tissue via a pair of platinum electrodes.

- The stimulus can deliver single pulses or repeated stimulation up to a frequency of 100 per sec. An audible click is produced when a pulse is delivered. The stimulator may also be triggered using the trigger switch attached to the kymograph spindle.

## B. The Need for Computerized Kymograph

Based on our experiment and observation, we notice that conventional kymograph has some weaknesses, such as: damping risk on muscle lever caused by friction between the lever and the kymograph's drum; therefore the result recorded is not accurate as it is supposed to be.

Hence, our proposed computerized-based kymograph could overcome the problem because the result recorded immediately afterward the sensor get the distance as muscle contraction result. We also argue that our device produces more visible graphical result and also easier to use than the conventional kymograph.



Fig. 1.   Conventional Kymograph (front view)



Fig. 2.   Conventional Kymograph (upper view)

## IV. Constructing Computerized Kymograph

This section discusses how computerized-based kymograph is developed. The device consists of hardware and software systems, as discussed below.

### A. Hardware System

The hardware system consists of two main parts, which are: (1) ultrasonic distance sensor and (2) microcontroller as an interface between the sensor and the computer. The sensor works in the following ways: A sonic pulse is emitted from the sensor; then when the pulse bounces off of an object, an echo is returned. The sensor is able to emit the pulse because of a transducer that converts between sonic, electrical and mechanical energies.

We use ultrasonic sensor because it provides an easy method of distance measurement and has easy interfacing. This sensor is perfect for any number of applications that require to perform measurements between moving or stationary objects [6]. Some application ideas that can use this sensor are: security systems, interactive animated exhibits, parking assistant systems, and robotic navigation. The key features in this ultrasonic sensor are [6]: (1) provides precise, non-contact distance measurements within a 2 cm to 3 m range; (2) ultrasonic measurements work in any lighting condition, making this a good choice to supplement infrared object detectors; (3) simple pulse in/pulse out communication requires just one I/O pin; (4) burst indicator LED shows measurement in progress; (5) 3-pin header makes it easy to connect to a development board, directly or with an extension cable, hence no soldering required.

We use the sensor to measure the distance of muscle contraction. The sensor measures the time required for the echo return, and returns this value to the microcontroller as a variable-width pulse via the same I/O pin. The sensor works as an input to the microcontroller by sending the data continuously and display the graphical image on the computer display (by setting timer initialization). Figure 3 shows the ultrasonic distance sensor, and figure 4 shows how the sensor works during the time.



Fig. 3. Ultrasonic distance sensor [6]



Fig. 4. Flowchart of Ultrasonic distance sensor

We also use microcontroller based on the AVR enhanced RISC architecture. By executing instructions in a single clock cycle, the ATmega8535 achieves throughputs approaching 1 MIPS per MHz allowing the system's designer to optimize power consumption versus processing speed [5]. Figure 5 shows the components of microcontroller.

Figure 6 demonstrates all stages required for the device as a complete system. The system starts whenever the frog's muscle is pinched (figure 10), the device will recognized a distance disparity using ultrasonic distance sensor. Having the distance data, the device will read and process it (figure 8). Those data processing and calculating are performed using Microkontroller ATMEGA 8535. Then, using serial port the data will be delivered and displayed to computer, as a data receiver and recorder (figure 8).

Fig. 5. Flowchart of Ultrasonic distance sensor [5]

*B. Software System*



Fig. 6. Software system in Computerized kymograph algorithm



Fig. 7. Computerized kymograph system

Figure 7 shows how software system in the device is connected to hardware system. Ultrasonic sensor is employed to measuring the distance caused by muscle contraction over a period of time.

The result is then being processed and calculated by means of microcontroller. The microcontroller also used as an interface device to display the graphical image on the computer's display. Figure 9 demonstrates how the sensor and microcontroller connected tocomputerized-kymograph.

*C. Device Testing*

In this section, we will discuss the device's result compare to the conventional kymograph. Figure 10 demonstrate a live frog's muscle to be tested by our computerized kymograph device.

These frog's muscles can be attached to kymograph's recording systems and stimulated by electrical shocks of varying strength, duration, and frequency. Recordings obtained from the sensor and processed by microcontroller. The result is then displayed on the computer.

Figure 11 displays the result on conventional kymograph; while figure 12 displays the result on computerized kymograph. We can see from both results that computerized kymograph (figure 11) shows more visible result than conventional kymograph, by providing muscle contraction value (figure 12).

Fig. 8.   Hardware system in Computerized kymograph



Fig. 9.   Sensor and microcontroller used in computerized-kymograph



Fig. 10.   Frog's muscle used in experiment

Fig. 11. Conventional kymograph result



Fig. 12. . Computerized kymograph result

## V. RESULTS AND DISCUSSION

Figure 12 shows the result of Computerized Kymograph, as our proposed device. From figure 11 and figure 12 above, we can compare the results of conventional kymograph (figure 11) and computerized kymograph (figure 12). It can be seen that our proposed device displays time value (x axis) and distance value (y axis); while the conventional one do not provide those value but graphics.

Our device also provides more accurate result because it reduces friction force resulted from the conventional kymograph's drum. Force of friction can be calculated by [7]:

$$F_f = \mu F_N \qquad \dots\dots\dots\dots\dots\dots\dots(1)$$

Where:

$F_f$ is the force of friction in N,
$\mu$ is the coefficient of friction, and
$F_N$ is the normal force in N.
The value of $\mu$ depends on surface that dealing with.

By recording and displaying the result directly on the computer, we argue that the muscle contraction result on our device displays more accurate result than the conventional one. This also will lead better result's interpretation. Moreover, our device provides values on result, for example, the result (figure 12) shows that the maximum contraction

achieved around 170 mm - 180 mm at time (seconds) 20, 90, 120, 150, 162, and 278.

## VI. CONCLUSION

Our proposed device, computerized kymograph demonstrates better result by providing muscle contraction value. It is also easier to use than the conventional kymograph. Recently, the device has been used in our Medical Faculty, Yarsi University, Jakarta, Indonesia as a device to learn about physiology. The experiment on our laboratory shows that most of the students can interpret better using our device than the conventional one.

As an improvement in the future, we need to assess the device from user experience perspectives and device's functionality

## VII. ACKNOWLEDGMENT

REFERENCES

[1] A. Mukherjee, B. Jenkins, C. Fang, R.J. Radke, G. Banker, B. Roysam, Automated Kymograph Analysis for Profiling Axonal Transport of Secretory Granules,

[2] I. Smal, I. Grigoriev, A. Akhmanova, W. Niessen, E. Meijering, Microtubule dynamics analysis using kymographs and variable-rate particle filters, Image Processing, IEEE Transactions on 19, no. 7, pp. 1861–1876, 2010.

[3] O. Welzel, D. Boening, A. Stroebel, U. Reulbach, J. Klingauf, J. Kornhuber, T. Groemer, Determination of axonal transport velocities via image cross-and autocorrelation, European Biophysics Journal no. 38, pp. 883–889, 2009.

[4] W. B. Ludington, W. F. Marshall, Automated analysis of intracellular motion using kymographs in 1, 2, and 3 dimensions, Vol. 7184, SPIE, 2009, URL http://link.aip.org/link/?PSI/7184/71840Y/1.

[5] http://www.atmel.com/images/doc2502.pdf accessed on 1 November 2013

[6] www.parallax.com/product/28015 accessed on 1 November 2013

[7] http://library.thinkquest.org/10796/ch4/ch4.htm accessed on 27 September 2013

# Comparative Study in Performance for Subcarrier Mapping in Uplink 4G-LTE under Different Channel Cases

Raad Farhood Chisab[1,2]

[1]Foundation of Technical Education, IRAQ
[2]Dept. of ECE, SHIATS (Deemed to be University)
Allahabad-211 007, UP, INDIA

Prof. (Dr.) C. K. Shukla

Prof. at Dept. of ECE, SHIATS,
(Deemed to be University) Allahabad-211 007,
UP, INDIA

*Abstract*—in recent years, wireless communication has experienced a rapid growth and it promises to become a globally important infrastructure. One common design approach in fourth generation 4G systems is Single Carrier Frequency Division Multiple Access (SC-FDMA). It is a single carrier communication technique on the air interface. It has become broadly accepted mainly because of its high resistance to frequency selective fading channels. The third Generation Partnership Project-Long Term Evolution (3GPP-LTE) uses this technique in uplink direction because of its lower peak to average power ratio PAPR as compared to Orthogonal Frequency Division Multiple Access (OFDMA) that is used for downlink direction. In this paper the LTE in general and SCFDMA will be discuss in details and its performance will be study under two types of subcarrier mapping which are localized and distributed mode also within different channel cases. The results show that the localized subcarrier mapping give lower bit error rate BER than the distributed mode and give different activity under miscellaneous channel cases.

*Keywords—LTE; SCFDMA; 4G; PAPR; BER; channel model*

## I. INTRODUCTION

Wireless communications is an emerging field which has seen enormous growth in the last several years. The unprecedented and ubiquitous use of mobile phone technology, rapid expansion in wireless local area networks (WLAN) and the exponential growth of the Internet have resulted in an increased demand for new methods of establishing high capacity wireless networks. As the wireless standards evolved, the access techniques used also exhibited increase in efficiency, capacity and scalability. The first generation wireless standards used Frequency Division Multiple Access (FDMA) or Time Division Multiple Access (TDMA) [1]. In wireless channels, FDMA consumed more bandwidth for guard to avoid inter-carrier interference (ICI) and TDMA proved to be less efficient in handling high data rate channels as it requires large guard periods to alleviate the multipath impact. 4G (4th Generation) mobile networks are evolving to provide a comprehensive IP-based integrated solution at an affordable price where voice, data and streamed multimedia can be given to users on an anytime, anywhere basis, and at higher data rates than previous generations. This will be achieved after the convergence of all types of wired and wireless technologies and will be capable of providing data rates between 100 Mbps and 1Gbps (both indoors and outdoors), with premium quality and high security.

High data rate calls upon an improved spectral efficiency. The Third Generation Partnership Project Long Term Evolution (3GPP-LTE) has been standardized for the emerging 4th generation (4G) wireless communications [2]. The OFDMA and SCFDMA are technique the most prominent candidates that are used in 4G mobile systems. The LTE decided to use the OFDMA for downlink and using the SCFDMA for uplink [3]. The choice of SCFDMA in uplink direction comes as a result of its ability to reduce the PAPR as compare with the OFDMA and also give lower BER in case of localized mode than the distributed mode.

## II. LONG TERM EVOLUTION (LTE)

LTE system is expected to be competitive for many years to come, therefore, the requirements and targets set forth for this system are quite stringent. The main objectives of the evolution are to further improve service provisioning and reduce user/operator costs. The parameter of LTE can be summarized in table 1 [4]. A key requirement for LTE is to make possible a seamless transition from current telecommunication systems. This can be made possible by reuse of the current spectrums, interoperability between current and upcoming system, reuse of existing sites and production competitively priced equipment. It gives the operators the ability to migrate to new systems with ease [5]. But this requires adoption of simplified system architecture, stringent limits on spectrum and usage of a new radio-access technology with better characteristics. Transmission parameters in LTE consist of frequency, space, and time to create transmission resources for carrying data [6].

All the LTE signals derive their timing from a clock operating at 30.72 MHz = 15 kHz × 2048. This is the timing required for the 2048 point discrete Fourier transform (DFT) specified for 20 MHz channels. Therefore, the basic time interval in an LTE physical channel is one clock period of duration $T_s = 1/(15k \times 2048) = 32.255\ ns$ per clock period. The LTE radio frame for downlink and uplink transmission is $307200 \times T_S = 10ms$ long. LTE supports two radio frame structures which are frequency division duplex FDD which uses type 1 frame structure and time division duplex TDD which is applicable to type 2 frame structure [7].

A radio frame consists of 10 sub frames $30720 \times T_S = 1ms$ in FDD and two half- frames $153600 \times T_S = 5ms$ in

TDD. A half-frame is divided into four subframes and a special subframe, or five subframes, based on downlink to uplink switch point periodicity. The TDD frame structure can be configured in seven different sub frame formats. The sub frames 0 and 5 and DwPTS (downlink pilot timing slot) are reserved for downlink transmission. The sub frame that appears after special sub frame as well as UpPTS (uplink pilot timing slot), is always assigned to uplink transmission. Each sub frame in both FDD and TDD has two slots of $15360 \times T_S = 0.5ms$.

TABLE I. THE IMPORTANT PARAMETERS FOR LTE-SCFDMA

| Parameters | value | | | | | |
|---|---|---|---|---|---|---|
| BW (MHz) | 1.25 | 2.5 | 5 | 10 | 15 | 20 |
| Resource Block | 6 | 12 | 25 | 50 | 75 | 100 |
| FFT Size | 128 | 256 | 512 | 1024 | 1536 | 2048 |
| $f_s$ (MHz) | 1.92 | 3.84 | 7.68 | 15.36 | 23.04 | 30.72 |
| Sample per slot | 960 | 1920 | 3840 | 7680 | 11520 | 15360 |
| No. of sub carrier | 76 | 151 | 301 | 601 | 901 | 1201 |
| Carrier spacing | 15 KHz | | | | | |
| (PRB) BW | 180 KHz | | | | | |
| No. of OFDM symbol/slot | 7 for normal CP and 6 for extended CP | | | | | |
| Full mobility | Up to 500 Km/h | | | | | |
| Capacity | > 200 User per cell | | | | | |
| Cell size | 5-100 Km | | | | | |

A resource element, consisting of one subcarrier during one OFDM symbol, is the smallest physical resource in LTE. Furthermore, as illustrated in Fig. 1, resource elements are grouped into physical resource block (PRB), where each physical resource block consists of a bandwidth equal to 180 kHz (12 consecutive subcarriers) in the frequency domain and one 0.5 ms (one slot) in the time domain [8].



Fig. 1. The resource block and subcarrier in LTE-SCFDMA

Depending on the cyclic prefix (CP) type, which is a copy of the last portion of the data symbol which is inserted in front of the same data symbol during the guard interval, LTE employed two types of cyclic prefix, namely normal CP and extended CP [9]. The duration of an extended cyclic prefix is 512 clock periods, $512 \times TS = 16.67$ $\mu$sec. In slots with seven symbols, the duration of a normal cyclic prefix is 160 clock periods, $160 \times TS = 5.21$ $\mu$sec, for the first symbol and 144 clock periods, $144 \times TS = 4.69$ $\mu$sec, for the other six symbols [7].

## III. SC-FDMA

There is considerable interest in the use of Single Carrier Frequency Division Multiple Access (SC-FDMA) as the uplink transmission scheme in the 3GPP-LTE standard. This interest is justified by the inherent single carrier structure of SC-FDMA, which results in reduced sensitivity to phase noise and a lower Peak-to-Average Power Ratio compared to Orthogonal Frequency Division Multiple Access OFDMA [9].

SC-FDMA, which utilizes single carrier modulation and frequency domain equalization, is a technique that has similar throughput and essentially the same overall structure as OFDMA [10]. One advantage over OFDMA is that the SC-FDMA signal has lower peak-to-average power ratio (PAPR) because of its inherent single carrier structure. SC-FDMA has attracted attention as an alternative to OFDMA especially in uplink communications where lower PAPR benefits the mobile terminal in terms of transmit power efficiency. SC-FDMA has been adopted as the uplink multiple access scheme for the 3rd Generation Partnership Project Long Term Evolution [11].

As shown in Fig. 2, the transmitter of an SC-FDMA system converts a binary input signal to a sequence of modulated subcarriers. At the input to the transmitter, a baseband modulator transforms the binary input to a multilevel sequence of complex numbers $x_n$ in one of several possible modulation formats. The transmitter next groups the modulation symbols $\{x_n\}$ into blocks each containing $N$ symbols. The first step in modulating the SC-FDMA subcarriers is to perform an $N$-point DFT to produce a frequency domain representation $X_k$ of the input symbols. The DFT equation is represented as [12]:

$$X_k = \sum_{k=0}^{N-1} x_n e^{\frac{-2\pi jkn}{N}} \quad k = 0, 1 \dots N - 1 \qquad (1)$$

It then maps each of the $N$ DFT outputs to one of the $M$ ($> N$) orthogonal subcarriers that can be transmitted. If $N = M/Q$ and all terminals transmit $N$ symbols per block, the system can handle $Q$ simultaneous transmissions without co-channel interference. $Q$ is the bandwidth expansion factor of the symbol sequence. The result of the subcarrier mapping is the set $\tilde{X}_l$ ($l = 0, 1, 2\dots, M$-1) of complex subcarrier amplitudes, where $N$ of the amplitudes are non-zero. As in OFDMA, an $M$-point IDFT transforms the subcarrier amplitudes to a complex time domain signal $\tilde{x}_m$. The Inverse discrete Fourier transform IDFT equation is represented as [12]:

$$\tilde{x}_m = \frac{1}{M} \sum_{l=0}^{M-1} X_l e^{\frac{2\pi jkm}{M}} \quad m = 0, 1 \dots M - 1 \qquad (2)$$

There are $M$ subcarriers, among which $N$ ($< M$) subcarriers are occupied by the input data. In the time domain, the input data symbol has symbol duration of $T$ seconds and the symbol duration is compressed to $\tilde{T} = T \frac{N}{M}$ seconds after going through SC-FDMA modulation.

There are two types of sub-carrier mapping which are localized and distributed mapping as shown in Fig. 3. In localized mapping the output from the DFT is mapped to a subset of consecutive subcarrier, confining only to a fraction of system bandwidth and the zero padding process is done either at the first or last, but the outputs of the DFT will be placed in the sequence order without any interchanging [13]. In

distributed mapping the output of the DFT is assigned, non-continuously to the sub-carrier, over the entire bandwidth and the zero padding is done equally over the entire bandwidth [14]. The data block consists of N complex modulation symbols generated at a rate $R_{source}$ (symbols/sec). The N-point FFT produces N frequency-domain symbols that modulate N out of M orthogonal sub-carriers spread over a bandwidth $W$. The sub-carriers mapping process can be shown in Fig. 4. Where $W$ can be defined as [15]:

$$W = M.F_0 \quad Hz \tag{3}$$

Where $F_0$ (Hz) is the sub-carriers frequency spacing. The channel transmission rate is:

$$R_{channel} = [M/N].R_{source} \quad (Symbol/sec) \tag{4}$$

The bandwidth spreading factor Q is given by:

$$Q = R_{channel}/R_{source} = M/N \tag{5}$$



Fig. 2.    The block diagram of the SCFDMA system

For LFDMA, the frequency samples after subcarrier mapping $\{\tilde{X}_l\}$ can be described as follows [13]:

$$\tilde{X}_l = \begin{cases} X_l, & 0 \le l \le N-1 \\ 0, & 0 \le l \le M-1 \end{cases} \tag{6}$$

Let $m = Q.n + q,$ where $0 \le n \le N-1$
and $0 \le q \le Q-1$ Then

$$\tilde{x}_m = \tilde{x}_{Qn+q} = \frac{1}{M}\sum_{l=0}^{M-1} x_l\, e^{j2\pi i \frac{m}{M}} \tag{7}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{l=0}^{n} x_l e^{j2\pi i \frac{Qn+q}{QN}} x_l \tag{8}$$

If q=0 then

$$\tilde{x}_m = \tilde{x}_{Qn} = \frac{1}{Q}\frac{1}{N}\sum_{l=0}^{N-1} X_l e^{j2\pi l \frac{Qn}{QN}} \tag{9}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{l=0}^{N-1} X_l e^{j2\pi l \frac{n}{N}} \tag{10}$$

$$= \frac{1}{Q} x_n = \frac{1}{Q} x_{(m)\,mod\,N} \tag{11}$$

If $q \neq 0, since\ X_l = \sum_{p=0}^{N-1} x_p\, e^{-j2\pi l\frac{p}{N}}$ then eqn. 6 can be expressed as follows:

$$\tilde{x}_m = x_{Qn+q} = \frac{1}{Q}\frac{1}{N}\sum_{l=0}^{N-1} X_l\, e^{j2\pi l\frac{Qn+q}{QN}} \tag{12}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{l=0}^{N-1}\left(\sum_{p=0}^{N-1} x_p\, e^{-j2\pi l\frac{p}{N}}\right) e^{j2\pi l\frac{Qn+q}{QN}} \tag{13}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{l=0}^{N-1}\sum_{p=0}^{N-1} x_p e^{j2\pi l\left\{\frac{(n-p)}{N}+\frac{q}{QN}\right\}} \tag{14}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{p=0}^{N-1} x_p \left(\sum_{i=0}^{n} e^{j2\pi l\left\{\frac{(n-p)}{N}+\frac{q}{QN}\right\}}\right) \tag{15}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{p=0}^{N-1} x_p\, \frac{1-e^{j2\pi(n-p)}e^{j2\pi\frac{q}{Q}}}{1-e^{j2\pi\left\{\frac{(N-P)}{N}+\frac{q}{QN}\right\}}} \tag{16}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{p=0}^{N-1} x_p\, \frac{1-e^{j2\pi\frac{q}{Q}}}{1-e^{j2\pi\left\{\frac{(n-p}{N}+\frac{q}{QN}\right\}}} \tag{17}$$

$$= \frac{1}{Q}\left(1-e^{j2\pi\frac{q}{Q}}\right)\frac{1}{N}\sum_{p=0}^{N-1}\frac{x_p}{1-e^{j2\pi\left\{\frac{(n-p)}{N}+\frac{q}{QN}\right\}}} \tag{18}$$

As can be seen from eqn. 9 and 16, LFDMA signal in the time domain has exact copies of input time symbols with a scaling factor of 1/Q in the N-multiple sample positions and in between values are sum of all the time input symbols in the input block with different complex-weighting.

Now, For DFDMA, the frequency samples after subcarrier mapping $\tilde{X}_l$ can be described as follows.

$$\tilde{X}_l = \begin{cases} X_{l/\tilde{Q}}, & l = \tilde{Q}.k\ (0 \le k \le N-1) \\ 0 & , otherwise \end{cases} \tag{19}$$

Where $0 \le l \le M-1$, $M = Q.N$, and $1 \le \tilde{Q} \le Q$

Let $m = Q.n + q\ (0 \le n \le N-1,\ 0 \le q \le Q-1)$

Then

$$\tilde{x}_m\left(= \tilde{x}_{Q.n+q}\right) = \frac{1}{M}\sum_{l=0}^{N-1} \tilde{X}_l\, e^{j2\pi l\frac{m}{M}} \tag{20}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{k=0}^{N-1} X_k e^{j2\pi\tilde{Q}k\frac{Qn+q}{QN}} \tag{21}$$

If q=0 then

$$\tilde{x}_m = \tilde{x}_{Q.n} = \frac{1}{Q}\frac{1}{N}\sum_{k=0}^{N-1} X_k e^{j2\pi\tilde{Q}k\frac{Qn}{QN}} \tag{22}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{i=0}^{n} X_k e^{j2\pi\tilde{Q}k\frac{n}{N}} \tag{23}$$

$$= \frac{1}{Q}\frac{1}{N}\sum_{i=0}^{n} X_k\, e^{j2\pi k\frac{\tilde{Q}n}{N}} \tag{24}$$

$$= \frac{1}{Q}\left(\frac{1}{N}\sum_{i=0}^{n} X_k\, e^{j2\pi k\frac{(\tilde{Q}.n)\,mod\,N}{N}}\right) \tag{25}$$

$$\frac{1}{Q} x_{(\tilde{Q}.n)\,mod\,N} = \frac{1}{Q} x_{(Q(m)\,mod\,N)\,mod\,n} \tag{26}$$

If $q \neq 0$, since $X_k = \sum_{p=0}^{N-1} x_p e^{-j2\pi k \frac{p}{N}}$ Eqn. 21 can be expressed as follows after derivation

$$\tilde{x}_m = \tilde{x}_{Q.n+q} = \frac{1}{Q}\left(1 - e^{J2\pi q \frac{\bar{Q}}{Q}}\right)\frac{1}{N}\sum_{p=0}^{N-1}\frac{x_p}{1 - e^{j2\pi\left\{\frac{(\bar{Q}n-p)}{N} + \frac{\bar{Q}q}{QN}\right\}}} \quad (27)$$



Fig. 3.   The two types of sub-carrier mapping

From a resource allocation point of view, subcarrier mapping methods are further divided into static and channel-dependent scheduling (CDS) methods. CDS assigns subcarriers to users according to the channel frequency response of each user [16]. CDS is of great benefit with localized subcarrier mapping because it provides significant multi-user diversity which leads to improved system capacity and performance [17]. For these reasons only LFDMA concept is proposed to use in the 3GPP-LTE specifications.



Fig. 4.   The process of sub-carriers mapping

The transmitter performs two other signal processing operations prior to transmission. It inserts a set of symbols referred to as a cyclic prefix (CP) in order to provide a guard time to prevent inter-block interference due to multi-path propagation. The transmitter also performs a linear filtering operation referred to as pulse shaping in order to reduce out-of-band signal energy. The receiver transforms the received signal into the frequency domain via DFT, de-maps the subcarriers, and then performs frequency domain equalization (that will be discussed later). The equalized symbols are transformed back to the time domain by means of an IDFT, and detection and decoding take place in the time domain [18]. After the Subcarrier de-mapping is done. The de-mapped signal is given to the IDFT to get the time domain signal back. The IDFT

output is given for QPSK or QAM demodulation. After the demodulation the receiver generate the final bit stream [19].

IV.   CHANNEL EQUALIZATION

SCFMA suffers from Inter-symbol Interference (ISI) if a transmission over a frequency selective channel is considered. Therefore for mobile radio applications, SCFDMA requires Equalization at the receiver. The only use of Guard Interval (GI) for equalization does not meet the challenges of the future mobile radio system because of the reduction of spectral efficiency [20].

Channel equalization is one of the key blocks in LTE receiver. It is one of the most important elements of wireless receivers that employ coherent demodulation. For practical LTE systems, it is important to have an equalization technique that is specifically designed for LTE pilots, and has low computational and hardware complexities. An equalizer within a receiver compensates for the average range of the expected channel amplitude and delay characteristic. Equalizer must be adaptive since the channel is generally unknown and time varying [21].

Equalizer is always used in both time and frequency domains in traditional communication system. In the time domain, for traditional FDM system, equalization is indispensable. Because equalizer is used to balance the channel characteristics in the receiver, equalizer produces the opposite characteristics of channel to offset ISI by time varying multi-path channel. But equalization is not a satisfactory method for OFDM system [22]. A possible way to reduce the complexity of linear equalization is to carry out the equalization in the frequency domain. The equalization is carried out block-wise with block size *N*. The sampled received signal is first transformed into the frequency domain by means of a size-*N* DFT. The equalization is then carried out as frequency-domain filtering [23].

The received signal is equalized in the frequency domain. After the equalization block the equalized signal is then transformed back to the time domain using the IFFT. The method of equalization, which is shown in Fig. 5, is done by the following steps:

Let E(m) where (m=0, 1, 2…$N_{FFT}$ -1) denote the equalizer coefficient for the m[th] sub carrier, the time domain equalized signal K(n) can be expressed as:

$$k(n) = \frac{1}{N_{FFT}}\sum_{m=0}^{N_{FFT}-1} E(m)G(m) e^{\frac{i2\pi mn}{N_{FFT}}} \quad (28)$$

Where $n = 0,1,2,\dots,N_{FFT} - 1$

The equalizer coefficients E(m) are determined to minimize the mean square error between the equalized signal and the original signal. The equalizer coefficients are computed according to the types of the frequency domain equalization (FDE) in two methods as follow [24]:

A.   The zero forcing (ZF) Equalizer is

$$E(m) = 1/H(m) \qquad m = 0,1,2,\dots,N_{FFT} - 1 \quad (29)$$

B.   The Minimum Mean Square Error (MMSE) Equalizer is

$$E(m) = H^*(m)/[|H(m)|^2 + (E_b/N_0)^{-1}] \qquad (30)$$

Where * denotes the complex conjugate, *H(m)* is the transfer function of the channel and $E_b/N_0$ is average energy-per-bit to noise power spectral density. Equalization will be used to eliminate the effect of ISI.



Fig. 5. The process of channel equalization

From Fig. 6 it can be noticed that the MMSE method is better than the ZF method and give lower BER compared with other method. Therefore, in all tests and simulations for channel models, the MMSE method will be use.



Fig. 6. the performance under two types of channel equalization

In the receiver side, OFDMA utilizes a simple equalizer per subcarrier after FFT. But, SC-FDMA utilizes a complex equalizer before sending the resultant to IFFT. IFFT removes the effect of the FFT in the transmitter. Notice that result of the IFFT is again a time domain signal; the time domain signal is sent to a single detector to create the bits. These differences in receiver side are illustrated in Fig. 7 in which we can see the equalizer simplicity of OFDMA against SC-FDMA. As you can see, SC-FDMA receiver is more complex than OFDMA, but in the transmitter simpler power amplifiers can be utilized to reduce the power consumption. These fortify the SC-FDMA as an uplink transmission scheme, since power efficiency and complexity is important for mobile stations but not in the base station [25].



Fig. 7. The equalization in OFDMA and SCFDMA

## V. THE WIRELESS CHANNELS SPECIFICATIONS

In wire-line communication, the data transmission is primarily corrupted by statistically independent Gaussian noise, as known as the classical additive white Gaussian noise (AWGN). In absence of interference, the primary source of performance degradation in such wire-line channels is thermal noise generated at the receiver. Reliable communication in wireless or radio channels, however, becomes a difficult task as the transmitted data is not only corrupted by AWGN, but also suffers from inter-symbol interference (ISI), in addition to (large-scale and small-scale) fading as well as interference from other users. To master the art of wireless communications, one must understand the propagation characteristics of a radio channel [26]. The fading in radio propagation can be classified into two groups; large-scale fading and small-scale fading as illustrated in Fig. 8. Large-scale fading manifests itself as the average signal power attenuation or path loss due to motion over large areas as shown in blocks 1, 2 and 3. Small-scale fading refers to the dramatic changes in the signal amplitude and phase that occur due to small changes in the spatial separation between the transmitter and the receiver. As indicated by blocks 4, 5 and 6 in Fig. 8, small-scale fading manifests itself in two mechanisms namely, time-spreading of the signal (or channel dispersion) and time-variant nature of the channel [27]. The signal time-spreading (signal dispersion) and time-variant nature of the channel may be examined in two domains, time and frequency, as indicated in block 7, 10, 13 and 16. For signal dispersion nature, we categorize the fading degradation types as frequency selective and frequency non-selective (flat) as illustrated in blocks 8, 9, 11, and 12. For time-variant nature, we categorize the fading degradation types as fast fading and slow fading, as shown in blocks 14, 15, 17, and 18.

Large-scale fading is responsible for path-loss in wireless communications and large-scale fading models typically find applications in mobile network planning and understanding free space wireless communication over large areas . In most practical wireless communication systems, the radio communications is far more complex than free-space situation, and is best explained by small-scale fading models [28].



Fig. 8.   the fading channel manifestations

Due to the reflection, diffraction and scattering by objects in the environment, transmitted signal propagates through different paths. Thus replicas of the transmitted signal arrive at the receiver with different time delays. This time delay variation is often quantified in terms of delay spread. Larger delay spread means there is a large variation in time delays of different multipath components [29].

The relation between delay spread and OFDM performance can be explained by the principle of frequency diversity. Delay spread is inversely proportional to coherence bandwidth, i.e. larger delay spread results in smaller coherence bandwidth. Coherence bandwidth is the bandwidth over which the channel is considered to be "flat". Within the coherence bandwidth, different signals experience the same channel frequency response. The fading characteristics within coherence bandwidth is flat, thus it is called flat fading. If the bandwidth of a signal is larger than the channel coherence bandwidth, the channel is considered as frequency selective channel. Signals with frequency difference more than coherence bandwidth experience different fading. Therefore, it is called frequency selective fading [30]. The effect of flat and selective fading can be shown in Fig. 9.

An important requirement for assessing technology for Broadband Fixed Wireless Applications is to have an accurate description of the wireless channel model. Channel models are heavily dependent upon the radio architecture. The profile of received signal can be obtained from that of the transmitted signal if we have a model of the medium between the two. This model of the medium is called channel model.



Fig. 9.   The effect of flat and selective fading channel on the signal

Channel models are essential tools for simulation and testing of wireless transmission systems. The literature is extensive on this topic, and many standards have recommended channel models for specific propagation environments. These models may characterize path-loss attenuation, shadowing and multipath effects [31]. In this paper some of these channel models will be studies in order to investigate the performance of the system under these channel models. more than one channel model will be discuss and investigate how the system will work under these channel models. These channel models are:

### A.  COST 207 channel models

The COST 207 model gives normalized scattering functions, as well as amplitude statistics for four typical environments which are rural area (RA), typical urban area (TU), bad urban area (BU), and hilly terrain (HT). The COST 207 model was presented as an outdoor wireless channel model. This model specifies power gains and time delays for four typical environments [32]. These parameters were evaluated by numerous measurements performed in many countries, including the United Kingdom, France, and Sweden [33]. COST 207 standards provided both the continuous time formula and discrete taps model. The performance of system under these channel models can be shown in Fig. 10 and Fig. 11. Their power distributions are characterized as follows [34]:

For Rural Area (RA):

$$P(\tau) = \begin{cases} \exp(-9.2\tau) & 0 \le \tau \le 0.7\mu s \\ 0 & otherwise \end{cases} \qquad (31)$$

For Typical Urban (TU):

$$P(\tau) = \begin{cases} \exp(-\tau) & 0 \le \tau \le 7\mu s \\ 0 & otherwise \end{cases} \qquad (32)$$

For Bad Urban (BU):

$$P(\tau) = \begin{cases} exp(-9.2\tau) & 0 \le \tau \le 5\mu s \\ 0.5\exp(5-\tau) & 5 \le \tau \le 10\mu s \\ 0 & otherwise \end{cases} \qquad (33)$$

For Hilly Terrain (HT):

$$P(\tau) = \begin{cases} exp(-3.5\tau) & 0 \le \tau \le 2\mu s \\ 0.1\exp(15-\tau) & 15 \le \tau \le 20\mu s \\ 0 & otherwise \end{cases} \qquad (34)$$

### B. COST 259 channel models

The COST 259 directional channel model (DCM) was developed by the European COST259 project. The COST 259 DCM is wideband and capable of providing channel impulse responses in both spatial and temporal domains. It can also provide these in vertical and horizontal polarization components. It operates at the frequency range from 0.45 to 5 GHz and bandwidth of less than 10 MHz.

The model is very general, and describes the joint effects of small-scale as well as large-scale effects; it covers different cases of macro-cells, micro-cells and Pico-cells. The environments identified so far in COST 259 and typical speeds for each channel type are given in Table II. One of the work items identified in COST 259 is to propose a new set of channel models which overcome the limitations in the GSM channel models, while aiming at the same general acceptance.

The main difference between the COST 259 model and previous models is that it tries to describe the complex range of conditions found in the real world by distributions of channels rather than a few typical cases. The probability densities for the occurrence of different channels are functions of mainly two parameters which are Environment and Distance [35].

TABLE II. DEFAULT SPEED FOR THE CHANNEL MODELS

| Channel model | Model speed |
|---|---|
| TUX | 3 Km/h |
| | 50 Km/h |
| | 120 Km/h |
| RAX | 120 Km/h |
| | 250 Km/h |
| HTX | 120 Km/h |

In the COST 259 model a large number of paths ensure that the correlation properties in the frequency domain are realistic. Path powers follow the exponential channel shapes. There are three types of channel models which are the Rural Area channel model (RAx), The Hilly Terrain channel model (HTx) and the Typical Urban channel model (TUx) [36]. The performance of the system under these channel models can be shown in Fig. 12.

### C. LTE channel models

The LTE standard adopts models based on the ITU-R M.1225 recommendation and the 3GPP TS 05.05 specification for GSM, widely used in the context of third generation mobile systems [37]. The ITU and 3GPP models are defined by

tapped-delay line (TDL) models, where each tap corresponds to a multipath signal characterized by a fixed delay, relative average power and Doppler spectrum. This model use the Pedestrian A and Vehicular A channels from [33], and the Typical Urban (TU) channel from [38], in order to model three reference environments characterized by a low, medium and large delay spread, respectively. Nevertheless, they were designed for a 5 MHz operating bandwidth, and an apparent periodicity appears in their frequency correlation properties for higher bandwidths [39].

The LTE channel models developed by 3GPP are based on the existing 3GPP channel models and ITU channel models. The extended ITU models for LTE were given the name of Extended Pedestrian-A (EPA), Extended Vehicular-A (EVA) and Extended Typical Urban (ETU). These channel models are classified on the basis of low, medium and high delay spread where low delay spreads are used to model indoor environments with small cell sizes while medium and high delay spreads are used to model urban environments with large cells. The high delay spread models are according to Typical Urban GSM model [40]. The performance of the system under these channel models can be shown in Fig. 13.

## VI. RESULTS AND DISCUSSION

The system (3GPP-LTE-SC-FDMA) based on FFT was simulated and run using MATLAB package version 7.12 (R2011a). The behavior of the proposed system was monitored under the parameters that effect on the performance of the system. These parameters are listed in table III.

The system was tested under three channel cases which are COST207, COST259, and LTE channel models. It can be noticed that the system have different responses under these channel models. The dominant thing that can be seen is that the system under localized subcarrier mapping is better than under distributed mode as shown in Fig. 10 to Fig. 13.



Fig. 10. The performance under COST 207 RA and TU channel models

Fig. 11. The performance under COST 207 BU and HT channel models



Fig. 12. The performance under COST 259 channel models



Fig. 13. The performance under LTE channel models

TABLE III.        THE PARAMETERS FOR SIMULATION OF SC-FDMA

| Parameters | Value |
|---|---|
| System bandwidth | 5 MHz |
| Modulation types | QPSK |
| Carrier Frequency ( fc ) | 2025 MHz |
| Sub-carriers spacing | 15 KHz |
| Sub-carriers mapping | Localized, Distributed |
| No. of sub-carrier | 256 |
| Channel equalization | ZF and  MMSE |
| Target BER | $10^{-3}$ |
| Channel estimation | Perfect |
| Channel Types | COST207, COST259, LTE channels |

## VII.    CONCLUSION

In this paper the system of SCFDMA was examined under different parameters but the important things is to study the two parameters which are the types of channel equalization and the types of subcarrier mapping. First, when we notice the behavior of system under two type of equalization which are zero forcing (ZF) and minimum mean square error (MMSE), we can notice and monitor the behavior of the system in the Fig. 6 and conclude that the MMSE equalization method was better than the ZF method and give lower bit error rate (BER) as compare with other method. Second, when notice the behavior of system through the Fig. 10 to Fig. 13, which specify the activity of system when changing the type of subcarrier mapping under different channel cases, it can be noticed that the system is run with better performance under the localized subcarrier mapping method for all types of channel models which are COST207, COST259 and LTE channel models and give lower bit error rate as compare with the other method which depends on the distributed subcarrier mapping. Also we can notice that the system give different activity during the different channel cases because each channel model has its own properties that effects on the system performance.

Finally, from all the results we can conclude that "First: MMSE equalization method is better than the ZF method, second: the localized subcarrier mapping is better than the distributed mode, third: the activity of the system is changes with different channel models"

### REFERENCES

[1]  Erik Dahlman, Stefan Parkvall, Johan Sköld and Per Beming, 3G evolution HSPA and LTE for mobile broadband, 2nd ed., Academic Press Elsevier, 2008.

[2]  S. Sesia, I. Toufik, and M. Baker, LTE the UMTS Long Term Evolution from theory to practice, 1st ed.,John Wiley and Sons Ltd., 2009.

[3] A. Jamalipour, T. Wada, and T. Yamazato, "A Tutorial on multiple access technologies for beyond 3G mobile networks," IEEE Communications Magazine, vol. 43, Issue 2, pp. 110 - 117, February 2005.

[4] 3GPP TS 36.211, Evolved Universal Terrestrial Radio Access (E-UTRA); Physical channels and modulation, 2009.

[5] 3GPP TR 25.913, 3GPP:Technical Specification Group Radio Access Network; Requirements for Evolved UTRA and Evolved UTRAN, 2008.

[6] S. Rappaport, Wireless communications: principles and practice, 2nd ed., Prentice-Hall, 2002.

[7] Ramjee Prasad, OFDM for Wireless communication systems, 2nd ed., Artech house universal Personal communication series, 2008.

[8] D. Astély, E. Dahlman, A. Furuskär, Y. Jading, M. Lindström, and S. Parkvall, "LTE: the evolution of mobile broadband," IEEE Communications Magazine, vol. 47, Issue: 4, pp. 44 – 51, April 2009.

[9] E. Dahlman, H. Ekström, A. Furuskär, Y. Jading, J. Karlsson, M. Lundevall, and S. Parkvall, "The 3G long term evolution - radio interface concepts and performance evaluation," IEEE Vehicular Technology Conference (VTC2006), pp. 137 – 141, May 2006.

[10] Mohamed Noune and Andrew Nix, "A novel frequency-domain implementation of tomlinson-harashima precoding for SC-FDMA," IEEE 69th Vehicular Technology Conference VTC, pp. 1 – 5, 2009.

[11] D. Haccoun and G. Begin, "High-rate punctured convolutional codes for viterbi and sequential decoding," IEEE Transactions on Communications, vol. 37 , Issue 11, pp. 1113 – 1125, 1989.

[12] Dhirendra Kumar Tripathi, S.Arulmozhi Nangai, R. Muthaiah, "FPGA implementation of scalable bandwidth single carrier frequency domain multiple access transceiver for the fourth generation wireless communication," Journal of Theoretical and Applied Information Technology, vol. 28 No.2, June 2011.

[13] H. G. Myung, "Single Carrier Orthogonal Multiple Access Technique for Broadband Wireless Communications," Ph.D. Dissertation, Polytechnic University, January 2007.

[14] W. H. Tranter, K. S. Shanmugan, T. S. Rappaport, and K. L. Kosbar, Principles of Communication Systems Simulation with Wireless Applications, 1st ed., Prentice Hall Professional Technical Reference PTR , 2004..

[15] M. A. Abd El-Hamed, M. I. Dessouky, F. Shawki, Mohammad K. Ibrahim, S. El-Rabaie, and F. E. Abd El-Samie, "Wavelet-Based SC-FDMA System," 29th National Radio Science Conference (NRSC 2012), pp. 447 – 460, 2012.

[16] Peng LI, Yu ZHU, Zongxin WANG, and Naibo WANG, "Peak-to-average power ratio of SC-FDMA systems with localized subcarrier mapping," IEEE Global Mobile Congress GMC2010, pp. 1 – 6, 2010.

[17] Weidong Wang, Yan Zhou, Yuan Sang, Xue Shen, Fan Li, and Yinghai Zhang, "A Ue-interfering area based inter cell interference coordination scheme in SCFDMA uplinks," 2nd IEEE International Conference on Network Infrastructure and Digital Content (IC-NIDC 2010), pp. 681 – 686, 2010.

[18] F. Classen and H. Meyr, "Frequency synchronization algorithms for OFDM systems suitable for communications over frequency selective fading channels," IEEE 44th Vehicular Technology Conference (VTC), pp. 1655 – 1659, 1994.

[19] Harri Holma and Antti Toskala, LTE for UMTS: OFDMA and SC-FDMA based radio Access,1st ed., John Wiley & sons, 2009.

[20] Sosth`ene Yameogo , Jacques Palicot, Laurent Cariou, "Blind time domain equalization of scfdma signal," IEEE 70th Vehicular Technology Conference Fall (VTC2009), pp. 1 – 4, 2009.

[21] 3GPP TS 36.212: Evolved Universal Terrestrial Radio Access (E-UTRA); Multiplexing and channel coding, 2009.

[22] Mehmet Kemal Ozdemir, Huseyin Arslan, "Channel estimation for wireless ofdm systems," IEEE Communications Surveys, 2nd Quarter, vol. 9, No. 2, pp. 18-48, 2007.

[23] D. Falconer, S. L. Ariyavisitakul, A. Benyamin-Seeyar, and B. Eidson, "Frequency domain equalization for single-carrier broadband wireless systems," IEEE Communications Magazine, vol. 40 , Issue 4, pp. 58 – 66, 2002.

[24] Yao Xiao, "Orthogonal frequency division multiplexing modulation and inter-carrier interference cancellation," MSc. Thesis, Louisiana State University, May 2003.

[25] 3GP TS 25.211, 3rd generation partnership project; technical specification group radio access network; Physical Channels and Mapping of Transport Channels onto Physical Channels (FDD), 2009.

[26] John Doble, Introduction to Radio Propagation for Fixed and Mobile Communications, 1st ed.,Artech House Publishers, 1996.

[27] Bernard Sklar, "Rayleigh fading channels in mobile digital communication systems part i: characterization," IEEE Communications Magazine, vol. 35 , Issue 7, pp. 90 – 100, 1997.

[28] Gordan L. Stuber, Principles of Mobile Communications, 2nd ed., Kluwer Academic Publishers, 2001.

[29] G.J. Foschini and M.J. Gans, "On limits of wireless communications in a fading environment when using multiple antennas," Wireless Personal Communications, vol. 6, pp. 311-335, March 1998.

[30] M. F. Pop and N. C. Beaulieu, "Statistical investigation of sum-of-sinusoids fading channel simulators," Global Telecommunications Conference, (GLOBECOM '99), vol. 1a , pp.419-426, 1999.

[31] S. Rajkumar, "Modelling of multipath fading channels for network simulation," PhD. Dissertation, Texas A&M University,2007.

[32] Peral Rosado, Lopez Salcedo, Gonzalo Seco , Francesca Zanier, and Massimo Crisci, "Evaluation of the LTE Positioning Capabilities under Typical Multipath Channels," 6th IEEE Advanced Satellite Multimedia Systems Conference ASMS, pp. 139 – 146, 2012.

[33] Ming-Xian Chang; Su, Y.T., "Blind joint channel and data estimation for OFDM signals in Rayleigh fading," IEEE 53rd Vehicular Technology Conference (VTC 2001), vol.2, pp. 791 -795, 2001.

[34] Ye Li, Leonard J. Cimini and Nelson R. Sollenberger, "robust channel estimation for ofdm systems with rapid dispersive fading channels," IEEE International Conference on Communications(ICC 98), pp. 1320 – 1324, 1998.

[35] 3G TR 25.943, 3rd Generation Partnership Project; Technical Specification Group (TSG) RAN WG4; Deployment aspects, Stockholm, Sweden, June 2001.

[36] ETSI TR 125 943 V4.0.0, Universal Mobile Telecommunications System (UMTS); Deployment aspects, 2006.

[37] Noman Shabbir, Muhammad T. Sadiq, Hasnain Kashif, and Rizwan Ullah, "Comparison of radio propagation models for long term evolution (LTE) network," International Journal of Next-Generation Networks, Vol.3, No.3, 2011.

[38] J. P. Dobbelsteen, "Mapping an LTE Baseband Receiver on a Multi-Core Architecture," MSc. Thesis, Eindhoven University Of Technology, 2009.

[39] ITU-R M.2135-1 International Telecommunication Union, Guidelines for evaluation of radio interface technologies for IMT-Advanced, 2008.

[40] Asad Mehmood and Waqas Aslam Cheema, "Channel Estimation For Lte Downlink," MSc. Thesis, Blekinge Institute of Technology, September 2009.

# A Competency-Based Ontology for Learning Design Repositories

Gilbert Paquette

CICE Research Chair, LICEF Research Center, TELUQ

Montreal, Canada

*Abstract*—Learning designs are central resources for educational environments because they provide the organizational structure of learning activities; they are concrete instructional methods. We characterize each learning design by the competencies they target. We define competencies at the meta-knowledge level, as generic processes acting on domain-specific knowledge. We summarize a functional taxonomy of generic skills that draws upon three fields of knowledge: education, software engineering and artificial intelligence. This taxonomy provides the backbone of an ontology for learning designs, enabling the creation of a library of learning designs based on their cognitive and meta-cognitive properties.

*Keywords*—*Learning Designs; Learning Objects Repository; Competency Referencin; Generic skills; Learning Design Ontology; Metadata for Learning Designs*

## I. INTRODUCTION

A search on the Internet reveals the importance given to competency-based learning and training [12]. Ministries of education, school boards and teacher training institutes use competency profiles to define school programs or required qualities from the teachers, especially in the use of technologies in education. Consulting companies present their expertise by enumerating competencies, marketing their services in this way. Other companies offer services or computerized tools to help their prospective customers define or manage the competence of their staff, looked upon as the main asset of an organization in a knowledge management perspective. Governmental agencies or professional associations use competency-based approaches to define conditions to the exercise of a profession and to plan their vocational training programs.

In the IMS-RDCEO specification [9], competencies are expressed as simple natural language sentences that state informally that a group of person has the "capacity" or the "knowledge" to do certain things. Competency profiles are in general loosely structured collections of such texts that are not always easy to interpret, communicate or use.

In our previous work [20,22,23] on the MISA Instructional Design method for eLearning, we have defined a structural definition for competencies using knowledge representation techniques. This definition is based on the interrelation of a meta-knowledge domain, where generic skills are described, and knowledge in an application domain to which generic skills are applied. Here we use the word "knowledge", not as a synonym for "concept", but for any intellectual structure (concept, procedure, principle, taxonomy, decision tree, sets of facts, etc.) that can be processed by a cognitive system.

In section II, we summarize the main elements of this competency model and its relation to other generic skills' taxonomies proposed in various fields. In section III, we define the backbone of a learning design ontology based on competencies and the generic skill component of a competency. In section IV, we develop an RDFS vocabulary to reference learning designs with a set of metadata that can be used to search and retrieve learning design from a repository of learning design objects.

## II. COMPETENCIES AND LEARNING DESIGNS

This section summarizes the basis of a competency definition and its relation to learning designs.

### A. Generic skills: integrating many views in a taxonomy

In order to solve classification, diagnosis or construction problems for example, it is necessary to mobilize corresponding classification, diagnosis or construction generic skills. Competencies can be defined by associating generic skills to some knowledge it can be applied to, demonstrating that this actor is able to solve a corresponding class of problems.

This first view sees generic skills as generic problem solving processes. The area of generic problems or tasks draws on the software engineering work of authors like by Chandrasekaran [9], McDermott [16], Steels [29] and Scheiber et al, 1993 [28]. Our generic skill taxonomy expand these various taxonomies.

Another view defines generic skills as active procedural meta-knowledge that can be applied to various knowledge domains. Procedural meta-knowledge is implicit in the work of researcher in science epistemology such as Thayse [30] , Popper [9] and Pitrat [24,25]. We make it explicit in the proposed generic skill's taxonomy.

A third view on generic skills is to be found in taxonomies of educational objectives elaborated by researchers in educational technology such as Bloom [3], Krathwohl et al [3], Romisowski [27], Gagné [7] and Merrill [17]. Our taxonomy integrates and extends these taxonomies.

### B. Competency as generic skill applied to knowledge

Integrating these viewpoints leads us to a structural definition of the notion of competency, as procedural meta-knowledge applied to specific knowledge. Romisowki [27] has

expressed very well the simultaneous acquisition of knowledge and meta-knowledge in the learning process: « The learner follows two kinds of objectives at the same time - learning specific new knowledge and learning to better analyze what he already knows, to restructure knowledge, to validate new ideas and formulate new knowledge ». This idea has been expressed concisely by Pitrat [24]: « meta-knowledge is being created at the same time as knowledge ». This is the essence of the notion of competency we are defining here.

This notion of competency, as a *generic skill applied to knowledge in an application domain* fits well also within the framework of action theory Bélisle et Linard [2], based on the work of cognitive science authors such as Vygotsky, Leontiev, Piaget, Searle and Bruner. The association between generic skills, seen as generic cognitive processes, and specific knowledge avoids an artificial separation between knowledge and know-how, integrating cognitive and meta-cognitive aspects that must be present together for thoughtful human action and learning.

This discussion leads us to a representation of competencies, as an association between a generic skill, to be represented graphically as a process model in a meta-knowledge domain, in relation to domain specific knowledge, also represented by a knowledge model.



Fig. 1.   Example of a generic skill model: Simulate a process

In many applications we have used the MOT graphic language [23] to represent both models. Figure 1 shows one example of generic skill, a simulation meta-process, that can be applied (through instantiation) to many specific process such as "Search the Internet" or "Extract a square root". The instantiated simulation process, *Simulate a search process on the internet*, provides the link between the meta-model and the application model (Internet processes), thus defining a competency: *to be able to simulate a search process on the Internet*. Figure 1 presents the input and output of the meta-process "Simulate a process", together with its component

subprocesses. Numbers refer to the generic skills in the taxonomy shown later in table I.

*C. Competencies and learning designs*

Now suppose we wish to design a course module to help a learner learn a process such as searching the Internet, or in another domain, a process to extract square roots. A good idea would be to use the simulation generic skill process model as a template. To do this, the main subprocesses of the simulation meta-process are transformed into learning activities. Then a second step is to instantiate the template with terms in the application domain, with Internet search terms or arithmetics terms.

The LD template derived from the generic process on figure 1 would contain at first the following learning activities.

- Activity 1: Consult the description of the process to be simulated and produce inputs to the process;
- Activity 2:  Select an applicable procesdure (or task) in the process;
- Activity 3:  Execute the selected task and produce its outputs;
- Activity 4: Check if the process is completed; if not, select an applicable procedure, repeat 2, 3 and 4; if completed, report excution trace and final output.

Then instatiating this LD template to the Internet search domain would provide the following learning design assignments.

- Activity 1: Read the Internet search process provided by the instructor on the course Web site and select possible search requests you could make;
- Activity 2:  Select a request  to search the internet;
- Activity 3:   Build and execute a search request to produce a list of Web pages;
- Activity 4: Is the result satisfactory ? ; if not, refine the search request and repeat 2, 3 and 4; if completed, report list of steps and final list of Web pages.

The important thing here is that the generic process provides the backbone of the learner's assignments in a learning design. In that way, we make sure that the learner exercises the right generic skill, here simulating a process, while working on the specific knowledge domain, thus building specific domain knowledge and meta-knowledge at the same time.

Other components of a generic skill model can also help choose the kind of learner assistance activities. For example, generic execution principles of procedures in the model could help a facilitator guide some learners who have difficulties with corresponding activities in the LD scenario.

### III.    AN ONTOLOGY OF LEARNING DESIGNS

A taxonomy of generic skills is here presented. It provides the backbone for an ontology of learning designs and, later on a set of metadata element for a repository of learning designs.

*A. The taxonomy of generic skills*

The taxonomy shown on Table 1 has matured through a long experimental process combined with Instructional Engineering tool building. A first version, close to Bloom's taxonomy was elaborate in 1993, integrated in the AGD instructional design support system and expanded within the various versions of the MISA method [18,19,20]. It has been used with experts in many organizations and was field-tested in various applications, in particular to build a complete program for professional training.

TABLE I.    GENERIC SKILLS' TAXONOMY

| Generic Skills Taxonomy Layers | | |
|---|---|---|
| 1 | 2 | 3 |
| Receive | 1. Pay Attention | |
| | 2. Integrate | 2.1 Identify<br>2.2 Memorize |
| Reproduce | 3. Instantiate / Specify | 3.1 Illustrate<br>3.2 Discriminate<br>3.3 Explicitate |
| | 4. Transpose/ Translate | |
| | 5. Apply | 5.1 Use<br>5.2 Simulate |
| Create | 6. Analyze | 6.1 Deduce<br>6.2 Classify<br>6.3 Predict<br>6.4 Diagnose |
| | 7. Repair | |
| | 8. Synthesize | 8.1 Induce<br>8.2 Plan<br>8.3 Model/ Construct |
| Re-invest | 9. Evaluate | |
| | 10. Self-manage | 10.1 Influence<br>10.2 Self-control |

The taxonomy expands on three layers from left to right, from the more general skills to more specific skills. The first layer groups four general information processing processes. The Second layer includes ten generic skills that can be found in educational objective taxonomies or software engineering problem types like those in KADS [4, 8], The third layer corresponds to more specialized skills that are widely used in instructional design methodology. We have described in [19] each of the generic skills in this taxonomy by its inputs and its products and by a detailed generic process such as the one presented on figure 1 for the *5.2 simulate* process.

It is possible to extend this specialization hierarchy to more layers. From layer to layer, we get more and more specialized generic skills until every aspect is totally instantiated in a particular application domain. If we go down the following chain:  Create – Analyze – Diagnose – Diagnose a biological problem – Diagnose a human heart problem, at a certain point,

the skill is no more generic. It becomes procedural knowledge within a specific knowledge domain.

The question whether generic skills are ordered from simple to complex is a delicate one. For generic skills in the first layer, it is quite straightforward: reception skills involves only attention and memory operations, needed at the other levels as well; reproduction skills are essentially instantiation operations from more general knowledge; creation skills produce new knowledge from more specialized ones, involving also some reproduction operations as components; and, finally, re-investment skills involve the explicit use of meta-concepts to evaluate and control the use of knowledge, thus embedding all other reception, reproduction and creation skills as components.

A simple-to-complex relation between two generic skills can be defined in the following way: A generic skill A is aid to be simpler than a generic skill B if the generic process representing A appears as a sub-process or an operation within the model of the generic process B. For example, to synthesize knowledge, we need to analyze components several times before we can combine them in creative ways, so analysis is simpler than synthesis according to this definition.

We have constructed process graphs for all the generic skills in layer 2, providing a validation of their simple-to-complex ordering. In addition, experimental studies with learners made by Martin and Briggs [15] have also validated partly this ordering hypothesis in the case of taxonomies for learning objectives in the cognitive domain [3] and in the emotional domain [13]. Both are embedded the generic skills' taxonomy on Table 1.

*B. Backbone of a learning design ontology*

Usign the MOT modeling software, we present the backbone of learning design ontology, corresponding to the generic skills' taxonomy presented on table 1.

The nodes of the model on figure 2 represent classes of learning designs. The specialization links (in black) between LD classes goes from right to left , while the simple-to-complex links go from top to bottom. We do not extend the simple-to-complex relation at the third layer of the generic skill's taxonomy because there is no evidence that would validate this relation. For example, a classification LD is neither simpler nor more complex than a diagnosis LD; both are just two brands of analysis LD.

*C. Specializing further the generic skills and LD templates*

The MOT software enables the association of submodels to any object in a graph such as the one on figure 2. Figure 3 shows such a submodel for the diagnosis LD class. This can be done also for any node in the graph on figure 2.

The graph on figure 3 presents three sets of diagnosis subclasses, each presenting orthogonal ways to specialize this LD class: according to the knowledge type, the required performance level and types of inputs in the process.

Fig. 2.   Specialization and complexity links in the LD taxonomy



Fig. 3.   Sub-classes of a LD class: the Diagnosis LD template

### D.  Knowledge type

The set of alternative diagnosis sub-classes on the right side of figure 3 refers to the type of knowledge objects in the application domains that are to be processed by the generic skill. For example, diagnosis processes for a system of components, or for a procedure, or for a rule-based system are all specializations of the general diagnosis process. The first one will attempt to generate and test components of the system; the second one will verify if the procedure outputs the expected results for different inputs, and the third one will check if the rules are consistent and cover all main cases. Theses situations refer to taxonomy of knowledge types and models such as the one presented in [23].

### E.  Performance level

The upper set of diagnosis sub-classes on figure 4 refers to the performance level to be deployed when performing a diagnosis. Combining attributes like reliability, complexity, autonomy and familiarity, we define performance class such as the ones on the figure. For example, if a person can perform a diagnosis only on certain occasions (unreliability), in simple and familiar cases, and with outside support, the performance level could be classified as very low skilled diagnosis or simple diagnosis awareness. The LD template should be modeled accordingly. At the other end of the spectrum, if the process is performed in a reliable way, autonomously, in complex as well as new situations, then the diagnose skill should be modeled as "expert diagnosis"

### F.  Cognitive, affective, social and psychomotor meta-domains

The left set of diagnosis sub-classes is based on the types of inputs and outputs of the diagnosis process. Are they of a cognitive, affective, social or psychomotor nature? These four meta-domains provide another way to specialize a generic skill or a LD template. This deserves more detailed explanations.

For example, the generic diagnosis process takes a component model of a situation as its input and returns a list of faulty elements. The diagnosis proceeds in such a way that a component is selected to check if it contains a faulty component. Then this component is decomposed into sub-components to find other hypothesis for faulty components. Each component is tested; its attributes are compared to some norm. If the component is faulty, it is added to the list; if not, a new hypothesis is generated and tested.

This generic process can be applied in various application domains, for example to diagnose a hardware system or to diagnose an emotional situation in a group. The difference between the two applications is the nature of the input and output of the diagnosis process. In the first case, it is applied to a model of a hardware system with components being pieces of equipment down to very small parts that can be deficient; the output is a list of faulty parts. In the second case, it is applied to an affective situation where the components are facts and opinions that people have expressed on a certain event that has caused guilt to occur in a person; the results can be acts that should not have been made, or opinions that are clearly misled.

In [19], we have built a complete table showing examples in the cognitive, affective, social and psychomotor meta-

domains for each of the 10 generic skills on the second layer of our taxonomy. It shows that this taxonomy can be used in each of the four meta-domains.

This work underlines that generic skills and LD templates are characterized by the operations they perform on some input, rather than according to the type of stimulus or response, whether they are cognitive acts, motor actions, affective or social attitudes. We are not claiming here any psychological theory on knowledge, skill and attitudes. Only that at a certain abstraction level, the same generic process can be applied to knowledge, action or socio-affective attitudes. We are claiming operational usefulness, not psychological truth.

We propose that in instructional engineering, it is important to integrate the meta-domains in the same framework. As underlined by Martin and Briggs *"This subdivision (between meta-domains) is relatively arbitrary because the psychologists and the educators agree that, in the reality of educational practice, no real separation between the cognitive, emotional and psychomotor states is possible.* [15] *"* Martin and Briggs quote in support to this assertion several other authors, notably some having produced important taxonomies for educational objectives such as Bloom [3] and Gagné [7].

## IV. METADATA FOR A REPOSITORY OF LD OBJECTS

Competencies provide the backbone for a learning design ontology templates and examples. They also provide a set of key metadata elements to describe LD objects. These metadata elements enable the construction of LD repositories where LD classes and subclasses, as well as instances (concrete learning scenarios) can be stored, search and retrieved to support the instructional engineering processes. Malone et al. [14] have proposed a similar repository for business processes.

### A. Repositories of LD objects

While working on documents to support the use of Educational Modeling Languages and the IMS Learning Design specification [10], we have underlined that "to support the reusability of good learning designs, it is essential that libraries of learning designs be made available as learning objects in one or more repositories" [21].

We proposed that the learning object repositories under construction in different countries should distinguish between "content object", "tool objects" and "process objects", the latter including generic and specific learning designs (or scenarios). Then a growing library of these learning designs could be reused and instantiated to particular knowledge domains. New learning design templates could be built by abstracting generic processes from the large body of existing scenarios, describing the learning material by generic principles or generic descriptions of learning resources, situating the resulting abstraction in the framework of a LD ontology.

We have applied these principles to build some learning design repositories. The first one was built using the IEEE Learning object metadata (LOM) specification. Figure 4 shows part of the repository in the PALOMA metadata editor

and research tool. It shows a list of folders grouping learning designs, one of them selected with its LOM record on the rightmost side of the figure.

A display of the taxonomy on figure 4 has been integrated in the LOM classification field enabling users to choose which generic skill could best describes the learning design. Other classifications can be included to select the type of knowledge, the meta-domain or the instructional strategy or delivery mode. Once learning designs have been referenced with such metadata, the repository can be queried with the classification entries in different ways to find LD objects according to the corresponding attributes.



Fig. 4. A LD repository – LD referenced using skill's Taxonomy

Such a repository can be used to register and retreive LDs using the LD taxonomy metadata. Figure 5 shows a class of LDs and eleven instances that were produced and referenced in the repository during a decomposition/aggregation process. This process was applied to an existing course on Artificial Intelligence labelled Inf-5100. The numbers on the figure show the order of operations in the process.

- (1) The INF-5100 course (level 3 in the LOM specification) was first modeled, referenced and integrated in the LD repository.

- (2) Using the MOT+LD graphic editor, the model was stripped of its content by deleting all content resources to obtain a level 3 pattern , also added to the repository.

- This pattern was then decomposed into five level 2 modular patterns, each added to the repository.

- (4) Using these level 2 patterns as activity structures, a new level 3 pattern (Course X) was aggregated and added to the repository.

- (5) Content items have been added to this level 3 pattern to obtain a new level 3 course in a completely different subject, such as political science, but using the same scenario structure as the initial course.

Fig. 5.   A set on LDs in a Learning Design repository.

The interesting thing in this process is the evolution of the generic skill's references for the various LDs. Looking at the 5 level-2 patterns at the bottom of figure 5, we have the first one aiming at a simple information RECEIVING skill. The next two are at the REPRODUCTION level in the taxonomy. The fourth one embodies a synthesis PRODUCTION skills, while the last one involves META-COGNITION skills. We observe that the initial IA course and the resulting political science courses both involve a progression of generic skill levels through the five level 2 modules. This seems to be a sound cognitive progression strategy.

### B.  RDFS Vocabulary for LD referencing in ISO-MLR

While the above repository had many interesting applications like this one, this project revealed many limitations of the LOM specification when it comes to referencing LD objects.

First the LOM classification for the types of learning resources mix all kinds of things such as exams, questionnaires, reference documents, tutorials and learning activities. This last concept is confused with the notion of instructional scenario or learning design. Second, the LOM section 7 provides 6 kinds of relationships between resources in general.  Some are useful for LDs (such as "has part"), others are two general (such as "is baiss for" or "requires"), while others are not really useful (such as "is based on"). We need more precise relationships based on the very structure of a learning design. Third, and more important, the LOM section 9, while providing the use of cutom-made classification had to be extended to include various taxonomies for competency, pedagogical strategy, delivery modes, évaluation modes, reusability criterias, among others. These taxonomies are not related. They do not form a LD ontology that could provide more expressivity and retrieval power.

We need a more flexible aproach based on the new developpements that are occuring in the Semantic Web [1,6], more precisely on the Web of Linked Data [1,8]. As proposed by the W3C, the basic structure of the Web of linked data is the RDF, the Resource description framework.

This basic RDF framework has been adopted by the International Standards Organization (ISO) and the International Electro technical Commission for their new standard ISO/IEC 19788 Metadata for Learning Resources, in short ISO-MLR [11]. The new standard is intended to provide optimal compatibility with both DC (Dublin Core) and the IEEE-LOM. It supports multilingual and cultural adaptability requirements from a global perspective.

We have begun the development of a LD (or scenario) referencing vocabulary aiming to extend ISO-MLR part 5, the "Pedagogical Elements" part of the standard. The vocabulary is described as a RDF schema (RDFS) vocabulary that provides the required scope and flexibility for a LD lightweight ontology. The SCEN vocabulary has three main parts: the LD concept, the LD taxonomies and the LD context.

The LD concept is shown on figure 6. A LD scenario is composed of the actors, the activities and the resources that compose the scenario. The structure of a LD scenario is defined by three RDF properties (not shown on the figure: first its URL, providing its location on the Web; second, its general structure, a choice between values such as free list of activities, sequence, hierarchy or network; and third, the format of its description that can be a narrative text, a template, a graph or a standard SCORM or IMS-LD manifest.



Fig. 6.   The learning design scenario components

The LD taxonomy part includes the cognitive strategy presented on figure 7. It includes the generic skill in the LD scenario, the knowledge type and the meta-domain.



Fig. 7.   A taxonomy and two properties for a LD's cognitive strategy.

Other taxonomies in the RDFS model are provided for the pedagogical strategy in a LD scenario, its collaborative mode, its delivery mode, its evaluation model and its reusability on the technical, content, context and accessibility dimensions. Contrary to the previous LOM model, these taxonomies are linked together to provide more search capability.

Finally, the LD context includes the intended audience, the scope (program, course or module) and the possible relationships to other scenarios, as shown on figure 8.



Fig. 8. LD context components: audience, scope and relationships

We retain 5 types of relationships between LD scenarios:

- *contains* (ex: the relation between course and module scenarios), which is transitive;

- *hasSchema* (ex: the relation between course and course pattern), which is functional;

- *isFormatOf* (ex: the relation between a LD graph and its narrative description), which is transitive and symmetric;

- *isVersionOf* (the relation between successive representations), which is transitive and symmetric;

- *hasPrerequisite* or pedagogical ordering, for example based on generic skills or knowledge content, which is transitive;

All the LD classes, the object properties and the data properties in theses models need to be precisely described. In order to improve search operations on the Web of linked data it is good practice to link these vocabulary elements to already described vocabularies such as ISO-MLR, DC or FOAF. Table

II provides such associations for part of the class elements on figure 6, 7 and 8. Other similar tables have been built to provide similar information for all scen classes and properties.

TABLE II.  A SAMPLE OF VOCABULARY ELEMENT DESCRIPTIONS

| Class Name | Sub-class Of | Definition |
|---|---|---|
| scen:Scenario | ISO_IEC_19788-1:2010::RC0002  dcmi :Collection  dct : MethodOfInstruction  foaf :Document | A scenario in a learning resource as defined in iso-mlr-1. It is also a Dublin Core collection and a method of instruction. Finally, it is a FOAF document |
| scen:Actor | ISO_IEC_19788-1:2010::RC0002  dct:Agent | An actor is also a learning resource in iso-mlr-1 It is also a Dublin Core agent, which is a person, an organization, or a software agent. |
| scen:Activity | dcmi:event | An activity is a Dublin Core event, that is a non-persistent event situated in time for duration. |
| scen:Resource | ISO_IEC_19788-1:2010::RC0002  foaf :document | A resource in a scenario is a iso-mlr-1 learning resources, as well as foaf:document |
| scen:Audience | ISO_IEC_19788-5:2010::RC0002  dcmi : Groupe | The audience of a scenario is a sub-class of the class mlr5:audience and also a group according to Dublin Core. |
| scen:Learning Structure | ISO_IEC_19788-5:2010::RC0003 | A learning structure is a learning program, a course or a learning unit that refines the class mlr5 :curriculum(fr) |
| scen:Cognitive Strategy | rdf:class | This is a generic auxiliary class that contains the values for generic skill, knowledge type and meta-domain. |

## V.  CONCLUSION

Reusable learning designs and LD repositories is large-scale initiative based on the important task of structuring important parts of meta-knowledge, the knowledge that applies to knowledge, more precisely the properties of generic skills that apply to various knowledge domains.

The graphic language and the few examples presented here have aimed at demonstrating some of the complex interrelations between generic skill's processes and specific domain knowledge. These associations seem a promising approach because they root learning designs in the rich relationship between specific domain knowledge and meta-knowledge. The most stimulating aspect of a generic skills taxonomy built at the meta-knowledge level is the opportunity

it provides to create an expandable and adaptable set of visual models to help solve that huge puzzle of learning environment engineering.

The implementation of these ideas into the Web of linked data still presents important challenges, especially on the usability dimension. Some use cases are promising but they will have to be thoroughly tested before we can claim sound and practical results have been achieved.

### REFERENCES

[1] Allemang D. and Hendler J. (2011) Semantic Web for the Working Ontologist – Effective Modeling in RDFS and OWL. 2nd Edition. Morgan-Kaufmann/Elsevier, Amsterdam.

[2] Bélisle C. et Linard M. (1996) Quelles nouvelles compétences des acteurs de la formation dans le contexte des TIC? Éducation permanente, no 127, 1996-2, pp. 19-47

[3] Bloom, Benjamin S. (1975) Taxonomy of Educational Objectives: the Classification of Educational Goals. New York: D. McKay.

[4] Breuker J. and Van de Velde W. CommonKads Library for Expertise Modelling. (1994) IOS Press, Amsterdam, 360 pages.

[5] Chandrasekaran B. (1987) Towards a Functional Architecture for Intelligence Based on Generic Information Processing Tasks. IJCAI-87 Proceedings, Milan Italy, pp 1183-1192

[6] Domingue, J., Fensel, D. et Hendler, J. A. (dir.). (2011). Handbook of semantic web technologies. Berlin, Allemagne : Springer-Verlag.

[7] Gagné, R. M. (1970) The conditions of learning (2nd ed.) New York, Holt, Rhinehart & Winston,

[8] Heath, T. et Bizer, C. (2011). Linked data: Evolving the web into a global data space. In Synthesis Lectures on the Semantic Web: Theory and Technology, 1(1), 1-136.

[9] IMS-RDCEO (2002) IMS Reusable Definition of Competency or Educational Objective – Best Practice (2002). 1.0 Final Specification, IMS Global Learning Consortium, Inc. Revision: 25 October 2002.

[10] IMS-LD (2003) IMS Learning Design. Information Model, Best Practice and Implementation Guide, Binding document, Schemas. Retrieved October 3, 2003, from http://www.imsglobal.org/learningdesign/index.cfm

[11] ISO-MLR (2013) ISO-IED 19788 Information technoogy – Learning, education and training – Metatda for learning resources multipart standard. http://en.wikipedia.org/wiki/ISO/IEC_19788

[12] LeBoterf G. L'ingénierie des compétences (2ème édition). Éditions d'organisation, Paris, 445 pages, 1999

[13] Krathwohl D.R., Bloom, B.S., and Masia, B.B. (1964) Taxonomy of educational objectives : The classification of educational goals. Handbook II: Affective domain. New York: Longman, 1964

[14] Malone T.W., K. Crowston, J. Lee, B. Pentland, C. Dellarocas, G. Wyner, J. Quimby, C. S. Osborn, A. Bernstein, G. Herman, M. Klein and E. O'Donnell. (1999) Tools for inventing organizations: Toward a handbook of organizational processes, Management Science 45(3) pp 425-443, March, 1999.

[15] Martin B.L. & Briggs L. (1986) The Affective and Cognitive Domains: Integration for Instruction and Research. Educational Technology Publications, New Jersey, 494 pages, 1986

[16] McDermott J.. (1988) Preliminary steps towards a taxonomy of problem-solving methods. In Marcus, S. (Ed), Automating Knowledge Acquisition for Expert Systems, pp. 225-255. Kluwer Academic Publishers, Boston, Mass.

[17] Merrill M.D. Principles of Instructionnal Design. Educationnal Technology Publications, Englewood Cliffs, New Jersey, 465 pages.

[18] Paquette, G. (2001) TeleLearning Systems Engineering – Towards a new ISD model, Journal of Structural Learning 14, 4pp. 319-354, 2001

[19] Paquette, G. (2002) Modélisation des connaissances et des compétences, un langage graphique pour concevoir et apprendre, 357 pages, Presses de l'Université du Québec, mai 2002.

[20] Paquette, G. (2003) Instructional Engineering for Network-Based Learning. Pfeiffer/Wiley Publishing Co, 262 pages.

[21] Paquette G., O. Marino, I. De la Teja, K. Lundgren-Cayrol, M. Léonard, and J. Contamines (2005) Implementation and Deployment of the IMS Learning Design Specification, submitted to the Canadian Journal of Learning Technologies (CJLT), http://www.cjlt.ca/

[22] Paquette G. (2007) An Ontology and a Software Framework for Competency Modeling and Management. Educational Technology and Society, Special Issue on "Advanced Technologies for Life-Long Learning", Volume 10, Issue 3, 2007 pp. 1-21

[23] Paquette G. (2010) Ontology-Based Educational Modelling - Making IMS-LD Visual, Technology, Instruction, Cognition and Learning , Vol.7, Number 3-4, pp.263-296, Old City Publishing, Inc.

[24] Jacques Pitrat J. (1991). Métaconnaissance, avenir de l'Intelligence Artificielle. Hermès, Paris, 1991.

[25] Pitrat J. (1993) Penser l'informatique autrement. Hermès, Paris, 1993.

[26] Popper K. R. (1967) The Logic of Scientific Discovery, Harper Torchbooks, New York, 1967..

[27] Romiszowski A. J. (1981) Designing Instructional Systems Kogan Page London/Nichols Publising, New York, 415 pages.

[28] Schreiber G., Wielinga B., Breuker J. (1993) KADS – A Principled Approach to Knowledge-based System Development. San Diego : Academic Press. 457 p.

[29] Steels L. (1990) Components of Expertise. AI Magazine, vol. 11, no 2, Summer 1990.

[30] Thayse, A. (1988) Approche logique de l'intelligence artificielle, Dunod, Paris, 1988

# Secure Undeniable Threshold Proxy Signature Scheme

Sattar J. Aboud

Department of Computer Science,
University of Bedfordshire,
UK

*Abstract*—**The threshold proxy signature scheme allows the original signer to delegate a signature authority to the proxy group to cooperatively sign message on behalf of an original signer. In this paper, we propose a new scheme which includes the features and benefits of the RSA scheme. Also, we will evaluate the security of undeniable threshold proxy signature scheme with known signers. We find that the existing threshold proxy scheme is insecure against the original signer forgery. In this paper, we show the cryptanalysis of an existed scheme. Additional, we propose the secure, undeniable and known signers threshold proxy signature scheme which answers the drawback of an existed scheme. We also demonstrate that a threshold proxy signature suffers from a conspiracy of an original signer and a secret share dealer, that the scheme is commonly forgeable, and cannot offer undeniable. We claim that the proposed scheme offers the undeniable characteristic.**

*Keywords—cryptography; digital signature; proxy signature; threshold proxy signature*

## I. INTRODUCTION

The proxy signature scheme is a method which allows original signer delegates his works to a designated person with a proxy signature key. The proxy signature key is generated by the original signer signature key which cannot be computed from the proxy signature key. The proxy signer can generate the proxy signature in a message on behalf of an original signer. Since Mambo *et al.* presented an idea of a proxy signature [1], various proxy signature schemes are suggested [2]. Based-on the type of delegation, a proxy signature is categorized into full delegation, partial delegation and delegation by warrant.

In full delegation, an original signer passes its private key as a proxy signature key to a proxy signer over a secure channel. In partial delegation, a proxy signer has the proxy signature key from a proxy signer secret key and a delegation key passed by an original singer. A delegation key is created by an original with the trap-door permutation of an original signer secret key. A proxy signature is dissimilar from an original and a proxy typical signature. In delegation by certificate, an original signer employs its typical signature to sign the warrant that records a kind of information delegated, an original signer and a proxy signer identities and a period of delegation. The signature of a warrant is a certificate that stops a passing of proxy power to a trusted authority.

The partial delegation can be altered into the partial delegation by warrant. A partial delegation by warrant can

offer sufficient security and efficiency. For simplicity, we denote that a partial delegation by warrant a proxy signature. Mambo *et al*. proxy signature scheme satisfy a characteristic of no one except an original signer and a proxy signer can generate the valid proxy signature on behalf of an original signer. In 2001, Lee *et al*. [3] enhanced a security characteristic of a proxy signature by create a valid proxy signature and someone else, even an original signer, cannot create a valid proxy signature. So, for a valid proxy signature, a proxy signer cannot repudiate signed a message and an original signer cannot repudiate delegated a signing authority to a proxy signer. Namely, a proxy signature scheme has a security characteristic of undeniable.

The present proxy signature systems have two drawbacks. First, a declaration of the valid delegation in a warrant is not practical since a proxy signer can generate the proxy signature and claim that the signing was released through a delegation phase. Second, even if a signer key is compromised and the delegated rights are misused; and an original signer needs to revoke a delegation before his strategy, he can make anything. Therefore, a revocation of delegated rights is the important matter of a proxy signature system. To solve the above difficulties, some proxy signature systems have been suggested. Sun indicated that time-stamp proxy signature system and its enhancement [4]. But Sun scheme cannot solve the second drawback. Seo *et al*. [5] suggested a proxy signature system to solve a fast revocation difficulty. The scheme uses the third trusted entity, entitled Security Mediator which is the online partially trusted server.

## II. RELATED WORKS

Based on a Shamir secret sharing scheme in 1979 [6].Zhang *et al*., in 1997 suggested a threshold proxy signature scheme [7]. In their scheme, the proxy signature key is shared among a subset of *n* proxy signers where at least *t* proxy signers can cooperatively sign documents on behalf of an original signer. To avoid argument regarding who is a proxy signer, Sun in 1999 [8] suggested the undeniable threshold proxy signature scheme with known signers. Sun scheme reduces Kim *et al*. scheme [2] drawbacks that a verifier is incapable to verify if a proxy group key is created by an authorized proxy group. In 2001 Hsu *et al*. [9] illustrated that Sun scheme is weak since any *t* proxy signers can get the private keys of other proxy signers. In 2003, Yang *et al*. [10] proposed an enhancement on Hsu *et al*. scheme. Yang *et al*. scheme is more efficient regarding the communication cost

and timing complexity. In 2004, Tzeng *et al.* [11] found that Hwang *et al.* scheme; malicious original signer can forge a threshold proxy signature without an agreement of the proxy signers. Tzeng *et al.* also built the undeniable threshold proxy signature scheme with known signers and claimed the suggested scheme enhanced a security of Hwang *et al.* scheme. In 2006, Yuan Yumin [12] introduced a threshold proxy signature scheme with non-repudiation and anonymity. Yuan Yumin claims that the scheme with any verifier can check if authors of a proxy signature belong to designated proxy group by the original signer, while outsiders cannot find the actual signers. In 2007, Qi Xie *et al.*, [13] claims that their scheme made an improvement of undeniable threshold multi-proxy threshold scheme with shared verification. In 2009, Hu and Zhang [14] presented a cryptanalysis and improvement of a threshold proxy signature scheme with undeniable. In 2012, Hwang *et al.,* proposed a scheme and claimed that its scheme eliminate the security leaks. But, in its scheme the improvement, a malicious original or proxy signer can forge a valid threshold proxy signature for any message by different ways. In this paper, we show the vulnerabilities of the Hwang *et al.,* scheme and proposed a new system that solves the existed problems.

The remainder of this paper is organized as follows. In Section 3, we will provide some notations and reconsider Pedersen threshold distributed key generation protocol [15]. In Section 4, we will analysis a security of Sun *et al.* threshold proxy signature scheme. In Section 5 we will describe the proposed scheme. Finally, conclusions are in Section 6.

### III. PRELIMINARIES

In this Section, we will provide some notations used by this paper and also reconsider Pedersen threshold distributed key generation scheme.

#### A. Notaions Used

In this section, we provide the notations which are used by this paper.

$p, q$ : Two large prime numbers where $q / p - 1$.

$g$ :    Generator of $Z_p^*$ its order is $q$

$O$ :    Original signer

$P_1, P_2, ..., P_n$ : The $n$ proxy signer

$d_O$ :    Private Key of an original singer $O$

$e_O$ :    Public key of an original signer $O$

$d_i$ :    Private Key of a proxy signer $P_i$

$e_i$ :    Public key of a proxy signer $P_i$

$h(.)$ :    Secure hash function.

$||$ :    Concatenation operation

$id$ :    The identity of the proxy signer

$m_w$ : A warrant which records information delegated an original signer and proxy signer.

#### B. Penderson Threshold Distributed Key Generation Protocol

Pedersen threshold distributed key generation scheme contains $n$ Feldman $(t, n)$ verifiable secret sharing schemes [16]. Suppose $(P_1, P_2, ..., P_n)$ are $n$ players. Pedersen scheme includes the following three stages.

*1) Every player $P_i$ arbitrarily selects a polynomial $f_i(z)$ over $Z_q$ of degree $t - 1$.*

$$f_i(z) = a_{i0} + a_{i1}z + a_{i2}z^2 + ... + a_{i,t-1}z^{t-1} \qquad (1)$$

$P_i$ Transmit $b^{a_{i0}}, b^{a_{i1}}, ..., b^{a_{i,t-1}}$. Then finds and passes $f_i(j) \bmod q$ to $P_j$ such that $j = 1, 2, ..., n$ where $j \neq i$ in the secure channel.

*2) Every $P_j$ check a validity of a share $f_i(j) \bmod q$ by verifying for $i = 1, 2, ..., n$ ,*

$$b^{f_i(j)} = b^{a_{i0}} (b^{a_{i1}})^j (b^{a_{i2}})^{j^2} ... (b^{a_{i,t-1}})^{j^{t-1}} \bmod p$$

When all $f_i(j)$ are checked to be certified, $P_j$ finds

$$x_j = \sum_{i=1}^{n} f_i(j) \bmod q \text{ as his share.}$$

*3) Assume* $f(z) = a_0 + a_1 z + a_2 z^2 + ... a_{t-1} z^{t-1} \bmod q$

$$= \sum_{i=1}^{n} f_i(z) \bmod q \quad . \text{ Where, } a_r = \sum_{i=1}^{n} a_{ir} \bmod q \qquad \text{for}$$

$0 \leq r \leq t - 1$ , and $x_i = f(i) \bmod q$  so $w = \sum_{i=1}^{n} x_i \bmod q$

when any $t$ secret shares, say $w_1, w_2, ..., w_t$ are Lagrange interpolating polynomial:

$$w = f(0) = \sum_{i=1}^{i=t-1} s_i \prod_{j=1, j\neq i}^{t-1} \frac{0-j}{i-j} \bmod q \qquad (2)$$

The validity of reconstructed private key $w$ can be checked by the following formula holds: $b^w = \prod_{i=1}^{i=n} b^{a_{i0}} \bmod p$

(3)

### IV. SIGNATURE OF THRESHOLD PROXY SIGNATURE SCHEME

We will describe two threshold proxy signature schemes which are follows:

**Sun Scheme**

The first scheme we will describe the Sun scheme as follows:

#### A. Description of Sum Scheme

First, we will describe Sun threshold proxy signature scheme as follows:

**Secret Share Generation Phase**

In this phase, a proxy group $(P_1, P_2, ..., P_n)$ should do the following:

*1) Create a group of private and public key pair* $(w, e_1) \in Z_q^* \times Z_p$.

1. Run Pedersen threshold distributed key generation protocol as described in Section 2.

*2) Every player $P_i$ uses*
$f_i(z) = d_i + a_{i0} + a_{i1}z + a_{i2}z^2 + ... + a_{i,t-1}z^{t-1}$

*3) The private key shared by a proxy group is* $w = \sum_{i=1}^{n} d_i$

*4) The related public key is* $e_i = \prod e_i \bmod p$.

*5) Gets a secret key share* $x_i = f(i) = \sum_{j=1}^{t} f_j(i) \bmod q$.

*6) Declare* $u_j = b^{g_j} \bmod p, j = 1, 2, ..., t$.

**Proxy Share Generation Phase**

In this phase, an original signer $O$ creates a proxy share as follows.

Step 1: Original Signer $O$

*1) arbitrarily selects* $r \in Z_q$

*2) find* $l = b^r \bmod p$

*3) Compute proxy* $k = d_O h(m_w \| l) + r \bmod q$.

*4) Allocate a proxy key $k$ between a proxy groups by implementing Feldman scheme.*

*5) Selects an arbitrarily polynomial of degree $t-1$:*
$f^-(z) = k + g_1 z + g_2 z^2 + ... + g_{t-1} z^{t-1} \bmod q$

*6) Finds and privately passes $k_i = f'(i) \bmod q$ to a proxy signer $P_i$ for $i = 1, 2, ..., n$*

*7) Declares $(m_w, l)$ and $v_j = b^{g_j}$ ($j = 1, 2, ..., t-1$)*

Step 2: Proxy Signer $P_i$

*1) Accepts $(k_i, m_w, l)$ when a formula*
$b^{k_i} = e_O^{h(m_w \| l)} l \prod^{t-1} v_i^{i^j} \bmod p$ *correct*

*2) Find $k_i = k_i + x_i h(m_w \| l) \bmod q$ as a proxy share.*

**Proxy Signature Generation Phase**

Suppose that $(P_1, P_2, ..., P_t)$ as an actual proxy group signs a document $m$ as follows:

*1) The $t$ proxy signer runs Pedersen threshold distributed key generation protocol for sharing value $c_O = \sum c_{i,O}$ using*
$f_i''(z) = (c_{i,O} + d_i) + c_{i,1}z + c_{i,2}z^2 + ... + c_{i,t-1}z^{t-1} \bmod q$

*2) Each $P_i$ for $i = 1, 2, ..., t$ gets the public key $y = b^{c_O} \bmod p$ and a private arbitrary value share*
$x_i' = f''(i) = \sum d_i + c_O + c_1 i + c_2 i^2 + ... + c_{t-1} i^{t-1} \bmod q$ *such that $c_j = \sum c_{ij}$ for $1 \leq j \leq t-1$*

*3) Each $P_i$ finds proxy signature share*
$s_i = x_i' y + k_i h(id \| m) \bmod q$

*4) Pass $s_i$ to proxy signers $P_j = (j = 1, 2, ..., t, j \neq i)$ in the secure channel.*

*5) Each $P_j$ can check a validity of $s_i$ by verifying when the following formula correct:*

$$b^{s_i} = \left[ e \left( \prod_{j=1}^{t-1} c_j^{i^j} \right) \left( \prod_{j=1}^{t} e_j \right) \right]^e$$

$$\left[ (l_{e_O}^{h(m_w \| l)} \prod_{j=1}^{t-1} v_j^{i^j}) \left( e_1 \prod_{j=1}^{t-1} u_j^{i^j} \right)^{h(m_w \| l)} \right]^{h(id \| m)} \bmod p$$

*6) Every proxy signer in actual proxy group can creates $s = f''(0)e + [f(0) + f'(0)]h(id \| m)$ by a Lagrange interpolation formula to $s_i$.*

*7) The proxy signature on $m$ is $(m, m_w, l, id, e, s)$.*

**Proxy Signature Verification Phase**

The verifier can identify an original signer and an actual proxy signers from $m_w$, and $id$, and validate a proxy signature by verifying when

$$b^s = \left[ l_{e_O}^{h(m_w \| l)} \prod_{i=1}^{n} e_i \right]^{h(id \| m)} \left( y \prod_{i=1}^{t} e_i \right)^y \bmod p \quad (4)$$

*B. Cryptanalysis of Sun Threshold Proxy Signature Scheme*

In this subsection, we illustrate that Sun scheme is weak against an original signer forgery. Since the malicious original signer can create the proxy signature on every document and claim that any $t$ proxy signers can be actual proxy signers of a proxy signature. Assume a message $m$; an original signer $O$ arbitrarily selects the proxy group (thus, $O$ selects $id$).

*1) Suppose that $O$ imitates proxy signers $(P_1, P_2, ..., P_t)$.*

*2) Then $O$ find $l = (\prod_{t} e_i)^{-1} g^a \bmod p$ where $e = (\prod e_i)^{-1} b^v$, such that $a \in Z_q, v \in Z_q$.*

*3) Then, $O$ finds:*
$s = (a + d_O h(m_w \| l))h(id \| m) + ve \bmod q \quad (5)$

*4) So $(m, m_w, l, id, e, s)$ is the valid proxy signature on message $m$ since*

$$b^s = b^{(a + d_O h(m_w \| l))h(id \| m) + ve} \bmod p$$

$$= b^a b^{d_O h(m_w \| l)h(id \| m)} (b^v)^e \bmod p$$

$$= \left( l_{e_O}^{h(m_w \| l)} \prod_{i=1}^{n} e_i \right)^{h(id \| m)} (y \prod_{i=1}^{t} e_i)^y \bmod p$$

*C. The Vulnerability of Sun Scheme*

With Sun $(t, n)$ threshold proxy signature system, the verifier checks a validity of a proxy signature and recognizes the real signers. Though, in this paragraph we illustrate that a proxy signer private key is not protected. The $(n-1)$ proxy signers in a group of $n$ can present their private keys to conspire a private key of a residue one. We so-call this attack a collusion attack.

In this attack, any $(n-1)$ proxy signers in a group of $n$ participants can masquerade a rest one. For instance, suppose that (3,5) threshold proxy signature system. A proxy signer $p_1,...,p_4$ aims to get a private key of a proxy signer $p_5$. Then, we can masquerade the authorized proxy signer $p_5$ to sign the document $m$. Any three proxy signers of $p_1,...,p_4$ can find $a_0$ using Lagrange equation since $a_0 = \sum_{i=1}^{5} d_i \bmod p$. So, proxy signers $p_1,...,p_4$ can appear the private keys to conspire a private key $d_5$ of a proxy signer $p_5$. We can masquerade a proxy signer $p_5$ to create the authorized proxy signature.

In the same manner $s_5, k_5$ and $k_5'$ is calculated by using Lagrange equation. In proxy signature issuing phase, we can masquerade $p_5$ to share the arbitrary number, and we can obtain the secret $s_5'$. By holding $s_5'$ and $k_5'$, we can get $\gamma_5$ and post it to other proxy signers of a proxy group. Then $T$ can be calculated and then, a proxy group can create the proxy signature $(m, T, l, m_w, id)$ for document $m$.

In a verification phase, a verifier can check a validity of a proxy signature and find $p_5$ as real signer of a proxy group. Actually, $p_5$ has never signed a document $m$, but cannot repudiate. Thus, in Sun scheme, a private key $x_i$ of a proxy signer $p_i$ can be compromised by collusion attack and hacker can masquerade authorized proxy signer $p_i$ to sign the document.

## Hwang *et al.* Scheme

This is the second threshold proxy signature scheme we are going to describe which is the Hwang *et al.* scheme [17] is the same as Sun threshold proxy signature scheme.

### D. Description of Hwang et al. Scheme

First, we will describe Hwang *et al.* threshold proxy signature scheme as follows:

**Secret Share Generation Phase**

In this phase, a proxy group should do the following:

*1) Creates group of private and public key pair* $(w, e_1) \in Z_q \times Z_p$ *as in Sun scheme.*
*2) Finds* $f_i(z) = d_i + a_{i0} + a_{i1}z + a_{i2}z^2 + ... + a_{i,t-1}z^{t-1}$ .
*3) The secret key shared by a proxy group is* $w = \sum_{i=1}^{n} d_i$
*4) The related public key is* $e_i = \prod e_i \bmod p$ .
*5) Gets the secret key share* $x_i = f(i) = \sum_{j=1}^{n} f_j(i) \bmod q$ .
*6) Declares* $u_j = b^{a_j} \bmod p$ *with* $j = 0,1,2,...t-1$ .

**Proxy Share Generation Phase**

In the phase, an original signer $O$ creates a proxy share as follows:

**Step 1: Original Signer** $O$

*1) Creates a proxy key* $k = h(m_w || l)d_0 + r \bmod q$ .

*2) Selects arbitrarily polynomial of degree* $t-1$ : $f'(z) = k + g_1 z + g_2 z^2 + ... + g_{t-1} z^{t-1} \bmod q$
*3) Finds and secretly posts* $k = f'(i) \bmod q$ to $P_i$ for $i = 1,2,...n$ .
*4) Declares* $(m_w, l)$ , *and* $v_j = b^{g_j} \bmod p$ *for* $j = 1,2,...t-1$

**Step 2: Proxy Signer** $P_i$

*1) Uses* $k_i = h(m_w || l)$ *when the following formula holds.* $b^{k_i} = y_0^{h(m_w, l)} l \cdot \prod v_j^{i^j} \bmod p$ .
*2) Finds* $k_{ij} = k_i' + x_i \cdot h(m_w || l) \bmod q$ .

**Proxy Signature Generation Phase**

We suppose $(P_1, P_2,...,P_t)$ are actual proxy group. So, the steps of this phase as follows:

*1) Creates a secret random share* $x_i'$ *as in Sun scheme.*
*2) Finds a single proxy signature* $s_i = x_i' e + k_i' h(id || m) \bmod q$
*3) Posts* $s_i$ *to the proxy signers* $P_j (j = 1,2,...t, j \neq i)$ *in the secure way.*
*4) Checks a validity of* $s_i$ *by verifying when the following formula holds:* $\left[ \left( l_{e_0}^{h(m||l)} \prod^{t-1} v_j^{i^j} \right) \cdot \left( e_1 u_0 \prod^{t-1} u_j^{i^j} \right)^{h(m_w||l)} \right]^{h(id||m)} b^{s_i} = \left[ e \left( \prod^{t-1} c_j^{i^j} \prod^{t} e_j \right) \right] \bmod p$ (6)
*5) Using a Lagrange interpolation equation with* $s_i$ , *every signer can create* $s = f''(0)e + [f(0) + f'(0)]h(id || m)$ .
*6) A proxy signature on* $m$ *is* $(m, m_w, l, id, e, u_0, s)$ .

**Proxy Signature Verification Phase**

*1) Verify a validity of a proxy signature from the following formula:*

$$b^s = \left[ lu_0 e_0^{h(m_w||l)} \prod_{i=1}^{n} e_i \right]^{h(id||m)} \left( e \prod_{i=1}^{t} e_i \right)^e \bmod p \qquad (7)$$

*2) When the formula holds, a proxy signature* $(m, m_w, l, id, e, u_0, s)$ *is valid.*

### E. Cryptanalysis of Hwang et al. Scheme Threshold Proxy Signature Scheme

In this subsection, we illustrate that Hwang *et al.* scheme is insecure versus universally forgery. The hacker can impersonate an original signer to forge the proxy signature on a message. Provided a message, an original signer, and the proxy group $(P_1, P_2,...,P_n)$, a hacker selects $(P_1, P_2,...,P_t)$ as actual proxy signers. Then, a hacker selects four arbitrary integers $a, v, \gamma \in Z_q^*$ and $e \in Z_p^*$. Then a hacker finds

$$l = \left( \prod_{i=1}^{n} e_i \right)^{-1} b^a \bmod p \qquad (8)$$

$$u_0 = (e_0^{h(m_w||l)})^{-1} b^v \bmod p$$

$$s = (a + v)h(id || m) + \gamma e \bmod q \qquad (9)$$

Therefore, $(m, m_w, l, id, e, u_0, s)$ is the valid proxy signature on message $m$, it convinces the following verification formula:

$$b^s = b^{(a+v)h(id\|m)+\gamma e} \bmod p$$

$$= (b^a b^v)^{h(id\|m)} b^{\gamma e} \bmod p$$

$$= \left[ lu_0 e_0^{h(m_w\|l)} \prod_{i=1}^{n} e_i \right]^{h(id\|m)} \left( e \prod_{i=1}^{t} e_i \right)^e \bmod p$$

## V. THE PROPOSED SCHEME

The suggested scheme combines theta $\theta(n)$ and an elimination of a computation of inverse in RSA scheme if we calculate a value of Lagrange coefficient. Also, we suggest an equation to find a result of message warrant $m_w$. Suppose that $N_O < N_i (i = 1,2,...,n)$.

### A. The Proxy Sharing Phase

The steps of the proxy sharing phase are as follows:

**Step 1: Proxy Generation**

The original signer $O$ must do the following:

*1) Find a group proxy signing key* $d_1 = d_O^{m_w} \bmod \theta(N_O)$
*2) Find the proxy verification key* $e_1 = e_O^{m_w} \bmod \theta(N_O)$
*3) Compute* $m_w = (P + T + r)^T \bmod \theta(N_O)$ *such that P is a validity period of proxy signature and T is a sum of identities of* $P_O = P_1, P_2,...,P_n$
*4) Declare* $(m_w, e_1, (m_w, e_1)^{d_O} \bmod N_O$

**Step 2: Proxy Sharing**

The original signer $O$ must do the following:

*1) Choose* $t-1$ *degree polynomial* $f(x) = d_1 + a_1 x + ..a_{t-1} x \bmod N_O$ *with* $a_1, a_2,...,a_{t-1}$, *are an arbitrary integers.*
*2) Compute a proxy singer* $P_i$ *partial proxy signing key* $k_i = f(i)$
*3) Pass* $((k_i)^{d_O} \bmod N_O, k_i)^{e_i} \bmod N_i$ *to proxy signer* $P_i$

**Step3: Proxy Share Generation.**

The proxy signer $P_i$ must do the following:

*1) Receive* $((k_i)^{d_O} \bmod N_O, k_i)^{e_i} \bmod N_i$
*2) Obtain* $((k_i)^{d_O} \bmod N_O, k_i)$ *by his secret key* $d_i$
*3) Verify a validity of* $k_i$ *and keeps it secret.*

### B. The Proxy Signature Issuing Phase

Suppose that $T$ indicate the group members including any $t$ proxy signers who desire to create the proxy signature on a message $m$ on behalf of $P_O$ cooperatively.

**Step 1: Proxy Signer** $P_i$

Every proxy signer $P_i$ uses a partial proxy signing key $k_i$ to do the following:

*1) Create a partial signature* $s_i = m^{k_i} \bmod N_O$
*2) Pass* $((s_i, i)^{d_i} \bmod N_i, s_i$ *to a combiner.*

**Step 2: The Combiner**

The combiner must do the following:

*1) Receive partial signature* $s_i$ *from* $P_i$
*2) Check the validity of a partial proxy signature by verifying if* $(s_i, i)^{d_i e_i} \bmod N_i = (s_i, i)$.
*3) Find* $v = \prod_{id_a, id_b \in T} id_a - id_b$ *such that* $a > b$
*4) Find* $\prod_{j=1, j \neq i}^{t} (id_i - id_j)$ *a factor of* $\prod_{id_a, id_b \in T} (id_a - id_b) \prod \frac{id_j}{id_i - id_j}$ *where* $\prod_{id_a, id_b \in T} (id_a - id_b)$. *Thus,* $L_i$ *are integer and combiner required not calculating inverse of* $\prod_{j=1, j \neq i} (id_i - id_j)$.
*5) Create a signature* $s = \prod s_i^{L_i} \bmod N_O$
*6) The result of proxy signature is* $(v, s)$.

### C. The Proxy Signature Verification Phase

The steps of this phase are as follows:

*1) A verifier can check a signature signed on behalf of an original signer by a formula* $s^{e_1} = m^v \bmod N_O$
*2) An original signer can distinguish a proxy signer from a signature by* $s_i^{d_i e_i} \bmod N_i = s_i$
*3) An original signer can trace proxy signers by* $e_i$.

## VI. COMPARISONS

We compare the running of five schemes, Hwang *et al.* [17], Kim *et al.*[2], Hwang *et al.*[11], Sun *et al* [8] and Hsu *et al* [9] with a performance of the proposed scheme. The proposed scheme is efficient and secure anti-disreputable conspiracy attacks. Table 1 shows a comparison of threshold proxy signature schemes relied on proxy needs every scheme.

TABLE I.       1 THE COMPARISION BETWEEN EXISTED SCHEMES AND PROPOSED SCHEME

| Security Features | Name of the Scheme | | | | |
|---|---|---|---|---|---|
| | *Kim* | *Hwang* | *Sun* | *Hsu* | *Proposed* |
| Proxy Protection | No | Yes | No | No | Yes |
| Unforgeability | Yes | Yes | No | No | Yes |
| undeniable | Yes | No | Yes | Yes | Yes |
| Known Signer | No | Yes | Yes | Yes | Yes |

## VII. CONCLUSION

In this paper, Sun threshold proxy signature scheme has been analysis. The scheme is based on discrete logarithm assumption. The security of Sun is undeniable threshold proxy signature scheme with known signers. We find that in Sun scheme, a malicious original signer can forge a valid proxy signature on any message without the agreement of the proxy group. We also suggest an efficient scheme which involves the characteristics and gains of the RSA cryptosystem which is a popular security scheme.

### References

[1] Mambo M., Usuda K., and Okamoto E., "Proxy Signatures for Delegating Signing Operation", Proceeding of 3rd ACM Conference on Computer and Communications Security, ACM Press, pp. 48-57, 1996.

[2] Kim H., Baek J., Lee B., and. Kim K, "Secrets for Mobile Agent Using Onetime Proxy Signature", Cryptography and Information Security 2001, Volume 2/2, pp. 845-850, 2001.

[3] Lee B., Kim H., and Kim K., "Secure Mobile Agent Using Strong Non-designated Proxy Signature," Proceeding of ACISP 2001, pp. 474-486, 2001.

[4] Sun M., "Design of time-stamped proxy signatures with traceable receivers", IEE Proceedings: Computers and Digital Techniques, 2000, vol. 147, no. 6, pp. 462-466.

[5] Seo S., Shim K., and Lee S., "A mediated proxy signature scheme with fast revocation for electronic transactions", Proceedings of the 2nd International Conference on Trust, Privacy and Security in Digital Business, Aug 22-26, 2005, LNCS 3592, German: Springer, 2005, pp. 216-225, 2005.

[6] Shamir A., "How to Share a Secret", Communications of the ACM, Volume 22, No. 11, pp. 612-613, 1979.

[7] Zhang K, "Threshold Proxy Signature Schemes, "Information Security Workshop", Japan, pp. 191-197, 1997.

[8] Sun H., "An Efficient Nonrepudiable Threshold Proxy Signatures with Known Signers", Computer Communications 22(8), pp. 717-722, 1999.

[9] Hsu C., and T. Wu, "New Nonrepudiable Threshold Proxy Signature Scheme with Known Signers", the Journal of Systems and Software 58(2001), pp. 119-124, 2001.

[10] Yang C., Tzeng S. and M. Hwang, "On the Efficiency of Nonrepudiable Threshold Proxy Signatures with Known Signers", Journal of Systems & Software 22(9), pp. 1-8, 2003.

[11] Tzeng S., Hwang M., and Yang C., "An Improvement of Nonrepudiable Threshold Proxy Signature Scheme with Known Signers", Computers & Security 23, pp. 174-178, 2004.

[12] Yuan Yumin, "A Threshold Proxy Signature Scheme with Non-Repudiation and Anonymity", Computer and Information Sciences-Proceedings of ISCIS 2006, 21st International Symposium, Istambul, Turkey, November 1-3, 2006.

[13] Qi Xie, Jilin Wang and Xiuyuan Yu, "Improvement of Nonrepudable Threshold Multy-Proxy Threshold Multi-Signature Scheme with Shared Verification", Journal of Electronics (China), Volume 24, 2007

[14] Hu, J., Zhang, J., "Cryptanalysis & Improvement of a Threshold Proxy Signature Scheme", Computer Standards & Interfaces, 2009.

[15] Pedersen T., "A Threshold Cryptosystem without Trusted Party", Proceeding of Advance in Cryptology-EUROCRYPTO'91, LNCS 547, Springer-Verlag, pp. 522-526, 1991.

[16] Feldman P., "A Practical Scheme for Non–Interactive Veriable Secret Sharing", Proceeding of 28th FOCS, IEEE, pp. 427-437, 1987.

[17] Hwang M, Lin I, and Lu K, "A Secure Nonrepudiable Threshold Proxy Signature Scheme with Known Signers", International Journal of Informatica, Volume 0, Number 0, 1-0, pp.1-14, 2012.

AUTHORS PROFILE

**Sattar J Aboud** is currently, a Visiting Professor in the Department of Computer at University of Bedfordshire, UK. He received his education from United Kingdom. Dr. Aboud has served his profession in many universities and he awarded the Quality Assurance Certificate of Philadelphia University, Faculty of Information Technology in 2002. Also, he awarded the Medal of Iraqi Council of Representatives for his conducting the first international conference of Iraqi Experts in 2008. His research interests include the areas of both symmetric and asymmetric cryptography, area of verification and validation, performance evaluation and e-payment schemes.

# Reduced Complexity Divide and Conquer Algorithm for Large Scale TSPs

Hoda A. Darwish, Ihab Talkhan
Computer Engineering Dept., Faculty of Engineering
Cairo University
Giza, Egypt

*Abstract*—**The Traveling Salesman Problem (TSP) is the problem of finding the shortest path passing through all given cities while only passing by each city once and finishing at the same starting city. This problem has NP-hard complexity making it extremely impractical to get the most optimal path even for problems as small as 20 cities since the number of permutations becomes too high. Many heuristic methods have been devised to reach "good" solutions in reasonable time. In this paper, we present the idea of utilizing a spatial "geographical" Divide and Conquer technique in conjunction with heuristic TSP algorithms specifically the Nearest Neighbor 2-opt algorithm. We have found that the proposed algorithm has lower complexity than algorithms published in the literature. This comes at a lower accuracy expense of around 9%. It is our belief that the presented approach will be welcomed to the community especially for large problems where a reasonable solution could be reached in a fraction of the time.**

*Keywords—Traveling Salesman Problem; Computational Geometry; Heuristic Algorithms; Divide and Conquer; Hashing; Nearest Neighbor 2-opt Algorithm*

## I. INTRODUCTION

Divide and Conquer is an algorithm method used in search problems. As the search problem increases this method proves to be one of the best in reaching quick solutions; not only does it breakdown the search problem for easier calculations, in some cases it also allows for parallelizing the search hence reaching faster results. It has come to our notice that not many or not enough tries were given to the Divide and Conquer method when it comes to the Traveling Salesman Problem (TSP). The trend in resolving TSP is for Local Search algorithms and Evolutionary algorithms. Most of the research targets enhancing the constraints and fitness functions of these 2 categories of algorithms to reach a better solution. In most cases, these enhancements affect computational complexity making the resulting algorithms unfeasible for large scale problems.

For TSP, eliminating the long paths between any 2 cities/points in advance enables us to find quickly a more optimum solution. By dividing the search space or plane into pieces, we are effectively eliminating the paths between cites at the 2 ends of the search space thus, decreasing the number of paths we need to search. The plane/space is divided into "Buckets" each holding a set of points that are within a specific distance from each other. In the most ideal situation of evenly distributed points, the Heuristic TSP would now need to find the path for $N/b$ points only where $b$ is the number of buckets.

In the case of NN 2-opt, finding the path of $N/b$ points requires a fewer number of iterations to reach a near optimum path and a much shorter run time. Accordingly, it is expected that the computational complexity of the Hashed Bucket algorithm will be of a much lower order of magnitude as we shall see in this paper.

The rest of this paper is organized as follows; Section II outlines the problem we are trying to address. Section III gives a briefing on TSPs and the current algorithms used for their resolution while Section IV presents a literature survey of related work. Section V describes our proposed solution. We then discuss the flow of our system in Section VII. Finally, our experimental results are presented in section VIII.

## II. PROBLEM DEFINITION

The traveling salesman problem asks the following question: Given a list of cities and the distances between each pair of cities, what is the shortest possible path that visits each city exactly once and returns to the origin city? The complexity of such a problem is NP-hard making it extremely unrealistic to solve optimally.

The problem addressed here is how to improve Local Search Algorithms specifically the Nearest Neighbor 2-opt using a spatial Divide and Conquer method to obtain a new hybrid faster Heuristic algorithm. This poses the challenge of deciding the correct search space division and how these space divisions impact the performance of the NN 2-opt.

## III. BACKGROUND

TSP is a very old problem with many references in literature as well as a long standing history. The first instance of the traveling salesman problem was documented by Euler in 1759. Euler wanted to address the problem of moving a knight to every position on a chess board exactly once as explained in [1]. The constraint set by Euler was that the knight must move according to the rules of chess and must visit each square exactly once.

### A. Types of TSP

There are 2 main characteristics of TSPs. Depending on these, the problem representation may use different data structures and different calculations.

*1) Symmetric vs. Asymmetric TSPs:* a symmetric TSP is a problem where the distance from point A to point B is equal the distance from B to A. Asymmetric TSPs is when the distances

from A to B and vice versa are not equal. For example: if we consider the effort needed to go up a hill higher than the effort needed to go down then we have an Asymmetric problem.

*2) Euclidean/Planer vs. 3 Dimensional:* problems that consider only 1 constraint for the distance between the TSP points/cities can be considered planer. The most famous example of that would be the Euclidean distance. Once we start considering geographical distances, time or monetary costs we find that we have more constraints and hence, more dimensions for the TSP problem representation.

### B. Complexity and Optimality

When we assess TSP algorithms, we look into optimality as well as complexity. Complexity is key given the number of permutations needed to calculate all possible paths; as we shall see in for exact algorithms in the following section. Yet in some cases, we can easily get an approximated path so it becomes necessary to measure the optimality of that path. By optimality, we mean how close it is to the real shortest path of the problem.

### C. Exact (Non-Heuristic) TSP Algorithm

Simply put, there is only 1 way of finding the most optimal path for a TSP: comparing all possible paths and picking the shortest. Unfortunately, this Brute Force seraph method is not realistic as it means we must enumerate all possible permutations for the points in the TSP. In other words, for a problem with N points, we would need to look through $N! - N$ Factorial – possible solutions. For example, a simple problem of 10 points would require passing through (and comparing) 3,628,800 possible paths. Such an approach would be impractical in real world situations where we would need to solve TSP for a huge number of points. The complexity of this approach is $O(n!)$ where n is the number of points. Algorithms with such complexity are called NP-hard. In the case of NP-Hard problems, other means of reaching a solution are required as we shall see in the next section.

### D. Heuristic TSP Algorithms

The traditional lines of attack for an NP-hard problem – when exact optimal methods are unfeasible – are the following:

- Devising "suboptimal" or heuristic algorithms, i.e., algorithms that deliver either seemingly or probably good solutions, but which may prove to be suboptimal.

- Finding special cases for the problem ("sub-problems") for which either better or exact heuristics are possible.

The TSP problem remains NP-hard even for the case when the cities are in the planer Symmetric Euclidean problem. Various heuristics and approximation algorithms have been devised specifically for TSP. Modern methods can find solutions for extremely large problems within a reasonable time and which are quite close to the optimal solution.

There are many types of Heuristic TSPs in the literature. Following is an overview of the main categories:

*1) Tour Construction Algorithms*: these algorithms gradually build a tour by adding a new city at each step. This approach is always quiet simple, but often too greedy. The first distances in the construction process are reasonably short, whereas the distances at the end of the process usually will be rather long. The most popular algorithm in this family is the Nearest Neighbor (NN). NN starts at some random city and then visits the city nearest to the starting city and then keeps visiting the nearest city that has not been visited so far until all cities are covered. It is a poor heuristic with the only simplicity as an advantage so it is normally used for small size problems.

*2) Iterative Local Search (ILS) Algorithms:* these start out with a complete solution at a certain optimality and iteratively try to change the features of the solution until a more optimal solution is found. For TSP, the initial complete solution can be a random tour through the problem with total cost S. The iterative changes would involve exchanging edges or paths between 2 or more cities and comparing the resulting tour of cost S' to S. If S' is a more optimal tour, we start iterating on that. If S' is worse than S, we discard that tour and begin iterating on other city pairs. A stopping criteria must be set in advance so that the algorithm doesn't iterate endlessly on all cases. There are many variations on the ILS, for example:

*a) 2-opt Heuristic algorithm:* this is the most basic of the ILS algorithms:

- Start with a given tour.

- Replace 2 links of the tour with 2 other links in such a way that the new tour length is shorter.

- Repeat until no more improvements are possible.

*b) 3-opt Heuristic Algorithm:*this is the same as the 2-opt but we pick 3 edges or links to replace instead of just 2 edges.

K-opt or Lin–Kernighan Heuristic Algorithm: this a generalization on the 2-opt and 3-opt algorithms that allows k-opt moves. It has many different constraints and modifications in an attempt to improve optimality and complexity. As explained in [2], the original algorithm as implemented by Lin and Kernighan in 1971, had an average running time of order $N^{2.2}$ and was able to find the optimal solutions for most problems with fewer than 100 cities. However, this algorithm is not simple because the number of operations to test all k-exchanges increases rapidly as the number of cities increases. In a naive implementation, the testing of a k-exchange has a time complexity of $O(N^k)$. Furthermore, there is no upper bound of the number of exchanges. Accordingly, the usefulness of general k-opt sub-moves usually depends on the candidate TSP. Unless it is sparse, it will often be too time consuming to choose k larger than 4. Another drawback is that k must be specified in advance and it is difficult to know what k to use to achieve the best compromise between running time and quality of solution. To overcome the drawbacks of the traditional LK algorithm, Lin and Kernighan introduced a powerful variable-opt algorithm: at each iteration, the algorithm examines – for ascending values of k – whether an interchange of k-links may result in a shorter tour. This continues until some stopping conditions are satisfied. Many other variations and enhancements can be found in [3].

*3) Evolutionary Algorithms:* As the name implies, evolutionary algorithms follow nature in an attempt to reach

the best solution for optimization problems. Genetic algorithms (GAs) are one of the most popular evolutionary techniques. Taken from nature, GAs use crossover and mutation to solve optimization problems. GAs are loosely based on natural evolution and use a "survival of the fittest" technique, where the best solutions survive and are varied until we reach a good result. The incorporation of the survival of the fittest idea provides a means of searching the problem space without enumerating every possible solution. A GA works by first 'guessing' a set of solutions and then combining the fittest solutions to create a new generation of solutions which should be better than the previous generation. We may also include a random mutation element to account for the occasional 'mishap' in nature. As [4] explains, the main disadvantages of GAs are premature convergence and poor local search capability. In order to overcome these disadvantages, evolutionary adaptation algorithms based on the working of the immune system have been devised. The interested reader can refer to [5], and [6] for more samples.

## IV. RELATED WORK

Being able to solve large scale TSPs has been of great interest to many. In this section, we give an overview of some proposed solutions and their usefulness for different types of TSP sizes.

*1) Medium Scale TSPs (500 to 3000 points):* In [7], the authors look into solving TSP problems with hybrid, iterative extended crossover operators for GA. The objective of the hybrid algorithm is to efficiently search for the optimum solution while maintaining the diversity of the cyclic paths composing the population. It is a kind of hybrid method which combines Edge Assembly Crossover (EAX) with Ant Colony Optimization. The algorithm was verified on test data of size up to 1173 cities. The optimal path was obtained but required 109 hours to calculate! In fact, the computational time increases exponentially with the increase in number of cities.

*2) Large Scale TSPs (5000+ points):* In [8], the authors consider a k-means partitioning algorithm to divide the initial TSP problem into multiple partitions to be solved separately then merging. The partitioned sub-problems are merged using Lin Kernighan algorithm. To partition the TSP problem, the authors represented the problem as a graph and used multilevel graph partitioning. Multilevel k-means graph partitioning reduces the size of the graph by collapsing vertices and edges as explained by [9]. It divides the graph into smaller graphs and then refines the partition during an "un-coarsen" phase to construct a partition for the original graph. For solving each sub-problem a greedy tour construction heuristic is used to get a good solution of individual small partitions. After solving each partition, step by step recreation of graph is carried out by simply adding each solved partition back to the graph. The algorithm was tested on TSPLIB and provided quite optimal TSP tours but no time complexity was clarified. It is known that the average LK complexity is $O(N^{2.2})$; by clustering and using LK in the coarsening phases of merging the partitions, it is clear that the complexity of this algorithm is definitely more than that of the proposed Divide and Conquer NN 2-opt.

Many other researchers have attempted to enhance the complexity of LK implementations and have reached $O(N^2)$ yet the tradeoff is extra memory of $O(N)$ making it again impractical for large scale problems.

*3) HW Parallelization of Large Scale TSPs:* Given the complexity of TSP algorithms, the speed up and execution time gained from increasing HW resources cannot be expected from normal software solutions so will not be compared with the algorithm proposed in this paper. It does show though that partitioning the problem still allows us to get relatively optimal tour solutions. The authors of ]10[ (2007) introduced the notion of "symmetrical 2-Opt moves" which allowed them to uncover fine-grain parallelism when executing the 2-opt local search optimization algorithm. Once the parallelism is apparent they use an FPGA (or FPGA simulator) to resolve each sub problem gaining an average speed up of 600%.

## V. PROPOSED APPROACH

The proposed algorithm depends on the theory that "the addition of shortest set of paths will yield shortest total path". Accordingly, if we have N points getting the shortest path for N/b points then consolidating the set of b paths will give us the shortest path through the N points. We divide the N points according to their proximity to each other in the search space using x and y dimensions.

For example, a square space of area $A^2$ will be divided into smaller areas called Buckets of area $A^2/b$ where b is the number of Buckets. All points in the same bucket are considered close in proximity and a heuristic TSP algorithm is used to get their shortest path. Given that the number of points in area $A^2/b$ is much less than the total area, the heuristic algorithm has a good chance of finding the optimal path in a much shorter execution time. Once all b paths have been obtained, merging them into a single path for the N points should yield the shortest total path. Fig. 1 shows a simple example where the search space/plane has been divided into 4 buckets.

This approach was inspired by the work done in [11] that is based on the Fixed-Radius Nearest-neighbor problem. The authors of [11] show that bucket hashing is very effective in the domain of electronic design automation specifically in chips of *millions* of transistors as it breaks the problem into manageable pieces for quicker resolution.

Fig. 1. Search Space division using Buckets

## VI. SYSTEM FLOW

The system is comprised of a set of functions that interact with each other. The flow of processing can be seen in the chart in fig. 2. The first step of the process is to parse the input file to get the points that make up the TSP problem. Using the input criteria for bucket size, a decision is made regarding the division of the search space. The following functions then handle the creation of the buckets and hashing the TSP points into the different buckets. Once the buckets are ready, we can then consider each bucket as a separate TSP problem for the regular heuristic NN 2-opt TSP algorithm and thus, the algorithm is run for each bucket. The final step in our flow is to merge the smaller bucket TSP solutions into 1 solution for the original input point thus providing 1 single path and its total cost.

The implementation explained above makes the assumption that all points are connected (in case the TSP doesn't fit this constraint, setting the distance between the unconnected points to infinity should automatically eliminate the path but this theory has not been tested here). We also assume that the input TSP is a Symmetric TSP and that the distance used is Euclidean.

The algorithm depends on finding the minimum path for each bucket and then merging the result. The sum of these local minima may not in fact result in the global minimum. This is one of the disadvantages of the Heuristic algorithms in general yet given that the hashed algorithm complexity is significantly lower we can run the hashed algorithm with different averages or bucket sizes and choose the minimum depending on the original problem size.

### A. Bucket Size Decision

The algorithm complexity depends heavily on the number of buckets used. Accordingly, we need a simple decider to use for dividing our search space. Heuristic TSP algorithms normally have an Average number of points that they can optimally get the shortest path for. Assuming the points are equally distributed, we use this average to divide the space according to the following equations:

*1)* Get the minimum possible number of buckets by dividing number of points on the input average.

*2)* Calculate the search space area A (maximum of x * maximum of y).

*3)* Get bucket width using eq. 1:

$$bucket\ Width = ceil(\frac{(maximum\ of\ x)^2}{minimum\ no.of\ Buckets}) \quad (1)$$

*4)* Get bucket length using eq. 2:

$$bucket\ Length = ceil(\frac{A}{minimum\ no.of\ Buckets * bucketWidth}) \quad (2)$$

### B. Path Merging

The other important step in the algorithm is the merging of the separate bucket solutions to form a single final path. An example of the bucket path merging is shown in fig. 3.

We have Start point "S" for the bucket and a Transition Forward point "TF" for moving to the next bucket in each individual bucket path. When we merge, we remove the path between the TF and the point following it in the bucket; instead we merge it with the start point of the successor bucket. On the way back, we remove the last leg of the path back to the bucket



Fig. 2. System Workflow

Fig. 3.   Path Merging

start point and instead move to the point following the TF in the previous bucket. In other words, we delete the "dashed" red lines and add the solid black lines.

## VII.   RESULTS AND ANALYSIS

The complexity of the simple NN 2-opt algorithm is O ($N^2$) where N is the input number of points of the TSP problem. In the hashed approach, we divide the space into buckets. Assuming that we have "b" buckets and that the points are evenly distributed on the buckets as "N/b" then the complexity of a single bucket is O($(N/b)^2$). To get the complexity of the entire Hashed algorithm we consider that we need to calculate the NN 2-opt for b buckets i.e. a complexity of O ($b*(N/b)^2$) which is equal O ($(N/\sqrt{b})^2$). As we shall see in the results, the above complexity is tangible numerically in the following example. If N = 493 and b = 9, the normal NN 2-opt would require 493² = 243,049 calculations/computations while the Hashed technique would provide $(442)^2/9$ = 27,005 computations which is 11% of the simple algorithm computations.

To show the effectiveness of the proposed algorithm and its ability to solve all general cases, some test cases from TSPLIB were used with focus on large samples. The results obtained are from running the system under Windows 7 and using Matlab 7.0. The system specifications are: Intel Core CPU (1.6GHz) and a 4GB RAM.

We show the test results for 3 different test samples. For each sample, we list the results of the normal NN 2-opt algorithm as well as the Hashed algorithm with the different values of the bucket decider. We can conclude that using 20% of the problem size as the bucket decider tends to give the least error %.

The error is calculated using the eq. 3.

$$E\% = \frac{Calculated\ Path\ Value - Optimal\ Path\ Value}{Optimal\ Path\ value} * 100 \quad (3)$$

For the results the execution time is provided in seconds. This documented execution time does not include the time spent in parsing the input file because this is the same function in both hashed and un-hashed algorithms. It is important to note that the NN 2-opt requires an input of the number of iterations to be used. We have kept this at a constant of 4 for both the NN 2-opt and the hashed buckets algorithm. More iterations should theoretically decrease the error but after some

test runs we found that 4 iterations are a suitable average as the added optimality is not proportional to the increase in time.

### A.  Sample 1: File d493.txt
- Number of Points: 493 with optimal path cost as per TSPLIB: 35,002.
- NN 2-opt without hashing has path cost = 36,099 and execution of 1.476s thus an Error % of 3.13%
- Results for the Hashed algorithm are in Table 1.

### B.  Sample 2: File rl5915.txt
- Number of Points: 5,915 with optimal path cost as per TSPLIB: 565,530.
- NN 2-opt without hashing (average results of 2 runs) has path cost = 591,715 giving an Error % of 4.63%. Its execution was 2004.5s which is equal to 33.5 minutes!
- Results for Hashing algorithm are in Table 2.

### C.  Sample 3: File rl11849.txt
- Number of Points: 11,849 with optimal path as per TSPLIB: 923,288.

TABLE I.          D493 RESULTS (AVERAGE OF 100 RUNS)

| Decider | Hashed Time | Hashed Cost | Hashed Error |
|---|---|---|---|
| 10% = 50 (15 buckets) | 0.168s | 44,788 | 27.96% |
| 20% = 100 (9 buckets) | 0.235s | 36,364 | 3.89% |
| 40% = 200 (4 buckets) | 0.261s | 38,774 | 10.78% |

TABLE II.          RL5915 RESULTS

| Decider | Hashed Time | Hashed Cost | Hashed Error |
|---|---|---|---|
| 10% = 600 ( 10 buckets) | 39.64s | 640,430 | 13.24% |
| 20% = 1200 (6 buckets) | 66.376s | 617,670 | 9.22% |
| 40% = 2400 (4 buckets) | 155.015s | 609,920 | 7.85% |

TABLE III.          RL11849 RESULTS

| Decider | Hashed Time | Hashed Cost | Hashed Error |
|---|---|---|---|
| 5% = 600 (30 buckets) 20 runs average: | 51.92s | 1,116,000 | 20.87% |
| 10% = 1200 (15 buckets) 20 runs average: | 251.14s | 1,053,300 | 14.07% |
| 20% = 2400 (6 buckets) 3 runs average: | 917.97s | 1,002,600 | 8.59% |

- NN 2-opt without hashing: the machine ran out of memory and thus no results were gained.
- Results for Hashing algorithm are in Table 3

## VIII.   CONCLUSION

As shown in the results section, the hashed bucket algorithm is very effective in reducing the overall execution time of large scale TSPs. Fig. 4 and fig. 5 show the comparison between algorithm methods and different decider values quite clearly.

We are able to reach a path in less than 10% of the time required for the original NN 2-opt. We understand that the tradeoff is in optimality yet a 9%~15% error is considered an acceptable margin for such a gain in execution speed.

We would also like to note that original NN 2-opt algorithm was unable to run on the limited specs of the machine after a certain size due to its memory consumption. Accordingly, another advantage of the algorithm is the possibility to reach results using limited memory and execution power. This begs the possibility that the algorithm would be useful in robots and applications that run on batteries (reduced power consumption) and limited size.

REFRENCES

[1] Z. Michalewicz, Genetic Algorithms + Data Structures = Evolution Programs, Springer-Verlag, 2nd edition, 1994

[2] Keld Helsgaun, "An effective implementation of the Lin-Kernighan traveling salesman heuristic", European Journal of Operational Research, 2000

[3] Keld Helsgaun, "General k-opt submoves for the Lin–Kernighan TSP heuristic", Springer and Mathematical Programming Society, 2009

[4] Donald Davendra, Traveling Salesman Problem, Theory and Applications, InTech, 2010

[5] Hirotaka Itoh, The Method of Solving for Travelling Salesman Problem Using Genetic Algorithm with Immune Adjustment Mechanism, 2010

[6] Oloruntoyin Sefiu Taiwo, Olukehinde Olutosin Mayowa & Kolapo Bukola Ruka, "Application Of Genetic Algorithm To Solve Traveling Salesman Problem", International Journal of Advance Research (IJOAR), Volume 1, Issue 4, April 2013

[7] Ryouei Takahashi, "Solving the Traveling Salesman Problem through Iterative Extended Changing Crossover Operators", 10th International Conference on Machine Learning and Applications, 2011

[8] Atif Ali Khan, Muhammad Umair Khan, & Muneeb Iqbal, "Multilevel Graph Partitioning Scheme To Solve Traveling Salesman Problem", Ninth International Conference on Information Technology- New Generations, 2012

[9] Chris Walshaw, "A Multilevel Lin-Kernighan-Helsgaun Algorithm for the Travelling Salesman Problem", Computing and Mathematical Sciences, University of Greenwich, Old Royal Naval College, Sep. 2001

[10] Ioannis Mavroidis, Ioannis Papaefstathiou, & Dionisios Pnevmatikatos, "A Fast FPGA-Based 2-Opt Solver for Small-Scale Euclidean Traveling Salesman problem", International Symposium on Field-Programmable Custom Computing Machines, 2007

[11] Hoda A. Darwish, Hoda N. Shagar, Yasmine A. Badr, et al., "A Hashing Algorithm for Rule-Based Decomposition in Double Patterning Photolithography", IEEE 22nd International Conference on Microelectronics (ICM), 2010

Fig. 4.   Performance/Time Comparison

Fig. 5.   Error Margin Comparison

# Spectrum Sharing Security and Attacks in CRNs: a Review

Wajdi Alhakami, Ali Mansour and Ghazanfar A. Safdar

Department of Computer Science and Technology, University of Bedfordshire

Luton, LU1 3JU, United Kingdom

*Abstract*—**Cognitive Radio plays a major part in communication technology by resolving the shortage of the spectrum through usage of dynamic spectrum access and artificial intelligence characteristics. The element of spectrum sharing in cognitive radio is a fundamental approach in utilising free channels. Cooperatively communicating cognitive radio devices use the common control channel of the cognitive radio medium access control to achieve spectrum sharing. Thus, the common control channel and consequently spectrum sharing security are vital to ensuring security in the subsequent data communication among cognitive radio nodes. In addition to well known security problems in wireless networks, cognitive radio networks introduce new classes of security threats and challenges, such as licensed user emulation attacks in spectrum sensing and misbehaviours in the common control channel transactions, which degrade the overall network operation and performance. This review paper briefly presents the known threats and attacks in wireless networks before it looks into the concept of cognitive radio and its main functionality. The paper then mainly focuses on spectrum sharing security and its related challenges. Since spectrum sharing is enabled through usage of the common control channel, more attention is paid to the security of the common control channel by looking into its security threats as well as protection and detection mechanisms. Finally, the pros and cons as well as the comparisons of different CR-specific security mechanisms are presented with some open research issues and challenges.**

*Keywords*—*Dynamic Spectrum Access; Spectrum Sharing; Common Control Channel; Cognitive Radio Networks*

## I. INTRODUCTION

Cognitive Radio (CR) [1] technology promises to intelligently solve the issues in conventional wireless technology related to their limited and under-utilised spectrum [2]. This problem has become an issue of great concern given the continued increase in wireless devices that use unlicensed bands to operate, which has resulted in overcrowding, leading to inefficient use of the spectrum [3-5]. Therefore, CR provides a resolution to spectrum inefficiency and the shortage on these bands by allowing CR users (secondary users (SUs)) to opportunistically access vacant spectrum space [6]. This results in providing great opportunities for a rising number of SUs to use these bands through an optimised approach for utilising radio resources [7-8].

radio networks' (CRN) technology has its own intrinsic fundamental approach and principles for dynamic operation within the environment, unlike in the conventional wireless approach, which is based on the static radio frequency

spectrum with fixed licensed users (primary users (PUs)) and fixed channels [9]. This indicates that the cognitive ability and reconfiguration capability are the core elements that make CR an advanced technology, which grants dynamic access to the unused spectrum for both licensed and unlicensed users through certain characteristics: adoption, awareness, modification, capability of learning, observation, and communication in realistic environments [10-16]. These characteristics provide reliable communication among CR users anytime and anywhere as a smart and intelligent choice to operate dynamically through artificial intelligence algorithms, such as spectrum sensing, spectrum sharing, and spectrum mobility [13, 17]. Moreover, they differentiate this new CR technology from existing wireless technologies. Due to these sophisticated features, the CR approach is known as Dynamic Spectrum Access (DSA) or Dynamic Spectrum Management (DSM) [8,18], in recognition of the potential to realise dynamically different paradigms within a network.

However, generally DSA is considered a big challenge to implement because of its dynamic behaviour and nature, such as different frequency, geographical location, and time of operation [19-20]. Also, SUs might utilise the licensed spectrum and encounter PUs who have diverse transmission characteristics. Moreover, in comparison to known security issues that exist in wireless networks, CRNs are more exposed to threats from targeted, intelligent malicious strategies [21-22]. This poses security challenges in preventing any definite or predictable risks from occurring.

As long as spectrum sharing is one of the fundamental aspects of the CR to provide access channels and sharing resources, this overview paper mainly focuses on the spectrum sharing security of the cognitive radio MAC layer. So far, most of the literature focuses on general aspects of CRNs security in spectrum sensing and spectrum mobility-related areas. But the security of spectrum sharing has received very little research coverage. It is very important to conduct thorough research to gain a broader and clearer overview of its techniques and security-related issues.

Therefore, this overview paper firstly provides details about the spectrum sharing classification, to show the differences of the mechanism, operation, and techniques. Subsequently, it focuses and gives detailed insights into the threats and attacks that are launched in the common control channel (spectrum sharing) part of MAC layer of CRNs. In addition, it investigates and includes the recent techniques that have been developed in this area in terms of protection and detection.

This paper is organised as follows: Section 2 briefly demonstrates the CR main functions and section 3 looks into the security challenges in cognitive radio's core functions, especially in spectrum sharing, i.e. common control channel security. Section 4 discusses common security threats to both traditional wireless and cognitive radio networks. It then concludes by outlining security threats specific to CR networks. Section 5 introduces the existing security methods for achieving secure communications in both centralised and ad hoc CRNs. Section 6 identifies some open research issues and challenges before the paper is concluded in section 7.

## II. COGNITIVE RADIO CORE FUNCTIONS

There are four fundamental functions which the CRN device must perform, as shown in Figure 1 and as stated below [8, 23]

*1) Spectrum sensing identifies the parts of the accessible spectrum and senses the presence of the PU operating in the licensed band.*

*2) Spectrum management determines the best channel to establish communication.*

*3) Spectrum sharing sets up a coordination access among users on the selected channel.*

*4) Spectrum mobility vacates the channel in case the PU is detected.*

One failure can easily affect and result in deterioration of the communication or introduce vulnerabilities to the network.



Fig. 1. Cognitive radio main functions

These embedded functions have a strong relationship between them for the process of establishing an efficient communication, considering the regulations and policies that govern CRNs. Each function influences another one by providing the necessary information required during the process of reaching a final decision. For instance, once the spectrum is sensed, in order to identify the available point of access, there are two possible decisions that can be taken: If the PU is detected then the process will be discontinued; if they are not, the obtained information will move forward to the next stage. The spectrum management function then decides and selects the proper channel for the communication. Once the channel is chosen, users are directed to access it by providing their information. During a successful communication, spectrum mobility remains ready for any changes that resulted from the appearance of a PU by a regular check of the spectrum

sensing, or from other alterations to the environment in terms of the current allocation that is provided by spectrum management and spectrum-sharing elements [8, 24].

As long as CRNs have a set of nodes that interact with each other using determined policies, regulations, and sophisticated protocols [25], they have different capabilities [22, 26] relating to the spectrum awareness of the network operation and spectrum context, defined regulations and policies, quality of service (QoS), and user requirements for requesting traffic load capacity, resilience, and security. This means that cognitive nodes are able to dynamically reconfigure themselves according to the current environment in order to transmit and receive on different frequencies, in addition to supporting a variety of transmission access technology schemas [2, 27]. Another capability is resource management, which plays an important role in collaborating to assign the vacant network spectrum management resources, whether these are internal to the current network or external to conventional wireless networks [8, 28].

Spectrum sharing generally can be classified into three major criteria, based on the network architecture, access technology, and allocation behaviour (Figure 2). Descriptions of these classifications as follows:



Fig. 2. Spectrum sharing classfications

The first technique is based on the network architecture, whether it is a centralised or distributed system (Figure 3). In centralised networks such as IEEE 802.22 cognitive radio, a base station governs and senses the free channel information from neighbours' nodes within range and performs the final decision on the availability of a channel. Unlike ad hoc CR networks, CR nodes generate and utilise a common spectrum allocation for the information exchange about available channels [8, 22]. Even though the centralised entity has the advantages of addressing better efficiency, the main drawback is that a single point of failure can be easily launched to the central entity [8]. More classifications can be added to ad hoc networks, classifying them into static and mobile networks. These apply in wireless sensor networks as a static form and in MANET (Mobile Ad-hoc Network) as mobile ad hoc networks in which a set of autonomous mobile terminals are liberated to move to other existing hybrid networks [29, 30] (more details about comparing the spectrum sharing mechanisms in both centralised and distributed architectures are discussed in [31]).

The second technique is based on allocation behaviour, whether it is cooperative or uncooperative. In the cooperative method, CR users are responsible for coordinating the

functionalities of the cognitive network in order to ensure the optimisation of the spectrum utilisation and improving network efficiency through the exchange of information. However, in non-cooperative systems, CR users are not responsible for coordinating the cognitive functionalities with other cognitive devices. Instead, they implement these functions on their own [24, 32]. The main difference between these two methods is relatively clear: the first approach essentially requires the exchange of information; hence a common control channel (CCC) is required to facilitate the information exchange. However, in the second approach, the cognitive nodes do the network functions tasks on their own without the need for any collaboration from other cognitive users. This would make the task more challenging and difficult for a cognitive user. In addition, this can affect the performance due to reasons like lower efficiency, slower sharing of spectrum resources allocation, and less reliability than the cooperative technique [8, 16, 24, 33].



Fig. 3.    Cognitive radio architecture

The last classification is access technology, whether it is an overlay or underlay approach [22, 24, 26, 34]. In the overlay approach a SU utilises the spectrum without sharing with a PU. This is in contrast to the underlay approach, in which both PUs and SUs utilise the licensed spectrum at the same time [35-38], with strict power control implemented by the CR users not to interfere with the PUs.

## III.    SECURITY CALLENGES IN CRN CORE FUNCTIONS

Due to the key differences in their specifications when compared to traditional wireless networks, cognitive radio networks face certain unique challenges in terms of their continued effective use and their vulnerability to outside attack. These particular characteristics of CRNs involve the need for additional implementation of specific functions, such as proper sensing protocols, correct decision making, appropriate switching, and the provision of sufficient access for the sharing of the resources required to operate each particular function. These challenges can be classified into four main areas, which will be described in greater detail in the following subsections:

### A.    Spectrum Challemges in Spectrum Sensing

The fact that spectrum sensing is responsible for sensing channels and the provision of accurate results means that CRNs must overcome certain specific challenges. The challenges broadly pertain to the ways in which a cognitive user detects and differentiates between PUs and SUs. This is of great importance as attackers may be able to emulate the signals of the PUs, thereby increasing the likelihood of false alarms being

triggered. In addition, the hidden node problem may be another issue that can lead to a failure to detect the PUs, which would result in unacceptable shadow fading [6, 39].

### B.    Spectrum Challenges in Spectrum Management

An incorrect decision made by the spectrum management is a significant issue that could arise relatively easily. Also, the inherent complexity of the protection techniques is a key requirement to providing reliable and secure transmission of information among users. It is possible for an attacker to easily forge or tamper with the transmitted information, which would affect the correctness of any decisions made by the spectrum management.

### C.    Security Challenges in Spectrum Mobility

The requirement for a seamless handoff from one channel to another also constitutes a significant challenge for cognitive users when an attacker launches a threat to hinder or prevent this integral and flawless switching by occupying the available channels. This kind of attack could potentially increase the waiting time involved in achieving a proper handoff. This increase is certainly unacceptable to the PUs, who wants to utilise their assigned channels.

### D.    Security Challenges in Spectrum Sharing

The dynamic environment in MANET network architecture leads to more challenges and security issues arising due to the lack of the central entity which usually provides security and key management among users [40]. The control channels selection in decentralised cognitive radio networks decreases the probability of successful communication among SUs due to authenticity and validity. As discussed in [11], SUs are the non-licensed users and attackers easily exploit them and by escalating their privileges, they might damage the spectrum and the traffic of the PUs as well. Moreover, without security, this issue becomes more critical when cognitive nodes use the spectrums only when PUs are not available or not using their licensed bands. Moreover, selecting data channel(s) for exchange of data among SUs without the authenticity of the SUs is another issue that needs to be addressed in CRNs, especially for maintaining the links if a PU returns to the licensed channel.

Much research has been conducted into developing security in centralised CRNs [1-3]. However, the issue is that no research has been carried out on addressing the authentication in decentralised CRNs and its requirements, especially providing authentication of confidentiality, non-repudiation, and integrity, which are considered the main security elements in cognitive radio technology.

## IV.    SECURITY THREATS

Although cognitive radio is similar to the traditional wireless network, using a wireless medium instead of a wire to transmit information, it faces different vulnerabilities, which has resulted in the discarding of the communication process among end users [41-42]. These vulnerabilities can leads to varied threats, which can be classified into two different categories: the first relates to common security threats in both conventional wireless and CR networks, and the second category is specific to CRN users.

### A. Common Security Threats in Conventional Wireless and CR Networks

In traditional wireless technology, radio channels are used to establish communication and transmit information between communicating nodes and access points (APs) or base stations (BSs). They are used in cognitive networks to address several similar functionalities. The transmitted information can be sensitive, such as the user's identity, the user's privacy, allocation and signaling information, as well as key information. However, an attacker using a range of techniques such as eavesdropping, forgery, and masquerading attacks can easily intercept the communication during the transmission process [9, 13]. An effective security mechanism must be applied to protect data transmission from malicious behaviour like eavesdropping and information tampering [29]. Therefore, as far as data protection is concerned, different security measurements can be used for protection, detection, and countermeasures based on wireless security protocols such as WEP, WPA, and WPA2 in conventional wireless networks and EAP, AES, and 3DES in WiMAX. These security protocols are designed with encryption levels of different strengths being used according to the importance of the information being secured. Figure 4 shows the most common threats in both traditional wireless and CR networks.



Fig. 4. Common security threats in conventional wireless and CR Networks

#### 1) Fake Attacks

In the infrastructures of wireless networks, BSs or APs act as central entities that are connected wirelessly to end terminals. In order to establish communication, some information is exchanged through a radio channel between the end terminal device and the central entity. This information includes the identity data belonging to the procedures of the network control, network services and network access. A malicious user can obtain this information by wiretapping and then pose as a legitimate user. The purpose of this fake attack is a malicious user accesses the network and obtaining a network service or to launch an attack against the network [13, 43-44]. Therefore, cryptographic encryption schemes are generally used to protect the transmitted messages.

#### 2) Information Tampering

This is a serious attack that causes change, modification, replacement, or deletion of the information before it is received at its intended destination [43], and that result in misleading the receiver, who can thus make a wrong decision. Alteration significantly affects message integrity, which is unacceptable for legitimate users and network policies. However, this type of attack generally occurs in a situation where a cooperative terminal is needed to forward the information [13, 45-46].

#### 3) Service Repudiation

In this attack, when the connection is achieved between two nodes, one user denies transmitting their information for two reasons: repudiation for the communication service to deny usage of the network, which requires payment for the network usage, and repudiation for the communication content to refuse the transmission of their content. For example, when transactions are made in a commercial process, the user refuses to pay. To overcome these issues, proof-of-origin evidence can be used against a particular individual for sending or receiving messages. Identity, authentication, and cryptography encryption schemas are presently used to prevent unpredictable or hidden issues arising [13, 47].

#### 4) Replay Attack

The key purpose of this attack at the MAC layer is to obtain effective information by intercepting and retransmitting the same signed information sent to a particular node over a period of time in order to build trust with the receiver. This gives an advantage to the attacker, granting them access to new useful information like user passwords, which then enables unauthorised access to resources and control network licenses, etc [13, 48-52]. Therefore, in order to overcome this attack, the timestamp procedure is recommended because of the message validation involved [52].

#### 5) Denial of Service and Information Interference

While electromagnetic waves are essential in order to gain wireless information from users, recent advanced hardware technologies can involve a higher transmitted power in the communication process at the physical layer. It is, therefore, possible for an attacker to use this transmitter power to block the ordinary transmission and create interference and noise in the communication procedure, thereby decreasing the capacity of the wireless BS resources and equipment. This can also lessen user access through a BS terminal. Therefore, the interference of information procedures is likely to have a critical social impact [53]. An example of this occurred in 2001, which the satellite communication service was interrupted due to the high power caused by locating a VSAT terminal [13, 50].

#### 6) Greedy Behaviour Attack

During the channel negotiation process in both centralised and decentralised multi-hop networks, an attacker intends to maximise their throughput of using a spectrum through manipulating and changing the parameters of the MAC layer protocol [54-57]. This is achieved by reporting false information regarding the available channel, which causes throughput collapse for other users. For instance, in decentralised networks, if a greedy user attempts to misbehave by starving the neighbouring node, the intermediate user will be affected and banned from transmitting its messages [13].

#### 7) Malicious and Selfish Behaviour Attacks

In malicious behaviour, the attacker makes other cognitive users to make handoff from the current channel. This generally causes degrading of the network performance [29, 41, 57-58]. However, in selfish behaviour, the attacker intends to maximise their throughput by using a spectrum to disturb the normal process [59].

*8) Black and Grey Hole Attacks*

Both black and grey hole attacks exist in decentralised networks, where an attacker pretends to be the destination node. Therefore, a sender can be easily deceived and start transmitting packets. The rate of dropping the transmitting packets is used to distinguish between these two attacks. In a black hole, the malicious user obtains all the transmitted packets; however, in the grey behaviour attack, a malicious user drops part of these transmitted packets [29, 60-66].

*B. Specific Security Threats in CR Networks*

Several potentially serious threats to network performance which increase spectrum availability to malicious users have been highlighted by researchers investigating CRN technology [9, 13, 67]. Moreover, due to the unique characteristics of CRNs, they are more exposed to security threats which are usually not faced by conventional wireless technology. Therefore, security mechanisms play an important role in maintaining the network that is potentially affected by these kinds of threats [13]. Malicious attacks are well known threats that target all layers in the CRN [9, 13] with their own behaviour, which can affect network performance by attacking a particular layer. Some of the main security threats related to CRNs are identified in Figure 5.



Fig. 5. Specific security threats in CRNs

*1) Security in Spectrum Sensing*

Spectrum sensing is a major aspect of CRNs environments, providing the spectrum information about the appearance of the PU and the available channels [12, 32-33, 68]. Therefore, it is subjected to the most prevalent attacks that bring the network performance down by reporting the false results of the PU detection. As long as the security in spectrum sensing is concerned with controlling the network operation, attackers have their own malicious behaviour strategies, focusing instead on degrading the network spectrum performance by causing collisions or occupying the spectrum. This can result in potential security vulnerabilities that enable denial of service (DoS) attacks to be launched easily [67]. Thus serious attacks can occur at this level of the spectrum, which are called primary users interference (PUI) and primary user emulation (PUE).

In PUE, an attacker can simulate a signal that resembles the signal of the PU, thereby misleading the SU [2, 12, 18, 58, 69-73]. In this case, the attacker has a chance to focus on the physical layer, pretending to be an authorised user by sending CR signals that are similar to PU signals, allowing them to deceive other SUs. This increases the availability of spectrum to the malicious user. The authors of [6, 41, 74] have proposed a simulation technique used by a malicious user, which

involves a multiple stage attack that demonstrates the general influence on the network performance and other special effects on the SUs. Additionally, the authors' experiment results showed how the relationship between the performance improvements can be associated with the bands' availability and vice versa. However, in PUI, the attacker breaks the rules of the CRN mechanism by affecting network performance through interfering with PUs within the network. This forces the PU to use spectrum with noise and unavailable frequency band [13]. This is also called a jamming message attack or lion attack, where an attacker transmits high signal power to disturb the PUs through TCP connection [9, 41, 43, 75].

Several researchers have investigated and proposed algorithms to detect malicious behaviours in cooperative sensing of the spectrum in order to improve security in this stage. A detection scheme based on a past test report obtained through calculating the suspected point of secondary users, and computing the value of trust behaviour mechanism, is proposed in [74]. The proposed algorithm is able to distinguish malicious from honest users within a network. However, [76] presented a data mining technique without needing priori information about a secondary user to detect misbehaviours. In addition, [67] explained that changing the spectrum modulation system strategy and protecting the location information of the PU, and using proactive techniques in transmission, can help to prevent DoS attacks at this stage.

*2) Security in Spectrum Management*

Spectrum management is considered to be the second task after obtaining the result from spectrum sensing. Once the available bands are allocated, spectrum management determines the proper spectrum for communications based on the desired characteristics for quality of service (QoS) [22]. However, this stage cannot be safe from attacks. A forgery attack or tampering attack is designed to attack this particular level of the network element and involves the attacker transmitting incorrect spectrum sensing information to the data collection centre in order to deceive the secondary user, encouraging the wrong decision from spectrum management, which enables the malicious user to utilise the channel with superlative adaptive purpose [13, 67].

*3) Security in Spectrum Mobility*

This stage refers to the mandatory process of seamlessly switching (handoff) from a current channel to another available one due to channel occupancy by the PU. With the appearance of the PU to utilise their assigned channel, a SU must vacate and select another available channel to initiate a new connection, resulting in greater energy consumption [22, 67, 77-78].

However, from a security perspective, the availability of spectrum is reduced when there are a large number of malicious users, and this limited availability affects other legitimate SUs, who are required to vacate the current channel due to the appearance of the PU and to select another available channel [53, 78]. Moreover, a failed handoff to a proper channel may occur when an attacker forces SUs to vacate the channel by pretending to be the PU. As a consequence, it results in slower communication and requires additional time to resume the process of the communication [18, 22, 69].

### 4) Security in Spectrum Sharing

As long as spectrum sharing is crucial to maintaining effective communication in traditional wireless networks through the application of the Medium Access Control (MAC) method, it is an area of great interest for a number of researchers, who have proposed different solutions for sharing the spectrum [80, 81-83]. These solutions include a non-dedicated common control channel [84], a hopping-based control channel [85] and a dedicated CCC, also known as a Dynamic Local Common Control Channel (DLCC) [86] (Figure 6). These approaches focus on achieving a proper level of sharing among cognitive users. In this paper, a brief explanation of the first two approaches has been given, while the third approach is the main one which is considered in detail.



Fig. 6.   Specific security threats in CRNs

#### a)   Non-dedicated CCC

In this approach, a predefined non-dedicated CCC is assumed among a set of SUs. Hence, a number of CRN MAC protocols are designed for predicting that a CCC is already recognised and allocated to those SUs. Industrial, scientific and medical (ISM) or underlay ultra-wideband that is identified as unlicensed band can be the appropriate place to implement a control channel for cognitive users in order to exchange the control information [78, 80, 88].

#### b)   Hopping-based Control Channel

This approach requires a predefined channel hopping sequence that is determined among SUs in order to achieve the hopping process over the existing licensed channels [87]. Both the cognitive sender and receiver necessitate time and channel synchronisation [5]. During this process, a proper channel is determined to be utilised to transmit data through exchange of control information between the sender and the receiver. Once successful control information is exchanged between both SUs, they end the hopping process and start with the second phase of transmitting data. After the completion of the data transmission phase, the synchronisation requests are recurred with the hopping sequence [23, 80].

#### c)   Dynamic Local CCC

The CCC technique is one of the methods used to facilitate the functional sharing process between two SUs in distributed cooperative CRNs.

In distributed cooperative systems, CCC is established between both the sender and the receiver for establishing a handshaking protocol [14, 54, 80, 82, 84, 90-91]. In addition,

CCC can be used to communicate with a base station through an existing centralised entity system [92]. It is also employed to include the related information that has resulted from the spectrum sensing. Due to these effective functionalities, a number of researchers believe that CCC designed procedures can play a major role in promoting the initial exchange of information processes among cognitive nodes.

However, from a security viewpoint, no spectrum sharing classifications, which are discussed in section 2, are secure against any malicious behaviour while they are not supported with security mechanisms for protection and detection (see table 1). Generally the attackers' intention is to determinate an effective strategy that exposes a predictable risk. For instance, when CCC is used in the cooperative method of decentralised CRNs for exchanging information about the available channels and the selected channel for data transmission between SUs, it is more prone to various attacks based on selfish and malicious behaviours [41-42]. Because it is regarded as a valuable structure for the attacker to access the channel and gain the most sensitive information, a key approach for some types of attackers involves applying a PUE attack. Moreover, it is more exposed to other attack types such as eavesdropping and DoS, which can be launched easily due to existing weaknesses within the MAC layer, where poor authentication and an existing lack of encryption mechanisms enable an attacker to detect available channels that they can occupy to forge or drop MAC frames, as shown in Figure 7 [41, 56, 90].



Fig. 7.   Malicious activities in decentralised CRNs

Another vulnerability in a CCC is where an attacker forges the transmitted packets to another path and causes collisions. As a consequence, this impedes the network performance and launches a DoS attack. Once a CCC is saturated by attackers, a large number of forged packets are generated to block the exchange of the control information, enabling DoS attacks to be easily launched against the network, hence affecting its performance.

Moreover, an author in [56] suggests that encryption must be applied between legitimate SUs for the exchange of control information; otherwise, it can be readable by attackers of other cognitive users. Also, it can protect the exchanged control information over the channel from predictable control channel hopping sequences, thereby preventing itself from being saturated [13, 92].

TABLE I.      OVERVIEW OF THE ATTACKS OCCURRING AT DIFFERENT CR FUNCTIONS

| Attack Name | CR function | Description |
|---|---|---|
| Forgery & Data tamper | Spectrum Sensing | Spectrum Management system makes wrong decision by receiving the attackers' sensing information |
| Overlapping | | An attacker impacts other networks by transmission to a specific network |
| Denial of Service | | An adversary user decreases the availability of the spectrum bandwidth by blocking the communication, through creating noise spectrum signals which cause interference with PUs |
| Lion or Jamming message | | An attacker transmits high signalling power to disturb the PU or the secondary user which results forcing the cognitive user to hop to different channel to utilise |
| Spectrum Sensing Data Falsification | | In collaborative spectrum sensing, a collaboration technique used among CR nodes to generate and utilise a common spectrum allocation for the exchange of information about available channels. However adversary node gives false observations information to other users. |
| Eavesdropping | Spectrum Sharing | Weaknesses within the layer due to the poor authentication and no existing encryption mechanisms |
| Denial of Service & masquerade | | Repetition of the frequent packets that result in overcrowding the channel which is being busy to be utilised by legitimated users |
| Selfish Behaviour or selfish masquerade attack | | an attacker does not follow the normal communication process for maximising their throughput, saving energy or gaining unfair beneficial access of using spectrums through injecting frequent anomalous behaviour |
| Key depletion | | An attacker attempts to break the cipher by repetition of the session key |
| Forgery Attack | | Lack of authentication mechanism leads to the occurrence of modification and forgery on MAC CR Frames which result in the launch of DoS attacks |
| Biased Utility | Spectrum Management | An attacker tries to reduce the bandwidth of other SUs in order to obtain more bandwidth by changing the spectrum parameters |
| False feedback | | An attacker secretes the incidence of the PU in order to disturb the information sensing of other SUs |

## V. RELTAED WORKS (EXISTING SECURE COMMUNICATION SCHEME IN CRNs)

Since the layers within CRNs have their own characteristics and parameters [74, 93], they are vulnerable and allow an attacker to make a decision to launch a specific attack for the purpose of degrading the whole network performance. In MAC layer frames, an adversary has a variety of aims for misbehaving and launching such an attack. For instance, a denial of the channel service is one of the serious threats that lead to the network degradation between both sender and receiver. This attack occurs when the attacker saturates the CCC till it becomes weak for attacking [8]. In addition, selfish behaviour is another example of an attack that can also exist in the MAC layer, in which an attacker does not follow the normal process of communication. Therefore, in order to provide a defence against these threats, security mechanisms are required in the MAC layer to provide authentication,

authorisation and availability (AAA) in the CRNs. Incorporating these security features can lead to the exchange of complete and reliable secure MAC frames among cognitive users [9, 13, 94]. Thus, several studies have been conducted for secure MAC protocols in CRNs [11, 14-15, 94-101]. They are classified into two categories, based on protection and detection techniques for addressing the security requirements and to defend the existing security issues in MAC protocols in CRNs.

### A. Protection Mechanism in CRNs

In general, a number of researchers [11, 15, 94-97, 99-100] have made efforts to address the security requirements and provide secure communication among SUs by applying different security mechanisms, such as authentication and authorisation access by different techniques, within a CRN. Their proposed procedures include digital signatures, certification authority (CA), and trust-based third parties entities like server and base stations. While these solutions may be effective in some ways, they have some drawbacks.

#### 1) Digital Signature

In [11, 15, 99] proposed different protection systems based on applying a digital signatures for protecting the network from DoS attacks and providing secure communication. Their approaches involve the activities of a CA, PUs, and both PUs' and SUs' base stations. However, the main differences of these mechanisms are that the BSs are connected to the CA using wire links in [15], while in the [99] approach, an asymmetric key scheme instead of a CA is mainly used.

#### 2) Certificate Authority

Another effective traditional approach-based CA on the application layer for achieving the same purpose of authentication is presented in [100, 101]. The proposed method uses both EAP-TTLS (for establishing a secure connection) and EAP-SIM (for authenticating the user) algorithms.

#### 3) Trust Values Procedures

Other techniques based on trust values procedures are proposed in [95-96] to address and analyse the issues within CRNs. Based on this, the trust value will be calculated, which leads to the decision that will either allow the current user to utilise the available licensed channel or not.

#### 4) Other Framework Architectures

Security for authentication and authorisation architecture frameworks have been proposed in [94, 97]. Both techniques require third-party entities for appropriate access policies to the spectrum. Authors in [94] use a technique based on processing user identification in the system and providing the user preferences to third parties according to privacy rules. Based on this, the user will be authenticated and then will determine whether or not a data port would be used. However the subsequent architecture in [97] consists of two layers, which are up-layers for authentication purpose and encryption techniques, while the physical layer is for securing and protecting the spectrum.

Overall, while these proposed mechanisms are effective in some way for protecting the networks from forgery and DoS attacks, they are not applicable in a decentralised environment

because a third-party node is incorporated in order to verify the identity and provide security key managements to end users. Therefore, the security and challenges in decentralised CRNs still arise and require defensive techniques for securing communication among cognitive users. Table 2 demonstrates the pros and cons of the proposed protection mechanisms.

TABLE II.　　PROTECTION MECHANISMS IN COGNITIVE RADIO NETWORKS

| Proposed Mechanism | pros/cons | Description |
|---|---|---|
| User identification | pro | Low complexity by generating two virtual ports for secure transmission: the first is for control traffic information and another is for data transmission which is blocked by default unless the user has been authenticated. |
| | con | It requires a third party to provide information like user preferences |
| Digital signature & certificate authority | Pro | Low complexity and using the basic architectures of symmetric and asymmetric key infrastructures. |
| | Con | It has not been simulated and tested to proof the security. It also does not work in Ad-hoc environment due to being based on centralised entities. |
| Certificate authority | Pro | Effective security mechanism due to identifying and verifying the user and the server respectively. |
| | Con | Requires a third-party to verify the user identity. Also the mechanism has not been simulated and tested to ensure security against malicious behaviours. |
| Trust values | Pro | It is an additional procedure that can be built on the top of other security techniques to increase the level of the protection and detection in term of secure communication. |
| | Con | Requires a third party procedure is to provide previous information of a node. Moreover, when a new node joins the network, the CA will not be able to provide reference for that particular user. Hence the mechanism does not operate in strong fixed level of the authentication for all cognitive users equally. |

### B. Detection schemes in CRNs

Authors in [14, 98,102] have focused on the detection mechanisms in CRNs. Their proposed techniques address a variety of attacks caused by malicious and selfish behaviours, and the pros and cons of these mechanisms are illustrated in table 3.

#### 1) Selfish behaviour

Selfish behaviour detection techniques for the CCC are proposed in [14, 103], where a puzzle punishment model is applied for bad behaviour activities in a situation where a receiver is asked for a new hidden channel that has not been included previously. Thus, the sender would be a suspicious case. Therefore, the receiver applies the puzzle punishment to detect whether the sender is a selfish node or not. If the sender node solves the puzzle, they will be considered as a legitimate user and communication will be resumed normally; otherwise, the communication will be disconnected. Another technique called Cooperative neighboring cognitive radio Nodes (COOPON) is applied among a group of neighbouring users to detect selfish nodes who broadcast fake channel lists. Consequently, neighbouring users can detect the selfish users

by comparing the transmitted channel list of the target user with their lists.

#### 2) Timing parameters

Another detection mechanism was proposed in [102]. They presented a mechanism that relies on timing parameters at MAC layer. When the negotiation phase is taking place, the node, which receives a request, sets up timing parameters for controlling the time interval. This forces the sender to transmit data without getting a higher rate. If the sender does not obey and sends packets more frequently, the receiver node takes action against the sender. Then the receiver node analyses the sender's misbehaviour and broadcasts the information over the current network.

#### 3) Anomalous spectrum usage attacks (ASUAs)

The others in [98] presented a cross-layer technique for CRNs for detecting ASUAs. Collecting the information on both the physical and network layers provides an awareness of the current spectrum. It operates against the PUE and jamming attacks to provide successful access to the spectrum.

TABLE III.　　DETECTION MECHANISMS IN COGNITIVE RADIO NETWORKS

| Proposed Mechanism | pros/cons | Description |
|---|---|---|
| Selfish activity | Pro | applied in both CCC and data channel which decreases the potential of misbehaviour in different stages of the network |
| | Con | focuses only on detecting selfish behavior and does not provide the complete secure communication between sender and receiver |
| timing parameter | Pro | Detecting misbehaving nodes during the negotiation phase. It helps to maintain the channel from getting saturated. |
| | Con | -Theoretical and has not been simulated and tested to provide the detection scheme results. -Weak against eavesdropping and forgery attacks especially once the FCL is not hidden which is exploited to launch Jamming attacks. |
| Anomalous Spectrum Usage Attacks | Pro | Combining both physical and network layers for detecting malicious users give a better achievement instead of selecting only a layer |
| | Con | Focuses only on the detection approach and does not consider a significant protection scheme against both jamming and PUE attacks mobility. |

### C. Comparisons of the Presented Schemes

Incorporating the security requirements; authentication, confidentiality, non repudiation and data integrity in CRNs can lead to the exchange of complete and reliable secure MAC frames among cognitive users [9, 13, 92]. For instance, while the proposed digital signature, trust value and certificate authority procedures are different in terms of their operations (see the protection schemes in section V), the security requirements are considered for providing defense against most of the MAC threats such as DoS, Forgery, eavesdropping and spoofing in centralised CRNs. In contrast, both puzzle punishment and COOPON approaches consider only selfish behaviour among the other MAC attacks such as DoS, forgery, eavesdropping, and spoofing in decentralized CRNs. However, they are effective in selfish behaviour's detection

due to the cooperation between a group of cognitive users which involve identifying selfish users in COOPON technique and demand of solving the puzzle to resume the communication in puzzle punishment system. Moreover, the timing parameter procedure easily addresses DoS attack due to the presence of the centralised entity, which controls the cognitive users' communication. Table 4 gives information about achieving the security requirements and addressing the MAC layer attacks for each proposed scheme in both centralised and decentralised CRNs.

TABLE IV.     COMPARISION OF THE PRPOSED SECURITY SCHEMES IN CRNs

| | Puzzle punishment | COOPON | Digital Signature | Trust value | CA | Timing parameter |
|---|---|---|---|---|---|---|
| Authentication | | | √ | √ | √ | |
| Integrity | | | √ | √ | √ | |
| DoS | | | √ | √ | √ | √ |
| Forgery | | | √ | √ | √ | |
| Eavesdropping & Spoofing | | | √ | √ | √ | |
| Confidentiality | | | √ | √ | √ | |
| Non-repudiation | | | √ | √ | √ | |
| Selfish | √ | √ | | | | |
| Architecture | Ad-Hoc | Ad-Hoc | Centralised | Centralised | Centralised | Centralised |

## VI.     OPEN RESEARCH AREAS AND CHALLENGES

As long as secure communication is crucial for the exchange of information between SUs, the primary security concerns in decentralised CRNs are authentication and data confidentiality. Compromising on these elements can potentially lead to the modification, forgery or eavesdropping of the MAC frames in CR networks, which could, in turn, increase the chance of DoS attacks that would adversely affect the performance of the network. However, these security factors in ad hoc CRNs have received relatively little attention in the literature, perhaps due to their complex nature and dynamic topology [104]. These must be investigated properly in order to meet the security needs of the CRNs' technology. Further research is required in order to support the security requirements, especially to provide authentication assurance for the authorised access. These requirements assist in maintaining secure communication and enable the provision of available resources in distributed multi-hop CR environments, while simultaneously avoiding external threats. Moreover, a proper high-level encryption method is required to support secure communication between end users, although due consideration should be given to the inherent power limitations of the devices. This issue is also important because of the lack of a central entity that provides security and key management to end users. Thus, the implementation of a secure CR MAC protocol must involve the design and implementation of a robust, secure system that can achieve authentication, availability, confidentiality, integrity, non-repudiation, anonymity, and authorisation for granting security demands.

This is of fundamental importance because CR users need to incorporate security by all possible means to ensure the protection of the relatively vulnerable network operations.

## VII.     CONCLUSION

Cognitive radio networks are a remarkable area for researchers due to their use of intelligent technology for providing a solution that utilises the available spectrum efficiently. However, security is a crucial aspect of CRNs to achieve successful communication between cognitive users. Due to some unique characteristics in CRNs, different new threats to CR functions exist, such as PUE and PUI in spectrum sensing, Tampering attacks in spectrum management, failed handoffs in spectrum mobility, and MAC threats like eavesdropping, forgery, and selfish behaviour attacks in spectrum sharing are other threats. Therefore, CRN is far more exposed to security threats than those facing the conventional wireless technology. This paper presented a comprehensive survey about the challenges and security in CRNs. The information is presented as a hierarchical structure, starting with challenges and then threats in spectrum sensing, spectrum management, and spectrum mobility. A major portion of the paper has been dedicated to spectrum sharing because it has been the main motivation behind this overview. Moreover, it introduced the spectrum sharing mechanisms: Non-dedicated CCC, hopping-based control channel and more details about the common control channel were chosen for investigation and highlighted the potential existing threats and vulnerabilities. The paper also highlighted several potentially serious threats to network performance in both centralised and ad hoc CRNs. As a result, the most recent detection and protection mechanisms were discussed in terms of their pros and cons and compared for the purpose of addressing the security issues in CRNs. Finally, some open research issues and challenges were presented, which must be met to ensure secure operation of CRNs.

For future work, a hybrid secure MAC protocol for CRN is proposed in [105]. The protocol is analysed and designed for addressing the security requirements, such as authentication, confidentiality, integrity, and non-repudiation. It also addresses most of the security issues in decentralised CRN, such as spoofing, eavesdropping, and forgery attacks. Therefore, the implementation stage of the proposed protocol is in progress in order to provide results that will be compared with others belonging to different secure protocols.

## REFERENCES

[1]    Cordeiro, C., Challapali, K., Birru, D., Shankar, N., (*2005*) "IEEE 802.22: The first worldwide wireless standard based on cognitive radios," in *New Frontiers in Dynamic Spectrum Access Networks,. DySPAN. First IEEE International Symposium on,* 2005, pp. 328-337.

[2]    Shin, K., Kim, H., Min. A., Kumar, A., (2010) "Cognitive radios for dynamic spectrum access: from concept to reality," *Wireless Communications, IEEE,* vol. 17, no. 6. pp. 64-74.

[3]    Wang, H., Qin. H., Zhu. L., (2008). "A survey on MAC protocols for opportunistic spectrum access in cognitive radio networks,". in

*Computer Science and Software Engineering, 2008 International Conference on*, pp. 214-218.

[4] Cao, L., Zheng. H., (2008) "Distributed Rule-Regulated Spectrum Sharing," *Selected Areas in Communications, IEEE Journal on,* vol. 26, no. 1, pp. 130-145

[5] Zhao, Q., Tong. L., Swami. A., Chen. Y., (2007) "Decentralized cognitive MAC for opportunistic spectrum access in ad hoc networks: A POMDP framework," *Selected Areas in Communications, IEEE Journal on,* vol. 25,no.3, pp. 589-600,.

[6] Chen, R., Park, J., Hou, Y., Reed, J., (2008) "Toward secure distributed spectrum sensing in cognitive radio networks," *Communications Magazine, IEEE,* vol. 46, no. 4, pp. 50-55.

[7] Lin, F., Hu, Z., Hou, S., Yu, J., Zhang, C., Guo, N., Wicks, M., Qiu, R., and Currie, K., (2011) "Cognitive radio network as wireless sensor network (II): Security consideration,"in *Aerospace and Electronics Conference NAECON, Proceedings of 2011 IEEE National*, pp.324-328.

[8] Baldini, G., Sturman, T., Biswas, A., Leschhorn, R., Godor, G., and Street, M., (2012)"Security Aspects in Software Defined Radio and Cognitive Radio Networks: A Survey and A Way Ahead," *Communications Surveys & Tutorials, IEEE,* vol. 14, no, 2. pp. 355-379

[9] Zhang, X. and Li, C., (2009) "The security in cognitive radio networks: A survey," in *IWCMC '09: Proceedings of the 2009 International Conference on Wireless Communications and Mobile Computing: Connecting the World Wirelessly,* Leipzig, Germany, pp. 21–24.

[10] Zhen-dong, W., Hui-qiang, W., Guang-sheng, F., Bing-yang, L., Xiao-ming, C., (2010) "Cognitive networks and its layered cognitive architecture," in *Internet Computing for Science and Engineering (ICICSE), Fifth International Conference on*, pp. 145-148.

[11] Sanyal, S., Bhadauria, R. and Ghosh, C., (2009) "Secure communication in cognitive radio networks," in *Computers and Devices for Communication. CODEC. 4th International Conference on,* , pp. 1-4.

[12] Yucek, T. and Arslan, H., (2009) "A survey of spectrum sensing algorithms for cognitive radio applications," *Communications Surveys & Tutorials, IEEE,* vol. 11, no. 1, pp. 116-130,.

[13] Tang, L., and Wu, J., (2012)"Research and Analysis on Cognitive Radio Network Security," *April 2012,* vol. 4, pp. 120-126,.

[14] Wu, H., and Bai, B., (2011) "An improved security mechanism in cognitive radio networks," in *Internet Computing & Information Services (ICICIS), 2011 International Conference*, pp. 353-356.

[15] Parvin, S., and Hussain, F.. (2011), "Digital signature-based secure communication in cognitive radio networks,". in *Broadband and Wireless Computing, Communication and Applications (BWCCA), 2011 International Conference on*, pp. 230-235.

[16] Gao, Z., Zhu, H., Li, Shuai., Du, S., Li, Xu., (2012) , "Security and privacy of collaborative spectrum sensing in cognitive radio networks," *Wireless Communications, IEEE* , vol.19, no.6, pp.106-112,

[17] He, A., Bae, K., Newman, T., Gaeddert, J., Kim, Kyouwoong., Menon, R., Morales-Tirado, L., Neel, J., Zhao, Y., Reed, J., Tranter, W., (2010) "A Survey of Artificial Intelligence for Cognitive Radios," *Vehicular Technology, IEEE Transactions on,* vol. 59, no. 4, pp. 1578-1592,.

[18] Chen, R., Park, Jung-Min,. Reed, J., (2008)"Defense against Primary User Emulation Attacks in Cognitive Radio Networks," *Selected Areas in Communications, IEEE Journal on,* vol. 26, no, 1. pp. 25-37.

[19] Datla, D., Wyglinski, A., Minden, G., (2009) "A Spectrum Surveying Framework for Dynamic Spectrum Access Networks," *Vehicular Technology, IEEE Transactions,* vol. 58, no, 8. pp. 4158-4168.

[20] Li, X., Chen, J., and Ng, F., (2009) "Secure transmission power of cognitive radios for dynamic spectrum access applications," in *Information Sciences and Systems, 2008. CISS 2008. 42nd Annual Conference on,* 2008, pp. 213-218.

[21] Burbank, J., (2008) "Security in cognitive radio networks: The required evolution in approaches to wireless network security," in *Cognitive Radio Oriented Wireless Networks and Communications, 2008. CrownCom 2008. 3rd International Conference*, pp. 1-7.

[22] Zhang, Y., Xu, G., Geng, X., (2008)"Security threats in cognitive radio networks," in *High Performance Computing and Communications, HPCC '08. 10th IEEE International Conference*, pp. 1036-1041.

[23] Domenico, A., Strinati, E., Benedetto, M., (2012) "A Survey on MAC Strategies for Cognitive Radio Networks," *Communications Surveys & Tutorials, IEEE,* vol. 14, no. 1, pp. 21-44.

[24] A. Umamaheswari., V. Subashini. and P. Subhapriya.,(2012)"Survey on performance, reliability and future proposal of cognitive radio under wireless computing," in *Computing Communication & Networking Technologies (ICCCNT), 2012 Third International Conference*, pp. 1-6.

[25] Kamruzzaman, S.,Alam, M., (2010). "Dynamic TDMA Slot Reservation Protocol for QoS Provisioning in Cognitive Radio Ad Hoc Networks". *World Academy of Science, Engineering and Technology* , 449-791

[26] Ji, Zhu., and Liu, K., (2007) "cognitive radios for dynamic spectrum access - Dynamic Spectrum Sharing: A Game Theoretical Overview," *Communications Magazine, IEEE,* vol. 45, no, 5. pp. 88-94.

[27] Akyildiz, I., Lee, W., and Chowdhury, K., (2009)"CRAHNs: Cognitive radio ad hoc networks," *Ad Hoc Networks,* vol. 7, no. 5, pp. 810-836, 7.

[28] [28] Wei, W., (2011)"The research of cognitive communication networks," in *Communication Software and Networks (ICCSN), 2011 IEEE 3rd International Conference*, pp. 1-5.

[29] Soleimani, M., and Ghasemi, A., (2011) "Detecting black hole attack in wireless ad hoc networks based on learning automata," in *Computer Sciences and Convergence Information Technology (ICCIT), 2011 6th International Conference* , pp. 514-519.

[30] Aboudagga, N., Refaei, M., Eltoweissy, M., Dasilva, L., Quisquater, J., (2005) "Authentication protocols for ad hoc networks taxonomy and research issues," in *Q2SWinet '05 Proceedings of the 1st ACM International Workshop on Quality of Service & Security in Wireless and mobile Networks,* ACM New York, NY, USA, pp. 96 - 104.

[31] Salami, G., Durowoju, O., Attar, A., Holland, O., Tafazolli, R., and Aghvami, H., (2011) "A Comparison Between the Centralized and Distributed Approaches for Spectrum Management," *Communications Surveys & Tutorials, IEEE,* vol. 13, no, 2. pp. 274-290.

[32] Ejaz, W., Hasan, N., Kim, H., and Azam, M., (2011) "Fully distributed cooperative spectrum sensing for cognitive radio ad hoc networks," in *Frontiers of Information Technology (FIT), 2011*, pp. 9-13.

[33] Wang, W.,(2009) "Spectrum sensing for cognitive radio," in *Intelligent Information Technology Application Workshops, 2009. IITAW '09. Third International Symposium*, pp. 410-412.

[34] Arkoulis, S., Kazatzopoulos, L., Delakouridis, C. Marias, G., (2008) "Cognitive spectrum and its security issues," in *Next Generation Mobile Applications, Services and Technologies, 2008. NGMAST '08. the 2$^{nd}$ International Conference*, pp. 565-570.

[35] Zhao, Q., Swami, A., (2007) "A survey of dynamic spectrum access: Signal processing and networking perspectives," in *Acoustics, Speech and Signal Processing,. ICASSP. IEEE International Conference*, pp. 1349-1352.

[36] Manosha, K., Rajatheva, N., Latva-aho, M., (2011) "Overlay/Underlay Spectrum Sharing for Multi-Operator Environment in Cognitive Radio Networks," *Vehicular Technology Conference (VTC Spring), 2011 IEEE 73rd* , pp.1-5, 15-18 May 2011,

[37] Senthuran, S.; Anpalagan, A.; Das, O. (2012), "Throughput Analysis of Opportunistic Access Strategies in Hybrid Underlay-Overlay Cognitive Radio Networks,"*Wireless Communications,IEEE Transactions on* , vol.11, no.6, pp.2024-2035, June 2012,

[38] Wyglinski, M., Nekovee, M., Hou, T., (2010). Cognitive Radio Communications and Networks: Principles and Practice.(2010 Elsevier)

[39] Akyildiz, I, Lee, W., Vuran, M., and Mohanty, S., (2008)"A survey on spectrum management in cognitive radio networks," *Communications Magazine, IEEE,* vol. 46, no 4, pp. 40-48.

[40] Yu, F., and Tang, H., (2010)"Distributed node selection for threshold key management with intrusion detection in mobile ad hoc networks ," *Springer Science+Business Media, LLC,* vol. 16, no. 8, pp. 2169–2178,

[41] León, O., Hernández-Serrano, J., and Soriano, M.,(2010) "Securing cognitive radio networks," vol. 23, pp. 633-652,

[42] Fragkiadakis, A., Tragos, E., Askoxylakis, I., (2012) "A Survey on Security Threats and Detection Techniques in Cognitive Radio Networks," *Communications Surveys & Tutorials, IEEE,* vol. 15, no. 1, PP, pp. 1-18.

[43] Sampath, A., Dai, H., Zheng, H., Zhao, B., (2007) "Multi-channel jamming attacks using cognitive radios," *Computer Communications & Networks. Proceedings of 16th International Conference*, pp. 352-357.

[44] Song, Y., Zhou, K., & Chen, X. (2012). "Fake BTS Attacks of GSM System on Software Radio Platform". *journal of networks, vol. 7, no. 2 , 7*, 275-281.

[45] Terence, J., (2011), "Secure route discovery against wormhole attacks in sensor networks using mobile agents," *Trendz in Information Sciences and Computing (TISC), 3rd International Conference on* , pp.110-115.

[46] Robles, R., Haas, J., Chiang, J., Hu, Y., Kumar, P., (2010), "Secure topology discovery through network-wide clock synchronization," *Signal Processing and Communications (SPCOM), International Conference*, pp.1-5, 18-21 July 2010,

[47] Rai, A., Tewari, R., & Upadhyay, S. (2010). "Different Types of Attacks on Integrated MANET-Internet Communication". *International Journal of Computer Science and Security* , vol. 4, no. 3, pp. 265-274

[48] Li, L., Kidston, D., Vigneron, P., Mason, P. (2011), "Replay attacks and detection in tactical MANETs," *Communications, Computers and Signal Processing (PacRim), IEEE Pacific Rim Conference on* , pp.226-231,

[49] Goyal, P. Batra, S., Singh, A., (2010). "A Literature Review of Security Attack in Mobile Ad-hoc Networks". *International Journal of Computer Applications* , vol. 9, no. 12, pp. 0975 – 8887

[50] Enneya, N., Baayer. A., Elkouttbi, M., (2011). "A Dynamic Timestamp Discrepancy against Replay Attacks in MANET". *Informatics Engineering and Information Science* ,vol 254, pp. 479-489

[51] Goyal, P., Parmar, V., & Rishi, R. (2011). "MANET: Vulnerabilities, Challenges, Attacks, Application". *IJCEM International Journal of Computational Engineering & Management, Vol. 11* , pp . 2230-7893.

[52] Baayer, A., Enneya, N., & Elkoutbi, M. (2012). "Enhanced Timestamp Discrepancy to Limit Impact of Replay Attacks in MANETs". *Journal of Information Security*, vol 3 , pp. 224-230.

[53] Jakimoski, G., and Subbalakshmi, K., (2008) "Denial-of-service attacks on dynamic spectrum access networks," in *Communications Workshops, 2008. ICC Workshops '08. IEEE International Conference*, pp. 524-528.

[54] Attar, A., Tang, H., Vasilakos, A., Yu, F. and Leung V., (2012) "A Survey of Security Challenges in Cognitive Radio Networks: Solutions and Future Research Directions," *Proceedings of the IEEE,*vol. 100, no. 12, pp. 3172-3186

[55] [55] Djahel, S., Abdesselam, F., Turgut, D., (2009) "An effective strategy for greedy behavior inwireless ad hoc networks," in *Global Telecommunications Conference,. GLOBECOM 2009.IEEE*, pp. 1-6.

[56] Zhu, L., and Zhou, H., (2008) "Two types of attacks against Cognitive radio network MAC protocols," in *Computer Science and Software Engineering, 2008 International Conference*, pp. 1110-1113.

[57] Guang, L., Assi, C.,(2006) "Mitigating smart selfish MAC layer misbehavior in ad hoc networks," in *Wireless and Mobile Computing, Networking and Communications, (WiMob'2006). IEEE International Conference*, pp. 116-123.

[58] Chaczko, Z., Wickramasooriya, R., Klempous, R., Nikodem, J., (2010) "Security threats in cognitive radio applications," in *Intelligent Engineering Systems ,14th International Conference*, pp. 209-214.

[59] Akkarajitsakul, K., Hossain, E., Niyato, D., Kim, D., (2011) "Game Theoretic Approaches for Multiple Access in Wireless Networks: A Survey," *Communications Surveys & Tutorials, IEEE,* vol. 13, no. 3, pp. 372-395.

[60] Kariya, D., Kathole, A., and Heda, S., (2012) "Detecting Black and Gray Hole Attacks in Mobile Ad Hoc Network Using an Adaptive Method," vol. 2, no. 1, pp. 2250-2459.

[61] Yi, P., Zhu, T., Liu, N., Wu, Y. Li, J.,( 2012) "Cross-layer Detection for Black Hole Attack in Wireless Network," vol. 8, no. 10, pp. 4101- 4109.

[62] Jhaveri R., Patel, S., and Jinwala, D.,(2012) "A novel approach for Gray Hole and Black Hole attacks in mobile ad hoc networks," in *Advanced Computing & Communication Technologies (ACCT), 2<sup>nd</sup> International Conference*, pp. 556-560.

[63] [63] Jhaveri R., Patel, S., Jinwala, D., (2012), "DoS attacks in mobile ad hoc networks: A survey," in *Advanced Computing & Communication Technologies (ACCT),Second International Conference*, pp. 535-541.

[64] Joshi, A., Agrawal, K., Arora, D. and Shukla, S., (2011) "Efficient Content Authentication in Ad-Hoc Networks-Mitigating DDoS Attacks," vol. 23, no. 4, pp. 0975 – 8887.

[65] Cai, J., Yi, P., Chen, J., Wang, Z., Liu, N., (2010)"An adaptive approach to detecting black and gray hole attacks in ad hoc network," in *Advanced Information Networking & Applications (AINA), 24th IEEE International Conference*, pp. 775-780.

[66] Xiaopeng, G., Wei, C., (2007) "A novel gray hole attack detection scheme for mobile ad-hoc networks," in *Network & Parallel Computing Workshops NPC, IFIP International Conference*, pp. 209-214.

[67] Mao, H., and Zhu, L., (2011)"An investigation on security of cognitive radio networks," in *Management and Service Science (MASS), 2011 International Conference*, pp. 1-4.

[68] Zhe, C., Guo. N., Qiu, C., (2010)"Demonstration of real-time spectrum sensing for cognitive radio," IEEE Communications Letters, *2010 - milcom*, vol, 14,. no. 10, pp. 323-328.

[69] Jin, Z., Anand, S., and Subbalakshmi, K., (2012) "Impact of Primary User Emulation Attacks on Dynamic Spectrum Access Networks," *Communications, IEEE Transactions on,* vol. 60, no, 9, pp. 635-2643

[70] Yuan, Z., Niyato, D., Li, H., Song, Z., and Han, Zhu., (2012)"Defeating Primary User Emulation Attacks Using Belief Propagation in Cognitive Radio Networks," *Selected Areas in Communications, IEEE Journal,* vol. 30, no, 10, pp. 1850-1860

[71] Zhou, X., Xiao, Y., Li, Y., (2011)"Encryption and displacement based scheme of defense against primary user emulation attack," in *Wireless, Mobile & Multimedia Networks (ICWMMN 2011), 4th IET International Conference*, pp. 44-49.

[72] Huang, L., Xie, L., Yu, H., Wang, W., and Yao, Y., (2010)"Anti-PUE attack based on joint position verification in cognitive radio networks," in *Communications and Mobile Computing (CMC), 2010 International Conference*, pp. 169-173.

[73] Anand, S., Jin, Z., and Subbalakshmi, K.,(2008) "An analytical model for primary user emulation attacks in cognitive radio networks," in *New Frontiers in Dynamic Spectrum Access Networks, 2008. DySPAN 2008. 3rd IEEE Symposium*, pp. 1-6.

[74] Wang, W., Li, H., Sun, Y., and Han, Z., (2009)"Attack-proof collaborative spectrum sensing in cognitive radio networks," in *Information Sciences and Systems, 2009. CISS 2009. 43rd Annual Conference*, pp. 130-134.

[75] Wu. Y, Wang. B., Ray, L., and Clancy, T., (2012)"Anti-Jamming Games in Multi-Channel Cognitive Radio Networks," *Selected Areas in Communications, IEEE Journal on,* vol. 30, no, 1, pp. 4-15,

[76] Li, H., Han, Zhu.,(2010)"Catching attacker(s) for collaborative spectrum sensing in cognitive radio systems: An abnormality detection approach," *New Frontiers in Dynamic Spectrum, IEEE Symposium*, pp. 1-12.

[77] Feng, W., Cao, J., Zhang, C., Liu, C., (2009) "Joint optimization of spectrum handoff scheduling and routing in multi-hop multi-radio cognitive networks," in *Distributed Computing Systems. ICDCS. 29th IEEE International Conference*, pp. 85-92.

[78] Song, Y., Xie, J., (2012) "ProSpect: A Proactive Spectrum Handoff Framework for Cognitive Radio Ad Hoc Networks without Common Control Channel," *Mobile Computing, IEEE Transactions on* , vol.11, no.7, pp.1127-1139, July 2012

[79] Akyildiz, I., Lee, W., Vuran, M., Mohanty, S., (2008)"A survey on spectrum management in cognitive radio networks," *Communications Magazine, IEEE,* vol. 46, no 4, pp. 40-48.

[80] A, H., Salameh, B., and Krunz, M., (2009) "Channel access protocols for multihop opportunistic networks: challenges and recent developments," *Network, IEEE,* vol. 23, no 4, pp. 14-19,.

[81] Brik, V., Rozner, E., Banerjee, S., Bahl, P., (2005) "DSAP: A protocol for coordinated spectrum access," *New Frontiers in Dynamic Spectrum Access Networks.1<sup>st</sup> IEEE International Symposium*, pp. 611-614.

[82] Ma, L., Han, X., and Shen, C., (2005)"Dynamic open spectrum sharing MAC protocol for wireless ad hoc networks," in *New Frontiers in Dynamic Spectrum Access Networks. DySPAN. First IEEE International Symposium*, pp. 203-213.

[83] Sankaranarayanan, S., Papadimitratos, P., Mishra, A. Hershey, S., (2005)"A bandwidth sharing approach to improve licensed spectrum

utilization," in *New Frontiers in Dynamic Spectrum Access Networks,. DySPAN 2005. First IEEE International Symposium*, pp. 279-288.

[84] Kondareddy, Y., Agrawal, P., Sivalingam, K., (2008), "Cognitive Radio Network setup without a Common Control Channel," *Military Communications Conference,.MILCOM. IEEE* ,pp.1-6, 16-19 Nov. 2008

[85] Lin, Z., Liu, H., Chu, X., Leung, Y., (2011), "Jump-stay based channel-hopping algorithm with guaranteed rendezvous for cognitive radio networks," *INFOCOM, 2011 Proceedings IEEE*., pp.2444-2452.

[86] Romero, E., Mouradian, A., Blesa, J., Moya, J., and Araujo, A., (2012)"Simulation framework for security threats in cognitive radio networks," *Communications, IET,* vol. 6, no. 8, pp. 984-990.

[87] Song, Y., Xie, J., (2012) "ProSpect: A Proactive Spectrum Handoff Framework for Cognitive Radio Ad Hoc Networks without Common Control Channel," *Mobile Computing, IEEE Transactions on* , vol.11, no.7, pp.1127-1139, July 2012

[88] Salameh, H.B.; Krunz, M.; Younis, O.; (2008), "Distance- and Traffic-Aware Channel Assignment in Cognitive Radio Networks," *Sensor, Mesh and Ad Hoc Communications and Networks, 2008. SECON '08. 5th Annual IEEE Communications Society Conference on* , vol., no., pp.10-18, 16-20 June 2008

[89] Shih, C., Wu, T., Liao, W., (2010) , "DH-MAC: A Dynamic Channel Hopping MAC Protocol for Cognitive Radio Networks," *Communications (ICC), IEEE International Conference*, pp.1-5.

[90] Safdar, G., and O'Neil, M., (2012) "A novel common control channel security framework for cognitive radio networks," *Int. J. of Autonomous and Adaptive Communications Systems,* vol. 5 No: 2, pp. 125 - 145.

[91] Kahraman, B., and Buzluca, F., (2010)"Protection and fairness oriented cognitive radio MAC protocol for ad hoc networks (PROFCR)," in *Wireless Conference (EW), 2010 European,* 2010, pp. 282-287.

[92] Bian, K. and Park, J.,(2006)"MAC-layer misbehaviors in multi-hop cognitive radio networks," *In Proceedings of the 2006 US-Korea Conference on Science, Technology and Entrepreneurship (UKC2006). National Science Foundation Under Grant CNS-0524052.*

[93] Ci, S., and Sonnenberg, J., (2007)"A cognitive cross-layer architecture for next-generation tactical networks," in *Military Communications Conference, 2007. MILCOM 2007. IEEE*, pp. 1-6.

[94] Prasad, N., (2008) "Secure cognitive networks," in *Wireless Technology, 2008. EuWiT 2008. European Conference*, pp. 107-110.

[95] Parvin, S., Han, S., Tian, B., Hussain, F., (2010), "Trust-based authentication for secure communication in cognitive radio networks,"in *Embedded and Ubiquitous Computing (EUC), IEEE/IFIP 8th International Conference*, pp. 589-596.

[96] Parvin, S., Hussain, F., (2012) "Trust-based Security for Community-based Cognitive Radio Networks",. 26th IEEE International Conference on Advanced Information Networking and Applications, pp. 518-525

[97] .Li Zhu; Huaqing Mao, "Unified Layered Security Architecture for Cognitive Radio Networks," *Power and Energy Engineering Conference (APPEEC), 2011 Asia-Pacific* , vol., no., pp.1,4, 25-28 March 2011

[98] Sorrells, C; Potier, P; Qian, L; Li, X., (2011) "Anomalous spectrum usage attack detection in cognitive radio wireless networks," in *Technologies for Homeland Security (HST), 2011 IEEE International Conference,* 2011, pp. 384-389.

[99] Mathur, C., Subbalakshmi, K., (2007), "Digital signatures for centralised DSA networks"*Consumer Communications & Networking Conference, CCNC. 4th IEEE*

[100] Zhu, L., Mao, H., (2010) "Research on authentication mechanism of cognitive radio networks based on certification authority," in *Computational Intelligence and Software Engineering (CiSE), 2010 International Conference on,* 2010, pp. 1-5.

[101] Zhu, L., Mao, H., (2011), "An Efficient Authentication Mechanism for Cognitive Radio Networks," *Power and Energy Engineering Conference (APPEEC), 2011 Asia-Pacific*, pp.1-5, 25-28 March 2011,

[102] Shaukat, R., Khan, S., Ahmed, A., (2008) "Augmented security in IEEE 802.22 MAC layer protocol,"in *Wireless Communications, Networking & Mobile Computing,. '08. 4th International Conference*, pp 1-4.

[103] Jo, M., Han, L., Kim, D., In, H.P., (2013) "Selfish attacks and detection in cognitive radio Ad-Hoc networks," *Network, IEEE* , vol.27, no.3, pp.46,50,

[104] Goyal, P., Parmar, V., & Rishi, R. (2011). "MANET: Vulnerabilities, Challenges, Attacks, Application". *IJCEM International Journal of Computational Engineering & Management, Vol. 11* , pp . 2230-7893.

[105] Alhakami, W.; Mansour, A.; Safdar, G.A.; Albermany, S., "A secure MAC protocol for Cognitive Radio Networks (SMCRN)," *Science and Information Conference (SAI), 2013* , pp.796,803, 7-9 Oct. 2013

# A Comparative Usability Study on the Use of Auditory Icons to Support Virtual Lecturers in E-Learning Interfaces

Marwan Alseid
Department of Software Engineering
Applied Science University
Amman, Jordan

Mohammad Azzeh
Department of Software Engineering
Applied Science University
Amman, Jordan

Yousef El Sheikh
Department of Computer Science
Applied Science University
Amman, Jordan

*Abstract*—**Prior conducted research revealed that the auditory icons could contribute in supporting the virtual lecturers in presence of full body animation while delivering the learning content in e-learning interfaces. This paper presents further empirical investigation into the use of these supportive auditory icons by comparing three different e-learning interfaces in terms of usability aspects; effectiveness, user satisfaction and memorability. The aim is to find out which combination of the tested multimodal metaphors is the best one in terms of utilizing the auditory icons to supplement the presentation of learning material by virtual lecturer. The first experimental e-learning interface incorporates a speaking virtual lecturer with full body gestures along with supportive auditory icons. The second experimental e-learning interface includes the use of virtual lecturer speech in the absence of his body and accompanied with the same auditory icons used in the first interface. However, the third interface is similar to the second one in terms of using the virtual lecturer's speech but without any additional auditory icons. The obtained results have shown that the inclusion of auditory icons could enhance the usability and learning performance of e-learning interfaces much better if combined along with speaking virtual lectures in the absence of any body animation.**

*Keywords—auditory icons; virtual lecturer; e-learning; usability; speech; avatar; multimodal interaction*

## I. INTRODUCTION AND MOTIVATION

The world has witnessed and still, a tremendous and accelerating development in the field of information technology and computer networks, which resulted in a quick and easy access to a huge amount of information including educational content. E-Learning describes the learning process that utilizes information and communication technology in delivering and managing the learning material. Recently, the majority of e-learning interfaces concentrate on user's visual channel to communicate information whereas other multimodal interaction metaphors could be involved to make use of other human sense in the interaction process. Previous related work demonstrated that speech and non-speech sounds as well as virtual lecturers could be beneficial in enhancing the usability of e-learning interfaces. Even though, the incorporation of multimodal interaction metaphors in this domain needs to be investigated further. The aim of the presented experimental study is to reveal the best utilization of auditory icon as supportive sounds to the virtual lecturer in e-learning interfaces. Therefore, three different e-learning interfaces have been developed and independently tested in terms of effectiveness, memorability and user satisfaction. Each of these interfaces involved different combination of speaking virtual lecturer and auditory icons to communicate the learning material about class diagram notation. The following sections present an overview of the relevant work in e-learning and multimodal interaction, the experimental e-learning interfaces, the experimental design, analysis and discussion of the obtained results.

## II. RELATED WORK

E-learning is the term that describes the learning process via information and communication technology [1, 2] in which a huge educational content could be easily and quickly accessed. As a result, utilizing this technology in delivering e-learning content has been and still investigated by researchers. Scheduled delivery is an example of the technology used in e-learning [3] where video broadcasting, remote libraries, and virtual classrooms have been used and constrained with time and place. This technology has been enhanced by the on-demand delivery platforms that facilitate anytime and anywhere learning in the forms of interactive training CD ROMs and web-based training. Compared to traditional learning, e-learning has many advantages such as offering more flexible learning in terms of time and location, enabling better adaptation to individual needs [4], facilitating online collaborative learning over the Internet [5] as well as increasing learners' motivation and interest about the presented material [6]. Nevertheless, it was found that students felt uncomfortable with computer-based learning and missed traditional face-to-face interaction with teacher. Therefore, users' accessibility and their attitude in regard to e-learning should be enhanced [7].

Multimodal interaction involves more than one human sense in human computer interaction and could be utilized to enhance the usability of user interfaces. It facilitates the use of different channels to communicate different information [8]. Also, it enables users to employ the most suitable communication metaphor to their abilities [9]. This multimodality in the interaction process was found to be helpful in enhancing the learning experience where visual, aural, haptic and other channels could be integrated in a

multimodal approach to deliver the learning content in e-learning interfaces.

Avatar is a multimodal metaphor that represents a human-like or cartoon-like computer-generated character [10] and utilizes both auditory and visual human senses in human-computer interaction. It has been used in interactive interfaces to communicate verbal and non-verbal information through facial expressions and body gestures [11]. Also, it was found that users' satisfaction and their ability to understand and remember the provided knowledge has been enhanced by the incorporation of speaking avatar [12]. In addition, it has been demonstrated by several studies that the use of avatars contributed positively in terms of facilitating the learning process and enhancing users' satisfaction in e-learning [13-15] and in e-book assessment interfaces [16]. Even though, these studies did not investigate the use of avatars along with auditory icons in e-learning interfaces.

Speech and non-speech sounds could be used to complement the visual output; however, sound is more flexible as it can be heard without paying visual attention to the output device. It was found that speech sounds could contribute with graphics and non-speech sounds (earcons and auditory icons) to enhance the usability of search engines interfaces [17] and e-government interfaces [18]. A study by Alseid and Rigas [19] investigated three different e-learning interfaces and demonstrated that the interface incorporating speaking virtual lecturer with full body gestures was found to be the most efficient, effective and satisfactory in communicating the learning content as opposed to the other two interfaces incorporating either single or two talking heads of facially expressive virtual lecturers. Even though, that study did not explore the incorporation of supportive non-speech sounds such as auditory icons. Therefore, a further experimental study has [20] been carried out by the same authors to investigate the non-speech sounds when used along with the speech of full body animated virtual lecturer during the presentation of learning content and found that earcons and auditory icons could be used beneficially in communicating auditory messages related to important parts of the learning content while being delivered by speaking virtual lecturer. However, that study involved only one group of participants who tested only one interface and did not investigate other combinations of non-speech sounds with virtual lecturers in e-learning interfaces. Therefore, this study aimed to investigate the use of auditory icons further. More specifically, it compares three experimental e-learning interfaces in terms of usability. Two of these interfaces incorporate auditory icons to support speaking virtual lecturer but one of them include the presence of full body gesture of the virtual lecturer whereas the other do not. The third interface employs the speech of virtual lecturer in the absence of his body gestures and any supportive auditory icons.

### III. EXPERIMENTAL E-LEARNING INTERFACES

Three different e-learning interfaces have been designed and built to serve as a basis for this study. These interfaces provide a command button to start the presentation of the learning material and pause/play functionalities to facilitate more control on the learning. The first interface; VSNS (Virtual lecturer with Speech and Non-Speech) includes an avatar as a



Fig. 1. Screenshot of VSNS (Virtual lecturer with Speech and Non-Speech) interface.

virtual lecturer with full body gestures and speaking naturally (recorded speech) with prosody to present the learning material in audio-visual format. Also, it includes the use of auditory icons to support the learning material presented by the virtual lecturer. The aim of including the avatar as mentioned earlier was to imitate the traditional interaction between the lecturer and the learner that usually take place in traditional classroom-based learning. In addition, VSNS interface offers textual and graphical representation of the communicated learning material placed in the background of the virtual lecturer who has been designed to simulate body gestures similar to those usually used by the human lecturer in the real classroom situation. Fig. 1 shows a screenshot of VSNS interface. In the second e-learning interface; RSNS (Recorded Speech with Non-Speech), the same textual and graphical information related to the presented learning material is placed in the middle of the interface (see Fig. 2) and explained similarly by the speech of the virtual lecturer but in the absence of his body. In addition, the same supportive auditory icons used in VSNS have been used similarly in RSNS. However, the third interface; RSON (Recorded Speech ONly) is identical to RSNS but without using any auditory icons.

Auditory icons are non-speech sounds come from the surrounding environment such as opening a window, closing a window, dropping a can and shredding a paper sounds. Such sounds have been empirically investigated to communicate



Fig. 2. Screenshot of RSNS (Recorded Speech with Non-Speech) interface.

different types of information and found to be beneficial in delivering such information and successfully enhanced the usability of user interfaces [21]. As explained earlier, two of the experimental e-learning interfaces (VSNS and RSNS) include the use of auditory icons. This inclusion aimed to bring users' attention to two key aspects of the learning material while being communicated by the virtual lecturer. These are the beginning and end of an important statement in the presented content where the sound of door opening has been used to indicate that the virtual lecturer is about to start mentioning an important statement, and the sound of door closing to indicate that this statement completed. Accordingly, representing these aspects by these auditory icons could provide natural mapping to facilitate remembering and interpreting its meaning successfully by users. The two sounds are played five times for five different statements in the presented content. Also, these sounds are played in pause intervals of virtual lecturer speech in order to avoid interference between both sounds. On overall, three types of multimodal interaction metaphors are incorporated in the experimental e-learning interfaces as shown in table 1.

TABLE I.    INCORPORATED MULTIMODAL INTERACTION METAPHORS

| Interfaces | Text and Graphics (Visual) | Speaking Virtual Lecturer with Full Body Gestures (Audio-Visual) | Speaking Virtual Lecturer without Body (Audio) | Auditory Icons (Audio) |
|---|---|---|---|---|
| VSNS | X | X | | X |
| RSNS | X | | X | X |
| RSON | X | | X | |

With respect to the learning material presented by the three tested e-learning interfaces, it is the same. It is a single lesson that contains introductory information about class diagrams including class diagram notations, what is meant by class and object, and how to differentiate between them. This content has been adapted from [22].

## IV.    EXPERIMENTAL DESIGN

The experimental work described in this paper aimed to examine three different usability aspects of the experimental e-learning interfaces. These aspects are the effectiveness, memorability, and user satisfaction.  The effectiveness is evaluated in terms of users' learning performance by the correctness of their answers to the learning activities prior and post experimentation with the tested e-learning interfaces. The memorability is evaluated by users' ability to recognize the sound (i.e. auditory icon) used and its meaning. However, the user satisfaction is assessed by users' responses to the satisfaction questionnaire. More specifically, the experiment aims to answer the following questions:

*1) Which of the experimental e-learning interfaces is more usable than the others in terms of effectiveness, memorability and user satisfaction?*

*2) Which of the experimental e-learning interfaces is the most effective one in terms of users' learning performance?*

*3) Would the participants be able to remember the meaning of the incorporated auditory icons successfully?*

*4) Which of the experimental e-learning interfaces is the most satisfactory to the participants?*

*5) Which is better to use in e-learning interfaces; a speaking and full body animated virtual lecturer with supportive auditory icons or without it?*

*6) Which is better to use in e-learning interfaces; the speech of the virtual lecturer along with supportive auditory icons or without it?*

### B.  Hypotheses

The main hypotheses in relation to the presented study are:

*1) There will be a difference in the usability of the tested e-learning interfaces in terms of effectiveness, memorability and user satisfaction.*

*2) There will be a difference in the users' learning performance before and after interacting with the experimental e-learning interfaces.*

*3) There will be a difference in users' evaluation of the involved auditory icons before and after experimentation.*

*4) The participants will be able to correctly remember the investigated auditory icons.*

### C.  Participants

In total, 45 participants involved voluntarily in this experiment and randomly assigned in equal proportions (N=15) to one of three independent groups each of which tested one of the experimental interfaces; VSNS, RSNS and RSON. All the participants used the experimental interfaces for the first time and were undergraduate students enrolled in information technology programs at the Applied Science University, Jordan. Fig. 3 shows that most of the participants (93%) in the three groups were 18-23 years old and first year undergraduate students (87%-93%). With respect to their gender, the number of females and males participants was approximately equal. Furthermore, their area of study was distributed as follows: Software Engineering (SE) with 27% of each group, Computer Information Systems (CIS) and Computer Networks Systems (CNS) with 27%, 20% and 27% for VSNS, RSNS and RSON respectively, and Computer Science (CS) with 20% of VSNS, 33% of RSNS, and 20% of RSON. The selection of participants was mainly based on their previous knowledge in the learning topic; class diagram notation.



Fig. 3.    Users' profiling.

In this regard, the majority of them (93%-100%) had no experience indicating that they relied on the communicated learning information to answer the learning performance questions after experimenting with the tested interfaces.

### D. Procedure

The same procedure has been followed with each user throughout the course of the experiment and the participation was on an individual bases. The experiment started with user-profiling and pre-experimentation tasks. Also, users of VSNS and RSNS have been provided with a short training in 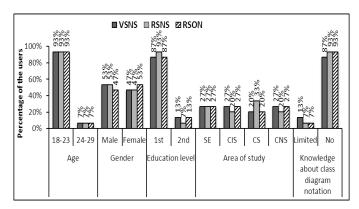which they had the opportunity to listen to the implemented auditory icons to insure their ability to understand and interpret each of these sounds when incorporated later in the experimental e-learning interfaces VSNS and RSNS. Next, the learning material about class diagram notation has been presented using one of the experimental e-learning interfaces. Once this presentation finished, the user has been instructed to carry out post-experimentation tasks as well as to provide any comments or suggestions.

### E. Tasks

Prior to experimentation with the assigned interface, each user has been requested to provide personal data in relation to age, gender, educational level and major. In addition, they have been requested to declare their previous experience in Computers, Internet, class diagram notation and e-learning applications. Furthermore, users involved in testing both VSNS and RSNS interfaces have been asked to express their points of views in terms of annoying, focus and helpfulness regarding the use of auditory icons in electronic learning in the absence of any interactive context. Then, users have been asked to answer a set of four questions about the learning material adopted in this experiment (pre-test). Two of these questions are recall activities in which the user needs to retrieve some information to be able to answer. The other two questions are recognition activities that provide the user with multiple answers to the required question and he/she is required to recognize the correct one.

Once the experimentation finished, the user has been instructed to carry out learning performance (post-test), memorability, and satisfaction tasks as well as to provide any comments or suggestions. The learning performance tasks asked the user to answer the same question of the pre-test in relation to the learning content delivered by the tested interface in order to measure the effectiveness of that interface as well as to measure how much learning gained by users. The memorability task aimed to evaluate users' ability to identify the sounds used to communicate key aspects of the presented content. More specifically, three different sounds were played and the user had to recognize which one has been used to communicate the start or end of an important statement in the lesson. This task has been applied only with users of both VSNS and RSNS interfaces. The satisfaction task aimed to obtain users' attitude towards the tested interface. In this task, the user has been instructed to fill a satisfaction questionnaire of 14 different statements on a 5-point Likert scale. The first 10 statements were the SUS questionnaire [23] whereas the remaining 4 statements were related to learning experience. For

VSNS and RSNS users, the experiment finished with additional task to obtain their opinions with respect to the implemented auditory icons in terms of annoying, focus and helpfulness.

### F. Variables

Three different types of variables have been considered which are: independent variables, dependent variables and controlled variables.

The independent variables represent the factors manipulated in the experiment and assumed to be the cause of the results. In this study, the presentation mode has been considered as the independent variable where the experimental e-learning interfaces offered three different modes for the presentation of the learning material. The first mode used text with graphics, a speaking virtual lecturer with full body gestures and auditory icons (VSNS interface). The second mode incorporated text with graphics along with the speech of virtual lecturer, and auditory icons (RSNS interface). However, the third mode included text with graphics and the speech of virtual lecturer (RSON interface).

The dependent variables are measured as a result of manipulating the independent variables. The dependent variables regarded in this study are as follows:

- Effectiveness: correctness of user's responses to the required learning performance activities; recall and recognition. In recall questions, partial or total correct answers have been considered whilst in the recognition questions, the answer had to be totally correct. The difference in users' learning performance between pre-test and post-test has been considered as well.

- Memorability: users' recognition of auditory icons has been measured by the number and percentage of users successfully recognized the non-speech sounds after being used in the experimental e-learning interfaces (VSNS and RSNS).

- User satisfaction: measured by users' responses to satisfaction questionnaire.

The controlled variables represent the external variables associated with the procedure of the experiment and could affect the obtained results. The controlled variables (known also as confounding variables) should be kept consistent throughout the experiment to avoid their influence on the dependent variables and so insure that the independent variables are the only cause of the experimental results. The controlled variables in this experiment are:

- Required tasks: the same tasks have been required from the participants.

- Presented learning material: the same learning content about class diagram notation has been presented by all experimental e-learning interfaces.

- Awareness of questions: none of the users were aware of the required learning performance questions in both pre and posttests.

- Procedure consistency: the same procedure has been followed during the execution of the experiment including measurement tools and used equipment.

- Familiarity with the interface: all participants were first-time users of the tested interfaces and provided with the same level of training prior to experimentation.

## V. RESULTS

The obtained experimental results have been analyzed in terms of different parameters including users' views regarding the auditory icons accompanied the virtual lecturer voice in both the absence and presence of interactive e-learning context. In addition, these parameters involved the effectiveness (learning performance), memorability and users' satisfaction. A more details are provided in the following subsections.

### A. Users' Evaluation of auditory icons

Before starting the interaction with the tested interface, users of both VSNS and RSNS have been requested to provide their opinions towards auditory icons when used to accompany the voice of virtual lecturer in e-learning interfaces. They have been asked to answer Yes or No if they think that it is annoying, could aid to focus, and helpful during the learning process. The same question has been repeated by the end of the experimentation. Fig. 4 shows that users' attitude positively changed towards auditory icons after being used in an interactive e-learning context within the tested VSNS and RSNS interfaces. The percentage of users who felt annoyed dropped from 67% to 27% and from 80% to 40% for VSNS and RSNS respectively which means that about half of them found it not annoying. This figure also demonstrates that the tested auditory icons did not substantially split users' attention away from the presented content where 73% of VSNS users and 67% of RSNS users believe that these sounds aided them to focus during the interaction compared to 20% and 7% respectively who think that it could enhance their concentration when it has been introduced in the absence of any interactive e-learning context. With respect to the helpfulness of the tested auditory icons in the learning process, it can be seen from Fig. 4 that users' opinion considerably changed after they tested the two interfaces. Prior to the experiment, 33% of VSNS users and smaller percentage (13%) of RSNS users thought that incorporating these non-speech sounds could help in enhancing their learning. This percentage increased remarkably to 80% (VSNS) and 73% (RSNS) after they have had the opportunity

to experience it interactively. In summary, the addition of auditory icons to the experimental e-learning interfaces was found to be neither annoying nor distracting and helpful to improve learning. These findings support the results of previous research [20].

### B. Learning Performance

One of the main concerns of the present study was the learning performance of the participants (effectiveness of the tested interfaces) which has been measured and compared in terms of the number of correctly answered questions related to the experimental content. Each user was required to answer the same 4 questions before (pretest) and after (posttest) interacting with the experimental interfaces, and therefore, the total number of questions in each case was 60. Fig. 5 shows the percentage of correctly answered questions by each group of users in both pretest and posttest. The obtained results demonstrated that the participants had a weak background about class diagram notation prior the experiment. It can be seen that the overall percentage of correct answers was 15% in VSNS and RSNS each compared to 10% in RSON. The participants achieved an average total score (i.e., the sum of correct answers out of 4) of 0.6 (SD = .63) in condition VSNS, 0.6 (SD = .51) in condition RSNS and 0.4 (SD = .63) in condition RSON. As expected, Kruskal-Wallis test revealed that no significant differences have been found between the experimental conditions in the pretest (H(2) = 1.64, p >.05), which means that all participants were at the same low level of knowledge with respect to the content communicated later on by the experimental interfaces. Mann-Whitney tests were used to follow up this finding and revealed that no significant differences were found between RSNS and VSNS (U = 109.5, r = -.26), between RSNS and RSON (U = 87, r = -.22), and between VSNS and RSON (U = 91.5, r = -.18). These results were found to be consistent with users' profiles where most of them had no or limited previous knowledge about the experimental learning material (refer to Fig. 3).

Confirming what has been initially hypothesized, the posttest provided different results. The overall percentages of correct answers were 75%, 60% and 50% in RSNS, VSNS and RSON respectively. In other words, condition RSNS resulted in a higher average total score (3, SD = 0.65) than condition VSNS (2.4, SD = 0.91) which in turn achieved a higher average total score than the RSON condition (2, SD = 1.13). The differences in posttest results were found to be significant according to Kruskal-Wallis test (H(2) = 8.42, p <.05) indicating that users' learning has been significantly affected by
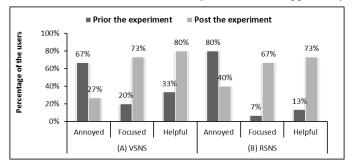
Fig. 4. Views of VSNS users (A) and RSNS users (B) about the tested auditory icons when used in both the absence and presence of interactive e-learning context.
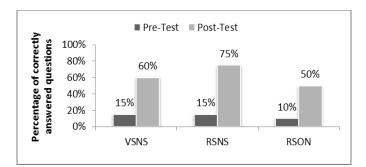
Fig. 5. Percentage of correct answers achieved by the users.

the way each of the tested interfaces used to deliver the learning content. Mann-Whitney follow up tests revealed that RSNS users significantly outperformed both VSNS users (U = 69, r = -.36) and RSON users (U = 51, r = -.49). However, no significant difference was found between VSNS users and RSON users (U = 82.5, r = -.24). Therefore, it can be concluded that if the learning material is presented using RSNS interface it will significantly increase users' learning performance (i.e., correctly answered questions) compared to presenting the same content using VSNS or RSON interfaces; however, users' learning performance will not be affected by VSNS or RSON interfaces.

For further analysis, the difference between posttest results and pretest results were computed to obtain a feedback on how much each of the tested interfaces affected on the learning of the users. It can be seen that the highest difference 60% was achieved by RSNS users followed by VSNS users (45%) and RSON users (40%).

Summarizing these findings, communicating the learning content using the speech of virtual lecturer in line with auditory icons (condition RSNS) could result in larger learning advantage in comparison to conditions VSNS and RSON.

### C. Remembrance of Auditory Icons

In order to make use of auditory icons in communicating auditory notifications related to key features of the learning content while being presented by the virtual lecturer, users should be able to successfully recognize it and interpret its meaning correctly. Therefore, by the end of interacting with the tested interfaces, users of both RSNS and VSNS were requested to carry out the memorability task to measure their remembrance of the tested auditory icons. Three sounds (auditory icons) have been played for each of the two key features and the user had to recognize which sound has been used to communicate which feature. This task has not been requested from RSON users because this interface did not incorporate any use of auditory icons. Fig. 6 shows the correctness rate of users' responses to this task. Although users of RSNS performed better, it can be seen that users of both interfaces achieved a high percentage of correct recognition. On overall, 97% of the tested auditory icons were correctly recognized by users of RSNS compared to 83% by users of VSNS. Also, users of RSNS were more able to correctly recognize the start of statement sound "opening a door" (100%) and the end of statement sound "closing a door" (93%) than users of VSNS (93% and 7% for "opening a door" and "closing a door" respectively). In brief, these results demonstrate that the tested auditory icons could be successfully interpreted and remembered by the users when used to indicate the importance of specific content delivered by a virtual lecturer in e-learning interfaces. This was found to be consistent with the findings of previous experiment [20] were similar results have been attained but with different experimental design in which the same auditory icons were used and tested in a single interface in addition to other auditory icons.

### D. User Satisfaction

At the end of experiment, users were required to respond to the satisfaction questionnaire composed of 14 statements each



Fig. 6. Users' successful recognition of the tested auditory icons.

of which had a 5-point Likert scale ranging from 1 representing strong disagreement to 5 representing strong agreement. The first 10 statements were adopted from SUS questionnaire [23] to obtain users' attitude towards the different aspects of the tested interfaces. For the analysis of results, the SUS scoring method has been used for the SUS statements, whereas the mode and median were calculated for the other 4 statements.

Findings showed that the RSNS interface scored the highest SUS satisfaction score (M = 84.17, SD = 21.79) compared to VSNS (M = 79.5, SD = 25.88) and RSON (M = 60.5, SD = 13.63). An analysis of variance (ANOVA) yielded significant difference in users' attitude towards the three interfaces $(F_{(2,)} = 5.32, p<.05)$. Also, the results of follow up pairwise comparisons found significant satisfaction difference between both RSNS and RSON (p<.05) and between VSNS and RSON (p<.05), however, not between RSNS and VSNS (p>.05). In other words, the interfaces that incorporated auditory icons were more satisfactory to the users than the interface which didn't use this kind of non-speech sounds.

In addition to the SUS statements, another 4 statements were included to obtain feedback from users regarding their learning experience attained during the interaction with the tested interfaces. More specifically, these statements investigated users' excitement and interest about the presented content (*S11- I was excited and interested about what has been presented in the lesson*), ease of identifying important parts of this content (*S12- It was easy to identify the important parts of the presented lesson*), and user's willing to use e-learning if presented similar to the tested interface (*S13- I would like to use e-learning once more if presented this way*). The last statement aimed to evaluate overall users' satisfaction (*S14- On overall, I am satisfied with the interface*). Users' responses to the additional four statements are shown in Fig. 7. The same level of users' ratings for S11, S12, and S13 statements could be observed in RSNS and VSNS with mode and median valued four where users of both interfaces expressed their agreement about these statements. However, users of RSON were neutral about the same statements. In other words, users of both RSNS and VSNS were more excited and interested, more capable to capture important parts of the content, and more willing to reuse these interfaces for e-learning compared to users of RSON. On overall (S14), users of RSNS were more satisfied than users of VSNS and users of RSON who were generally satisfied in spite of their neutral impressions regarding S11, S12, and S13. To summarize, both RSNS and VSNS provided the users with more enriching learning experience in comparison to RSON.

Fig. 7.   Users' rating of the additional four statements S11, S12, S13, and S14 in the satisfaction questionnaire.

## VI.   Discussion

The experimental study reported in this paper investigated three different e-learning interfaces each of which incorporated different multimodal approach to communicate the learning content to the participants. The first interface (VSNS) involved the use of virtual lecturer's speech with the presence of animated body gesture and accompanied by the sounds of auditory icons. The second interface (RSNS) was similar to VSNS but with the absence of virtual lecture's body. However, the third interface (RSON) included only the speech of the virtual lecturer. Otherwise, the three interfaces were similar to each other in regards to visual appearance. Our goal was to identify among these interfaces which one provides better use of auditory icons in accompanying the virtual lecturer speech. The obtained results have been used to compare these ways of presentation in terms of effectiveness (learning performance), memorability and user satisfaction where the difference among the three experimental interfaces with respect to these usability attributes has been predicted in the research hypotheses.

The experimental results revealed that the tested interfaces were significantly different in terms of users' learning performance as well as users' satisfaction.  In addition, the results of multiple comparisons among the three interfaces demonstrated that the RSNS was the most effective in communicating the learning content to the participants and this has been reflected on their ability to achieve the highest learning performance in terms of correctly answered questions. Also, RSNS was found to be the most satisfactory presentation to the participants. Presenting the learning content using the RSNS interface enabled the users to fully concentrate on the delivered content shown on the screen and explained by the speech of virtual lecturer and as a result enhanced their understanding of the presented information . At the same time, using auditory icon sounds contributed to capture users' attention to the important statements spoken by the lecturer and helped them to identify the most important parts of the learning content particularly if we know that most of the participants who tested RSNS stated that these sounds did not annoy them, helped them to focus and were helpful in their learning (see Fig. 4). This can be attributed to the fact that using "opening a door" and "closing a door" sounds helped the users to establish natural mapping between the communicated information and familiar sounds from everyday life and each of these sounds

transmit only one meaning and used consistently throughout the tested interfaces. As well, users of RSNS were more able to remember these sounds compared to the users of VSNS (see Fig. 6). This has been supported by users responses to the satisfaction questionnaire where RSNS users found themselves excited and interested about what has been presented, capable to easily identify the important pats of content, and like to learn from similar e-learning interfaces (see Fig.  7). As a result, they were significantly more satisfied compared to users of VSNS and RSON.

On the other hand, the VSNS interface came in the second place in terms of the investigated usability attributes. Previous research [18] has proven that e-learning interface which include virtual lecturer speaking the learning material with the presence of full body gesture along with supportive auditory non-speech messages (similar to VSNS) could significantly enhance the usability of e-learning interfaces. That research, however, tested only that interface by one group of users without comparing it with another interface that make use of virtual lecturer and non-speech sounds such as RSNS. The findings of the current study demonstrated that VSNS-like interfaces performed lower in usability evaluation when compared to RSNS. Although the VSNS interface enabled the users to be engaged in learning environment similar to the real face-to-face interaction take place in the traditional class rooms, it seems that the presence of full body animation contributed to split users' attention away from the presented learning content and overloaded their visual channel moving their eyes between two visual metaphors; virtual lecturer and his background content, and as a result they achieved lower learning performance (see Fig. 5) giving that they expressed positive impressions about the used auditory icons (see Fig. 4) and were able to remember it successfully (see Fig. 6) as well as were satisfied with the interface and attained learning experience (see Fig. 7) like RSNS users.

Similar to RSNS, the RSON interface enabled the users to hear the spoken explanations of the presented content and watching that content at the same time which contributed to reduce users' visual overload and keeping them involved better in cognitive processing of the communicated content compared to VSNS. Even though, users of RSON attained the lowest number of correctly answered questions. This can be attributed to the contribution made by the auditory icons used in RSNS

and VSNS which improved users' concentration and attention towards the presented content. Also, users of RSON were found to be significantly less satisfied comparable to VSNS and RSNS users. On overall, the obtained results confirmed the experimental hypotheses and suggest that combining auditory icons with virtual lecturer speech in the absence of body gestures is much better in enhancing the usability and learning performance of e-learning interfaces than combining auditory icons with virtual lecturer speech in the presence of body gesture.

## VII. CONCLUSION

This paper described further empirical investigation into the use of auditory icon to communicate supportive auditory messages related to the learning content while being delivered by the virtual lecturer speech. The main aim of this investigation was to identify the best combination of these two metaphors when incorporated in e-learning interfaces in addition to other visual metaphors. In order to achieve this aim, three different e-learning interfaces have been built and experimentally tested by three independent groups of users each of which examined one of the experimental e-learning interfaces in terms of usability attributes; effectiveness (learning performance), memorability of auditory icons, and user satisfaction. The obtained results revealed that incorporating auditory icons with the speech of full-body animated virtual lecturer could enhance the usability of e-learning interfaces much better compared to the remaining tested combinations of multimodal interaction metaphors.

### REFERENCES

[1] S. Alexander, "E-learning developments and experiences," Education and Training, vol. 43, pp. 240-248, 2001.

[2] D. Yu, W. Zhang, and X. Chen, "New Generation of E-Learning Technologies," Proceedings of the First International Multi-Symposiums on Computer and Computational Sciences (IMSCCS'06), 2006.

[3] R. Hamilton, C. Richards, and C. Sharp, "An examination of e-learning and e-books," 2001.

[4] F. Mikic and L. Anido, "Towards a Standard for Mobile E-Learning," *International Conference on Systems and International Conference on Mobile Communications and Learning Technologies, 2006. ICN/ICONS/MCL 2006. International Conference on Networking*, pp. 217-222, 2006.

[5] A. P. Correia and P. Dias, "Criteria for evaluating learning web sites: how does this impact the design of e-learning," Actas da II Conferência *Internacional de Tecnologias da Informação e Comunicação na Educação: Desafios Challenges*, pp. 521-528, 2001.

[6] G. Theonas, D. Hobbs, and D. Rigas, "Employing Virtual Lecturers' Facial Expressions in Virtual Educational Environments," International Journal of Virtual Reality, vol. 7, pp. 31-44, 2008.

[7] W. L. Johnson, S. Kole, E. Shaw, and H. Pain, "Socially intelligent learner-agent interaction tactics," *Proceedings of the International Conference on Artificial Intelligence in Education*, 2003.

[8] N. B. Sarter, "Multimodal information presentation: Design guidance and research challenges," International Journal of Industrial Ergonomics, vol. 36, pp. 439-445, 2006

[9] A. Dix, G. Abowd, J. Finlay, and R. Beale, *Human-Computer Interaction (3rd Edition)*. Prentice Hall, 2004

[10] R. Sheth, "Avatar Technology: Giving a Face to the e-Learning Interface," *The eLearning Developers' Journal*, 2003.

[11] J. Beskow, "Animation of Talking Agents," *Proceedings of AVSP*, vol. 97, pp. 149-152, 1997.

[12] M. Alotaibi and D. Rigas, "The Role of Avatars with Facial Expressions to Communicate Customer Knowledge," *International Journal of Computers, NAUN*, vol. 3, pp. 1-10, 2009.

[13] J. Holmes, "Designing agents to support learning by explaining," Computers & Education, vol. 48, pp. 523-547, 2007.

[14] R. Moreno and R. Mayer, "Interactive Multimodal Learning Environments," Educational Psychology Review, vol. 19, pp. 309-326, 2007.

[15] L. A. Annetta and S. Holmes, "Creating Presence and Community in a Synchronous Virtual Learning Environment Using Avatars," International journal of instructional technology and distance learning, vol. 3, pp. 27-43, 2006.

[16] D. Rigas and A. Algahtani, "An Investigation of Multimodal Metaphors in E-Book Assessment Interfaces", HCI International 2013-Posters' Extended Abstracts, 567-571, 2013.

[17] A. Ciuffreda and D. Rigas, "A usability Study of multimodal interfaces for the presentation of Internet Search Results," *International Journal of Computers, NAUN*, vol. 2, pp. 120-125, 2008.

[18] D. Rigas and B. Almutairi, "Investigating the impact of combining speech and earcons to communicate information in e-government interfaces", Human-Computer Interaction. Interaction Modalities and Techniques, 23-31, 2013.

[19] M. Alseid and D. Rigas, "Three Different Modes of Avatars as Virtual Lecturers in E-learning Interfaces: A Comparative Usability Study", Open Virtual Reality Journal, Bentham Open, ISSN: 1875-323X, volume 2: p.p. 8-17, 2010.

[20] M. Alseid and D. Rigas, "The Role of Earcons and Auditory Icons in the Usability of Avatar-Based E-Learning Interfaces", Proceedings of the 4th International Conference on Developments in E-Systems Engineering, p.p. 276-281, 2011

[21] W. W. Gaver, "The SonicFinder: An Interface That Uses Auditory Icons," *Human-Computer Interaction*, vol. 4, pp. 67-94, 1989.

[22] T. C. Lethbridge and R. Laganiere, *Object-oriented software engineering*: McGraw-Hill Education, UK, 2001.

[23] J. Brooke, "SUS: a" quick and dirty" usability scale," *Usability evaluation in industry*, pp. 189-194, 1996

# Probabilistic Monte-Carlo Method for Modelling and Prediction of Electronics Component Life

T. Sreenuch, A. Alghassi and S. Perinpanayagam

Integrated Vehicle Health Management Centre
Cranfield University
Bedford MK43 0AL, UK

Y. Xie

Shanghai Aircraft Design and Research Institute
Commercial Aircraft Cooperation of China
Shanghai 201210, P. R. China

*Abstract*—Power electronics are widely used in electric vehicles, railway locomotive and new generation aircrafts. Reliability of these components directly affect the reliability and performance of these vehicular platforms. In recent years, several research work about reliability, failure mode and aging analysis have been extensively carried out. There is a need for an efficient algorithm able to predict the life of power electronics component. In this paper, a probabilistic Monte-Carlo framework is developed and applied to predict remaining useful life of a component. Probability distributions are used to model the component's degradation process. The modelling parameters are learned using Maximum Likelihood Estimation. The prognostic is carried out by the mean of simulation in this paper. Monte-Carlo simulation is used to propagate multiple possible degradation paths based on the current health state of the component. The remaining useful life and confident bounds are calculated by estimating mean, median and percentile descriptive statistics of the simulated degradation paths. Results from different probabilistic models are compared and their prognostic performances are evaluated.

*Keywords—Prognostics; Monte-Carlo Simulation; Remaining Useful Life*

## I. INTRODUCTION

Nowadays in order to provide warning and predict failures to avoid catastrophic failure of products and systems, there has been an increasing tendency in monitoring the ongoing "health" of them [1]. The insulated gate bipolar transistor (IGBT) modules play an increasing significant role in on-wing for avionics system, such as communication system in autonomous working, radar system and navigation system and so on. The failures of IGBT components can degrade the efficiency of the systems or result in system failures [1]. In general, IGBT modules have several thousand hours' lifetime expectancy [2], but in order to analyze failures from several of them, the lifetime of the modules needs to be reduced. Therefore the process that causes it to fail must still operate the module within its specifications, but in a greatly reduced time frame.

IGBT accelerated aging system is to design and implement a system capable of performing robust experiments on gate controlled power transistors to induce and analyze prognostic indicators [3]. The main goal for the development of experiment system was to identify precursor parameters for device failure. Precursor parameters are parameters of the device that change with time wherein the change can be mapped to degradation in the device. Once the precursor parameters are identified, suitable diagnostic and prognostic algorithms can be implemented using these parameters to provide early warning of failure and predict remaining useful life [4].

Hence, a comprehensive approach to the development of a prognostic framework for IGBTs is required, there is a necessity to develop methods to predict the remaining useful life (RUL) of IGBTs to prevent system stoppage and costly failures. Prognostic is a technology under ongoing development. The technology aims towards high technology sectors, for example the automotive or aerospace industries, for ensuring safety and customer satisfaction. Most modern vehicles monitor their systems to ensure correct operation. If a fault is detected or predicted the user of the vehicle is usually notified before the fault has had a detrimental effect on the vehicle. Modern vehicles also monitor their usage and change their service intervals accordingly. The reliability of IGBTs directly affect the reliability and performance of these vehicle system. In recent years, series of research work about IGBT reliability, failure mode and aging analysis has been carried out widely, and a suitable prognostic method for IGBT and an efficient algorithm for predicting the IGBT RUL become increasingly important.

As electronic components have an increasingly consumption in new generation aircrafts and vehicles, and the amount of electronic failure will also become significant. Fault diagnosis and prognostic, estimation of remaining useful life and health management have vital roles to avoid catastrophic failure, improve aircraft reliability, reduce maintenance cost and increase performance [5]. This paper bases its study on IGBT for development of algorithms for estimating remaining useful life of components, and it is considered to contribute to the prognostic technology development in integrated vehicle health management (IVHM) field and advance the electronic components prognosis.

There are several approaches that have been developed for electronic prognostics. The issues unaddressed in previous IGBT prognostics studies will form the basis for the motivation of the current study. [6] used a system model approach to estimate the remaining useful life of lithium ion batteries. The battery was represented by a lumped parameter model. The parameters of the model were calculated using relevance vector machine (RVM) regression on experimental data. An extended Kalman filter and particle filter algorithms were used to determine the battery RUL.

[7] describes the use of prognostic cells to predict failure in integrated circuits. The prognostic cell was developed to fail

prior to the circuit on the same chip for all realistic operating conditions. Prognostic monitors in the test cell experienced the exact environment that the actual circuit experienced, but at an accelerated rate, thereby providing failure prediction. [8] used the data-driven approach to detect anomalies of notebook computers by monitoring performance parameters and comparing them against the historical data using Mahalanobis distance.

A physics-based prognostic approach was used by [7] in the development of a diagnostic system based on a virtual system. Using a Virtual Test Bed (VTB), system faults found in a real world system were simulated along with a normally operating real world system. For the development of a fault diagnostic system for a brake-by-wire system, [9] used a similar fuzzy system approach. Six failure modes of the system were identified and three measurement points chosen. The input signals were processed on a segment-by-segment basis, by a feature extraction process, and by fault detection.

The aim of this paper is to develop a prognostic approach that is applicable to power electronic components and computational efficient embeddable in a low power device. Here, IGBT is used a case study. The outline of this paper is as follows: Section I describes the background of the prognostic is described and summarizes the current research work on IGBT. Section II describes IGBT accelerated aging experiments, IGBT aging data. The aging data are processed and the degradation profiles of IGBT are analyzed. In section III, the maximum likelihood method is utilized to computing the parameters of IGBT degradation models. In section IV, Monte Carlo simulation method and IGBT degradation models are used to predict the RUL, and the algorithm of IGBT prognostic is developed. The RUL prediction results are analyzed in section V and the error and root mean square error are analyzed to compare the efficiency of different models in predicting the RUL. Section VI concludes the paper and the future works are discussed.

## II. IGBT DEGRADATION PROFILE

### A. Aging Experiments

The IGBT accelerated aging experiments are designed to study the aging characters of the IGBT and develop the algorithm of prognostic for prediction of the remaining useful life. The IGBT degradation data set is acquired from the aging process system, which is provided by the AMES laboratory of NASA [10]. The data set can be used to design and develop prognostic algorithms for semiconductor components such as IGBTs which have increasingly been used in modern multiple vehicle systems. IGBT accelerated aging experiments belong to the project in NASA to investigate the degradation characterizations of electronic components [11], as electronic components have an increasingly consumption in new generation aircrafts and vehicles, and the amount of electronic failure will also become significant. Fault diagnosis and prognostic, estimation of remaining useful life and health management have a vital role to avoid catastrophic failure, improve aircraft reliability, reduce maintenance cost and increase performance.

IGBT accelerated aging experiments are based on the aging platform which induces the degradation and electronic faults into the test system. Prevalently, four kinds of accelerate aging methods are widely used in accelerated aging experiments, which are thermal cycling, hot carrier injection, electrical over stress and time dependent dielectric breakdown stimulus [12]. The IGBT functional failure such as die solder degradation and wire lift were brought by the thermal cycling accelerate aging approach. Hot carrier injection could accelerate electrons and holes pass into gate oxide, which could result in the increase of IGBT threshold voltage. IGBT condition mutation and lighting could be caused by the electrical overstress due to the excessive voltage, current or power. The breakdown of IGBT gate oxide will happen when the charge injection exceeds the threshold which is caused by accumulating of the temperature in the gate oxide when it is being operated [12]. Accelerated aging approaches such as thermal cycling and electrical overstress are used in IGBT accelerate aging experiments to speed up the degradation and failure of the IGBT in experiments environments which simulate the scenarios of industrial practical application. Precursor parameters, such as collector voltages, collector currents, gate voltages and currents, and environmental parameters such as temperature are monitored and recorded to be utilized for IGBT diagnosis and prognosis research [13].

The experiment data and measurements are shared in the website of NASA as an open database which can be used to develop prognostic algorithms available to academic and industrial researchers [10]. IGBT accelerated aging data set are measurements and sensor data collected from IGBT accelerated aging experiments platform shown in figure 1. The data set includes the measurements being recorded from IGBT experiments (or operating) environment and survey data representing the deterioration of IGBT in the experiments. This data set contains mass data from thermal overstress aging experiments, including several parameters being recorded continuously such as collector current, collector voltage, gate voltage, package temperature etc. [14]. These data and parameters were monitored and recorded constantly until the IGBT failure in accelerate aging experiments. The data set were formatted in a data array which could be read by MATLAB to facilitate analysis and processing for the data in the subsequent research and investigation.

Figure 2 depicts a process of the prognostic algorithm development used in this paper. Firstly, IGBTs are tested in accelerated aging experiments following standard experiment procedures in an environmental simulation scenario to accelerate aging and failure. The monitoring data and experiments parameters are recorded and collected to transport into the software platform which are used as data formation and data storage. IGBT diagnostic and prognostic investigation will based on these data set. Prognostic algorithm for RUL prediction will be developed, and Monte Carlo simulation is used in the prognostic algorithm which will be described in more detail in the subsequent sections.

## B. Degradation Data

The aim of data processing is to gain useful information from the data with the approach of analysis and sorting. Collector emitter voltage is selected as a precursor parameter for the IGBT aging prognostic in this paper [15]. The profile of the $V_{CE}$ collected from the aging experiment is presented in figure 3. The collector emitter voltage of the IGBT presents a monotone increasing in the whole aging process and the $V_{CE}$ is also presents a fluctuation and oscillation during this process, but the $V_{CE}$ falls quickly at the end of the aging process when the IGBT becomes to fail. The whole aging process is more than 10000 time units. The whole aging process is more than 10000 time units.

The aging data of raw $V_{CE}$ as a precursor parameter are processed by low-pass filtering, and its filtered profile is shown in figure 4. It can be seen that $V_{CE}$ presents an increase step by step during the whole IGBT aging process. The data is now clean and more suitable for the analysis. The variation of $V_{CE}$ in the whole aging process could be separated into 7 stages, and the values of $V_{CE}$ for each phase are discretely different. Seven IGBTs were used in the accelerated aging experiment. It can be seen that the degradation stages are clearly separated from each other. The time duration where $V_{CE}$ stays in each stage are computed and listed in table I. The $V_{CE}$ voltage value is approximately 2.45V at the starting of the aging process. The degradation states can be determined by



Fig. 1.    IGBT accelerated aging experiments hardware [10].



Fig. 4.    Collector-Emitter Voltage after K-Mean Clustering.



Fig. 2.    Process of IGBT prognostic algorithm development.



Fig. 3.    Collector-Emitter Voltage Profile.

the level of $V_{CE}$. In this particular IGBT, $V_{CE}$ increases about 0.5V discrete step at each degradation phase.

TABLE I.        IGBT DEGRADATION DATA SETS.

| IGBT No. | IGBT degradation process | | | | | | |
|---|---|---|---|---|---|---|---|
| | 1st Phase | 2nd Phase | 3rd Phase | 4th Phase | 5th Phase | 6th Phase | Failure |
| 1 | 670 | 970 | 1368 | 2369 | 3793 | 5079 | 10740 |
| 2 | 827 | 827 | 1389 | 1389 | 2306 | 3208 | 5075 |
| 3 | 1099 | 1099 | 1905 | 2490 | 2900 | 3889 | 6799 |
| 4 | 894 | 894 | 1733 | 2384 | 2789 | 3887 | 5141 |
| 5 | 927 | 1055 | 2115 | 2544 | 3388 | 7449 | 10285 |
| 6 | 578 | 578 | 1560 | 2109 | 3403 | 4236 | 12164 |
| 7 | 750 | 1631 | 2755 | 3001 | 3757 | 5079 | 6861 |

## III.    MAXIMUM LIKELIHOOD ESTIMATION

### A. Degradation Model

Table II records the run-to-failure degradation process of 7 IGBT samples used in the accelerated aging experiment. The columns are the time durations of each degradation phases. It can be seen that an IGBT will degrade and undergo 6 degradation phases before it eventually fail. Each phase will last for a period of time before the degradation progresses further to the next phase. Take the first IGBT for example, the operational use life of IGBT-No.1 is 5079 unit time, and the duration of the IGBT stayed in its first degeneration phase is

670 unit time, then the IGBT degraded into the second degradation phase and stayed 300 unit time before its further degradation to step into the third degradation phase. And so on, until the IGBT had stayed for 1286 unit time in the last phase, the IGBT continued degraded and completely failed.

In this paper, the occurrence of degradation (or time duration of each degradation phase) is assumed to be random and uncorrelated to other degradation phases. Therefore, 6 independent stochastic process models could be built to represent the degradation phases which follow the random probability distribution. In this paper, Gamma, Exponential and Poisson distributions are used in modeling the degradation process.

TABLE II.        IGBT DEGRADATION PHASE DURATION.

| IGBT No. | Duration of Each Phase | | | | | | |
|---|---|---|---|---|---|---|---|
| | *1st Phase* | *2nd Phase* | *3rd Phase* | *4th Phase* | *5th Phase* | *6th Phase* | *IGBT Life* |
| 1 | 670 | 300 | 398 | 1001 | 1424 | 1286 | 5079 |
| 2 | 827 | 0 | 562 | 0 | 917 | 902 | 3208 |
| 3 | 1099 | 0 | 806 | 585 | 410 | 989 | 3889 |
| 4 | 894 | 0 | 839 | 651 | 405 | 1098 | 3887 |
| 5 | 927 | 128 | 1060 | 429 | 844 | 4061 | 7449 |
| 6 | 578 | 0 | 982 | 549 | 1294 | 833 | 4236 |
| 7 | 750 | 881 | 1124 | 246 | 756 | 1322 | 5079 |

The IGBT degradation and failure are considered to be random, and hence the duration time $(T_i)$ of the degeneration phase is considered to be a random variable, see figure 5. The y-axis in the figure represents the collector emitter voltage of the IGBT, and the x-axis represents the aging time of the experiment. The figure indicates that the duration time $(T_i)$ in which the IGBT was measured in different volt of $V_{CE}$ is a random variable and it could be represented using the probability density functions summarized in table III.

*B. Modelling Parameters*

In this paper, the Maximum Likelihood method is used to estimate the parameters listed in table III for the Gamma, Exponential and Poisson distribution models. In order to estimate the model parameters, the duration time of 7 IGBTs in degradation process are used as statistical samples, and the six degradation phases are considered to be uncorrelated stochastic processes.

For the Gamma distribution, there are two modelling parameters $\kappa$ and $\theta$ to be estimated. It is assumed that the duration of each degradation phase are uncorrelated and follow the Gamma probability distribution defined in table III. Maximum Likelihood Estimation (MLE) is generically formulated as

$$L(\kappa, \theta) = \prod_{i=1}^{n} f(x_i, \kappa, \theta) \qquad (1)$$

TABLE III.        DISTRIBUTION FUNCTIONS AND MODELLING PARAMETERS [16].

| Models | Density Functions | Parameters |
|---|---|---|
| Gamma | $f(T_i = x) = x^{k-1}\left(e^{\frac{-x}{\theta}} \Big/ \Gamma(k)\theta^k\right)$ | $\kappa, \theta$ |
| Exponential | $f(T_i = x) = \lambda e^{-\lambda x}$ | $\lambda$ |
| Poisson | $f(T_i = x) = \dfrac{e^{-\lambda}\lambda^k}{k!}$ | $\lambda$ |

from which $\kappa$ and $\theta$ can be analytically estimated using [17]

$$\hat{\theta} = \frac{1}{\kappa N}\sum_{i=1}^{N} x_i \qquad (2)$$

$$\hat{\kappa} \leftarrow \kappa - \frac{\ln(\kappa) - \psi(\kappa) - s}{\frac{1}{\kappa} - \psi'(\kappa)} \qquad (3)$$

Table IV summarizes the parameters for 6 uncorrelated degradation phases obtained from MLE.

TABLE IV.        MLE FOR GAMMA PROBABILITY DISTRIBUTION.

| Parameters | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 | Model 6 |
|---|---|---|---|---|---|---|
| $\kappa$ | 25.77 | 0.0883 | 9.3663 | 0.2818 | 5.1662 | 3.3289 |
| $\theta$ | 31.8 | 211.67 | 88 | 1754.4 | 167.3 | 450.2 |

*There is only one modelling parameter $\lambda$ for the Exponential distribution model. Similar to Gamma distribution, the duration of each degradation phase are assumed to be uncorrelated. In this paper, $\lambda$ can be estimated using an analytical MLE solution [17]*

$$\hat{\lambda} = \frac{1}{\bar{x}} = \frac{n}{\sum_{i=1}^{n} x_i} \qquad (4)$$

and the estimated parameters are listed in table V.

TABLE V.        MLE FOR EXPONENTIAL PROBABILITY DISTRIBUTION.

| Parameter | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 | Model 6 |
|---|---|---|---|---|---|---|
| $\lambda$ | 820.7 | 187 | 824.4 | 494.4 | 864.3 | 1498.7 |

FIGURE 5: IGBT Degradation Model.

For Poisson probability distribution, similar to Exponential, $\lambda$ is only the modelling parameter that needs to be estimated. $\lambda$ can be analytically derived, and its MLE can be calculated using [17]

$$\hat{\lambda} = \frac{1}{\bar{x}} = \frac{n}{\sum_{i=1}^{n} x_i} \tag{5}$$

Table VI summarizes the parameters for 6 uncorrelated degradation phases obtained from MLE.

TABLE VI. MLE FOR POISSON PROBABILITY DISTRIBUTION.

| Parameter | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 | Model 6 |
|---|---|---|---|---|---|---|
| $\lambda$ | 820.7 | 187 | 824.4 | 494.4 | 864.3 | 1498.7 |

## IV. PROGNOSTIC APPROACH

In section III, the IGBT degradation models have been developed based on the probabilistic distributions and tuned using the data obtained from the accelerated aging experiments. Based on the degradation profiles shown in figure 3 and 4, the degradation process can be observed by tracking the $V_{CE}$ measurement values. The profile indicates that $V_{CE}$ monotonically increases in discrete steps. In this paper, Monte Carlo simulation is utilized to generate the degradation paths to represent the time durations the IGBT stays in different degradation phases.

Figure 6 shows a block diagram of the prognostic algorithm developed in this paper. The aging data sets were used to train the degradation model depending on what probability distributions are employed. In this paper, the RUL of the IGBT component is predicted by the mean of simulation. The Monte Carlo simulation is used to propagate degradation paths into the future. The RUL can either be the mean or median RUL of the multiple propagated paths. The $V_{CE}$

measurements provide regular updates to the determination or confirmation of the current degradation phase. For every measurement updates, the Monte Carlo simulation is re-run to generate new degradation paths based on the updated measurement and the RUL is then re-calculated from the newly generated paths.

Figure 7 shows an example of IGBT degradation paths. It is assumed that the duration $T_i$ follows Gamma, Exponential or Poisson probability distributions. $S_i$ is the ending time of the degradation phase of the sequence i. t is the elapse time started from the beginning of the experiment. 6 stochastic models are built based on the 6 degradation phases, where the related MLE parameters are summarized in table IV-VI. The RUL prediction could be calculated using the equation

$$RUL_P = S_f - t \tag{6}$$

where $RUL_P$ means the predicted RUL by either Gamma, Exponential or Poisson models, and $S_f$ is the predicted IGBT failure time in the aging process. $S_f$ is also the ending time of the last degeneration phase. The precursor parameter $V_{CE}$ is used to determine in which degradation phase the IGBT is. When the IGBT is stay in the first degeneration phase, then:

$$S_f = \sum_{i=1}^{6} T_i \tag{7}$$

where $T_i$ is generated by Monte Carlo simulation and it means the duration time of relevant phase. So the predicted RUL could be represented as:

$$RUL_P = S_f - t = \sum_{i=1}^{6} T_i - t \tag{8}$$



Fig. 6. RUL Prediction Process.

Fig. 7.  Probability Model for RUL Prediction.

When the IGBT stays in the ordinal of i degeneration phase, because the past generation phase is monitored by measurements which is the $V_{CE}$. So $S_{i-1}$ is a known number could be surveyed from the $V_{CE}$ data recording. Then the predicted RUL could be calculated as:

$$RUL_{P_i} = S_f - t = \left(S_f - S_{i-1}\right) - \left(t - S_{i-1}\right) \tag{9}$$

Then

$$RUL_{P_i} = \sum_{i}^{6} T_i - \left(t - S_{i-1}\right) \tag{10}$$

Where if the simulated Ti generated by Monte Carlo simulation is larger than (t-$S_{i-1}$), it means that the component is still and will be continue in this degradation phase, or the component has ended the generation phase and begin to jump



Fig. 8.  Combined Model for RUL Prediction.

into the next degradation phase.

Note that Monte Carlo simulation is used to generate the duration time of each phase ($T_i$). In order to improve the accuracy of prediction, a large number of multiple runs is needed, which is 500 in this paper. The RUL prediction result uses the mean value and the median value of the distribution.

A combination of statistical properties can be also used to improve prognostic accuracies. One can be based on the duration time of the degradation phase, and the other can be based on the ending time of the degradation phase. Figure 8 illustrates the use of a combined model for predicting the RUL. prediction. The distribution of these two kind of probability models are combined to predict the IGBT RUL. When the IGBT stays in the ordinal of i degeneration phase, the ending time of the phase is $S_i$, and $S_i$ could be represented as:

$$S_i = S_{i-1} + T_i \tag{11}$$

Where $T_i$ is the duration time of this degeneration phase which will be generated by Monte Carlo simulation, and $S_{i-1}$ is the ending time of the prior degradation phase. So $S_{i-1}$ is a known number. $T_i$ is simulated by Monte Carlo. It follows the model distribution based on the duration time of the degradation phase. The probability of $S_i$ is different from the distribution of $T_i$, which means two probability distribution for $S_i$ has been established. Combining these two distributions to predict the RUL is the main solution in this model.

## V. RESULTS

### A. RUL Predictions

The IGBT RUL prediction results are expressed by a series of polylines. The RUL prediction is a continue process from the beginning of the IGBT running to the end of the process when IGBT is becoming failed. The sensor data are recorded at each moment of the IGBT degradation process and the RUL prediction are also carried out at each moment of the degradation process. Hence, the RUL prediction happens through the whole process of the IGBT degradation experiment.

Figure 9 shows an example (i.e. IGBT sample number 1) of RUL prognostic results. The result was computed based on Gamma distribution model. The straight blue, green and yellow dash lines are used as the real and ±10% deviation RUL, respectively. They are used as baselines to indicate how well the prognostic algorithm performs during the test. In figure 9, the red and blue scatter plots are the mean and median values of the RUL prediction. The green and yellow plots are the 90 and 10 percentiles of the Monte Carlo simulated degradation paths. At the beginning of the rendering test, the RUL prediction is lower than the real RUL value, however as the predicted RUL slowly converges to the real value as the operating time is towards the end of component life.

Fig. 9.   Example of RUL Prediction Results using Gamma Distribution.

Figure 10 shows an example (i.e. IGBT sample number 1) of RUL prognostic results based on Exponential distribution model. The red, blue green and yellow scatter plots are the mean, median, 90 and 10 percentiles of the simulated degradation paths. Similar to Gamma distribution model, the predicted RUL slowly converges to the real RUL value as the $V_{CE}$ measurements are rendered towards the end of component life. However, the 10 and 90 percentile bounds are significantly different from the results obtained in the Gamma distribution case. The width of these bounds is too wide and practically become meaningless information-wise in this case. In Exponential distribution, $\lambda$ , i.e. mean, is the only parameter in the model. This model lacks of additional parameter that represents the statistical description equivalent to the standard deviation. This explains the practically irrelevant of the 10 and 90 percentile bounds if these values to be derived based on the Exponential distribution model.



Fig. 10. Example of RUL Prediction Results using Exponential Distribution

An example (i.e. IGBT sample number 1) of Poisson based RUL prognostics results is shown in figure 11. Similarly to Gamma and Exponential results, the predicted RUL slowly converges to the real RUL value as the operating cycles close to the end of component life. The degradation paths are linear with sudden changes reflected the discrete change in the degradation state updated from the $V_{CE}$ measurement. In contrast to Exponential distribution, the 10 and 90 percentile scatter plots lie very close to the mean and median values. These bounds are unrealistic close and practically do not provide meaningful information in terms of confident in the RUL prediction.



Fig. 11. Example of RUL Prediction Results using Poisson Distribution.

Figure 12 shows an example (i.e. IGBT sample number 1) of RUL prognostic results based on the combined (Gamma) distribution model. The combined model gives the prognostic results, i.e. mean and median, as accurate as the Gamma model. However, the 10 and 90 percentiles are much tighter than the results obtained from the Gamma distribution. In this case, these confident bounds lie close to the ±10% deviation of the real RUL. Hence, the 10 and 90 percentile bounds calculated using the combined model practically provides more meaningful confident intervals in comparison to the Gamma, Exponential and Poisson distribution models.



Fig. 12. Example of RUL Prediction Results using Combined Model

*B. Error Analysis*

The errors between the predicted and real RUL values reflect the performance of the IGBT prognostic approach. In this paper, the prediction (or prognostic) error is defined by

$$E_r = RUL_R - RUL_P \qquad (12)$$

where $E_r$ is the error value between the predicted and real values, $RUL_R$ is the real RUL value of an IGBT and $RUL_P$ is the predicted value obtained from the prognostic algorithm. Using equation (12), the prognostic accuracy can be quantitatively calculated using

$$P_{pr} = \frac{E_r}{RUL_R} = \frac{RUL_R - RUL_P}{RUL_R} \qquad (13)$$

In this paper, root mean square error (RMSE) is used to measure the prognostic performance of different probability distribution models. The RMSE can be calculated using the following equations:

$$\text{MSE} = \frac{1}{n}\sum_{i=1}^{n}\left(\widehat{Y_i} - Y_i\right)^2 \tag{14}$$

$$\text{MSE} = \frac{1}{T_f}\sum_{i=1}^{T_f}\left(RUL_{P_i} - RUL_{R_i}\right)^2 = \frac{1}{T_f}\sum_{i=1}^{T_f}(Er_i)^2 \tag{15}$$

$$\text{RMS} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}\left(\widehat{Y_i} - Y_i\right)^2} = \sqrt{\frac{1}{T_f}\sum_{i=1}^{T_f}(Er_i)^2} \tag{16}$$

Figure 13 and 14 summarize the mean and median based RMSEs of different probability distribution models tested against 7 IGBT data samples. For the mean based RUL, the Poisson model performs better than other models on most of the test samples; Its RMSE is less than other models. However, for the median based RUL, the Exponential model has less RMSE values in comparison to other models. Comparing the RMSE of mean and median predicted RULs, the IGBT test sample number 1 and 7 have the smallest RMSE for combining model. The rest of other IGBT test samples have similar prognostic performance in terms of how different probabilistic models performed in relation to each other.



Fig. 13. RMSE for Mean Based of Predicted RULs.



Fig. 14. RMSE for Median Based of Predicted RULs.

## VI. CONCLUSIONS

The main contribution of this paper is the development and implementation of a prognostics framework for IGBTs, and a prognostic algorithm with Monte Carlo simulation and with the collector emitter voltage as a precursor parameter is developed. According to the IGBT failure mechanism and degradation characterization, the IGBT degeneration models are built. Monte Carlo simulation method and the precursor parameter, collector emitter voltage ($V_{CE}$), are integrated to develop the prognostic algorithm on predicting the IGBT RUL.

Gamma, Exponential, Poisson distribution and the combining distribution models are established, and Monte Carlo simulation is utilized in the algorithm to computing the IGBT remaining useful life. The collector emitter voltage ($V_{CE}$) is used as the precursor parameter used in the prognosis.

Comparing with the results of RUL prediction with different models, the mean value of the RUL prediction and the median value of the RUL prediction presents a different accuracy. Different models also perform their preference to different IGBT on the RUL prediction. The implementation of the developed prognostics framework could be applied to provide advance warning of failures thereby preventing costly power electronics system downtime and failures.

The combining model can perform much more efficient RUL prediction results in some IGBTs, and the combined model in this paper is only based on the Gamma distribution, much more combining models based on different probability distribution could be established and implemented in IGBT RUL prediction, and a comparative analysis between these combining models is beneficial to the IGBT prognostic.

References

[1] J. Celaya, B. Saha and P. Wysocki, "Prognostics for Electronics Components of Avionics Systems," in IEEE Aerospace Conference, Big Sky, MT, 2009.

[2] J. Celaya, A. Saxena, C. Kulkarni, S. Saha and K. Goebel, "Prognostics Approach for Power MOSFET under Thermal-Stress Aging," in Reliability and Maintainability Symposium (RAMS), Reno, NV, 2012.

[3] M. Pecht and R. Jaai, "A prognostics and health management roadmap for information and electronics-rich systems," Microelectronics Reliability, pp. 317-323, 2010.

[4] N. Patil, D. Das, K. Goebel and M. Pecht, "Identification of Failure Precursor Parameters for Insulated Gate Bipolar Transistors (IGBTs)," in INTERNATIONAL CONFERENCE ON PROGNOSTICS AND HEALTH MANAGEMENT, Denver, CO, 2008.

[5] B. Saha, J. Celaya, P. Wysocki and K. Goebel, "Towards Prognostics for Electronics Componenets," Aerospace conference, IEEE, no. 2009, pp. 1-7, 2009.

[6] B. Saha, K. Goebel and S. Poll, "Modeling Li-ion Battery Capacity Depletion in a Particle Filtering Framework," in Annual Conference of the Prognostics and Health Management Society, San Diego, 2009.

[7] M. Pecht, Prognostics and Health Management of Electronics, New York, NY: Wiley-Interscience, 2008.

[8] N. Patil, D. Das and M. Pecht, "Aprognostic approach for non-punch through and field stop IGBTs," Special section on International Seminar on Power Semiconductors 2010, vol. 52, no. 3, pp. 482-488, 2010.

[9] Y. Murphey, A. Masur and B. Chen, "A Fuzzy System for Fault Dignostics in Power Electronics Based Brake-by-wire System," in Annual Meeting of the North American Fuzzy Information Processing Society, 2005, New York.

[10] NASA, "NASA," 24 June 2013. [Online]. Available: http://ti.arc.nasa.gov/tech/dash/diagnostics-and-prognostics.

[11] J. Celaya, P. Wysocki, V. Vashchenko and S. Saha, "Accelerated Aging System for Prognotics of Power Semiconductor Devices," in Autotestcon, 2010 IEEE, Orlando, FL, 2010.

[12] G. Sonnenfeld, K. Goebel and J. Celaya, "An Agile Accelerated Aging, Characterization and Scenario Simulation System for Gate Controlled Power Transistors," in IEEE AUTOTESTCON, Salt Lake City, UT, 2008.

[13] N. Patil, D. Das, K. Goebel and M. Pecht, "Failure Precursors for Insulated Gate Bipolar Transistors (IGBTs)," IEEE TRANSACTIONS ON RELIABILITY, pp. 271 - 276, 2009.

[14] J. Celaya, P. Wysocki and K. Goebel, "IGBT accelerated aging data set," NASA Ames Prognostics Data Repository, Moffett Field, CA, 2009.

[15] A. Alghassi, S. Perinpanayagam and I. Jennions, "A Simple State-Based Prognostic Model for Predicting Remaining Useful," in Power Electronics and Applications (EPE), 2013 15th European Conference, Lille, France, 2013.

[16] J. Lawless, Statistical Model and Methods for Lifetime Data, John Wiley & Sons, 2011.

[17] R. B. Millar, Maximum Likelihood Estimation and Inference, John Wiley & Sons, 2011.

# An Automated approach for Preventing ARP Spoofing Attack using Static ARP Entries

Ahmed M.AbdelSalam

Information Technology Dept.
Faculty of Computers and
Information, Menofia University
Menofia, Egypt

Wail S.Elkilani

Computer Systems Dept.
Faculty of Computers and
Information, Ain Shams University
Cairo, Egypt

Khalid M.Amin

Information Technology Dept.
Faculty of Computers and
Information, Menofia University
Menofia, Egypt

*Abstract*—**ARP spoofing is the most dangerous attack that threats LANs, this attack comes from the way the ARP protocol works, since it is a stateless protocol. The ARP spoofing attack may be used to launch either denial of service (DoS) attacks or Man in the middle (MITM) attacks. Using static ARP entries is considered the most effective way to prevent ARP spoofing. Yet, ARP spoofing mitigation methods depending on static ARP have major drawbacks. In this paper, we propose a scalable technique to prevent ARP spoofing attacks, which automatically configures static ARP entries. Every host in the local network will have a protected non-spoofed ARP cache. The technique operates in both static and DHCP based addressing schemes, and Scalability of the technique allows protecting of a large number of users without any overhead on the administrator. Performance study of the technique has been conducted using a real network. The measurement results have shown that the client needs no more than one millisecond to register itself for a protected ARP cache. The results also shown that the server can a block any attacker in just few microsecond under heavy traffic.**

*Keyword—component; layer two attacks; ARP spoofing; ARP cache poisoning; Static ARP entries*

## I. INTRODUCTION

The evolving of computer networks, and the variety of its services and applications, has increased the users need for LANs [1] and the security of LANs also become a more concern. An essential part of successful communication between users within LAN is the Address Resolution Protocol (ARP) [2]. ARP is specified in RFC 826 [3] to allow hosts to resolve network layer address (IP) to datalink layer address (MAC) [3]. Although the importance of ARP protocol for communication in LAN, it formulates the most dangerous attacks threating LANs. ARP spoofing or ARP cache poisoning are the two main attacks threating the ARP protocol operations [4].

ARP Spoofing is a hacking technique to send fake ARP request or ARP reply, ARP spoofing problem comes from the way the ARP protocol works [5]. Since the ARP protocol is a stateless protocol that receives and processes ARP replies without issuing ARP request [6], the ARP cache can be infected with records that contain wrong mappings of IP-MAC addresses. ARP spoofing can be used to launch one of two different attack categories [7]: Denial of Service (DoS) attacks or Man in the Middle (MITM) attacks

Several solutions have been proposed to mitigate the ARP spoofing, but each has its limitations [7]. The solutions have been classified into five different categories [8]:

- Modifying ARP using cryptographic techniques

These solutions add some cryptographic features to the ARP protocol, but will not be compatible with the standard ARP and affect the protocol performance.

- Kernel-based patching

The technique adds a patch to the operating system kernel in order to prevent ARP spoofing attacks, but the problem is that not all operating systems can be patched and it may become incompatible with the standard ARP protocol.

- Securing switch Ports

Use the switch port security or Dynamic ARP inspection (DAI) option to prevent ARP spoofing. However its ability of preventing ARP spoofing easily, the cost of implementing such solution may not be acceptable by most of the organizations.

- ARP spoof detection & protection software

Programs or tools developed to prevent ARP spoofing attacks, but the experimental results have shown there ineffectiveness in protection.

- Manually configuring static ARP entries

The most basic and effective way to prevent ARP spoofing [1] [6] [9] is adding static ARP entries at each host. However this solution cannot be easily managed and cannot scale well specially in organizations that have large number of users and require a heavy workload on the network administrator.

In this paper, a scalable technique to prevent ARP spoofing attacks, which automatically configures static ARP entries is proposed. It overcomes the problems of the solutions of the techniques that use static entries. The remaining part of the paper is organized as follows: Section II surveys background and related work. Section III shows the details of the proposed method. The experimental results are discussed in section IV. Finally section V concludes the paper.

## II. BACKGROUND AND RELATED WORK

### A. Address resolution Protocol

ARP is specified in RFC 826 [3] as a protocol that provides dynamic mapping from an IP address to the corresponding MAC address to grant successful communication between users within LAN. ARP messages are classified as request and reply message. When a user has packet to transmit, it will send a broadcast ARP request asking about the MAC address for a certain IP. The machine, recognizing the IP address as its own address, returns an ARP reply containing its MAC address. The mapping will be saved in the device ARP cache [2].

### B. ARP cache

ARP cache is a table of recently resolved IP addresses and their corresponding MAC addresses. The ARP cache is checked first before sending an ARP Request frame. ARP cache entries can be dynamic or static [4].

- *Static ARP cache entries:* are permanent and manually added records using a TCP-IP utility. Static ARP cache entries are used to provide ARP requests for commonly used local IP addresses. The problem with static ARP entries is that they have to be manually updated when network interface equipment changes [10].

- *Dynamic ARP cache entries:* are entries learned by ARP protocol and have a time-out value associated with them to remove entries from the cache after a specified period of time [10].

### C. ARP spoofing

ARP is a stateless protocol that uses ARP replies to update ARP cache using wrong or spoofed mappings [11]. As mentioned before, ARP spoofing can be used to launch [7] either one of the following attack categories:

- *Man in the Middle (MITM):* An attacker deceives both ends of communication and fills their ARP cache with wrong IP-MAC mapping. As a matter of fact, it inserts itself between the two ends of communication. Hence, it will gain a copy of every bit sent between them.

- *Denial of service (DoS):* An attacker fills the ARP cache of victim with wrong IP-MAC mapping, so every packet sent from victim will be sent to the wrong MAC.

### D. Related Work

As mentioned previously, solutions attempting to prevent ARP spoofing attack using the static ARP cache entries are very efficient. Yet this category of solutions has some major problems [7] [8]: (1) overhead required for manual configuration of static entries, (2) Limited scalability for large networks, and (3) Ability to work in static and DHCP based networks.

In the following, we will survey several methods belonging to this category along with their drawbacks.

The DAPS (Dynamic ARP spoof Protection System) technique suggested in [8] is a solution to ARP spoofing that snoops DHCP packets and use them as vaccines. Yet this technique doesn't scale well for those network that use static IP addressing scheme and also vaccines will be invalid if DHCP starvation attack occurs.

In [12], the NIDPS (Network Intrusion Detection and Prevention System) technique is suggested have a server collecting IP-MAC mappings from users using small agents. These mappings will be then used as static ARP entries to correct any wrong mapping detected. However, agents aren't authenticated to the server. Moreover, it detects only attacks from its LAN segment. Also, the server examines every packet going in or out the LAN segment. Finally, it waits for the attack to occur and then try to solve it.

Xiangning et al. [13] has proposed a technique that expands the snort preprocessors plug-ins by adding an ARP detection module. The proposed technique doesn't scale well in large networks due to the need of manual configuration of the static mappings at the server. It also doesn't work in DHCP based networks.

A solution to ARP spoofing using a server is proposed in [14]. The server will get mappings for the network users from the DHCP server. It replies also to ARP requests. Unfortunately, this solution works only in DHCP networks. Also, it is not compatible with the standard ARP. Moreover, if DHCP starvation occurs, all the server information will be invalid.

Ai-zeng Qian [15] proposed a technique to prevent ARP spoofing by using static ARP entries but the technique still doesn't work with dynamic networks using DHCP addressing. The administrator must assign all IP addresses along with their MAC to the server so it will be not visible for large scale network.

A method is suggested in [16] to solve ARP spoofing problem using snort IDS and static ARP entries. Yet, it still needs the administrator to add the static mappings manually. Also, it works only in static networks.

Table 1 compares the proposed algorithm with the previous solutions, the comparison criteria includes if it works in DHCP and static networks, ability to prevent attacks, scalability, and if manual or automatic configuration of the static entries is used.

TABLE I.    COMPARISON BETWEEN THE PROPOSED ALOGORITHM AND PREVIOUS SOLUTIONS

| Technique | Static | DHCP | Prevention | Scalability | Automatic |
|---|---|---|---|---|---|
| *Proposed* | ✓ | ✓ | ✓ | ✓ | ✓ |
| *DAPS [7]* | | ✓ | ✓ | | ✓ |
| *NIDPS [12]* | ✓ | ✓ | | ✓ | ✓ |
| *Xiangning [13]* | ✓ | | ✓ | | |
| *Ortega [14]* | | ✓ | ✓ | | ✓ |
| *Ai-zeng [16]* | ✓ | | ✓ | | |
| *Xiangdong [17]* | ✓ | | ✓ | | |

### III.    THE PROPOSED METHOD

The proposed technique is a client-server protocol that prevents ARP spoofing by automatically configuring static ARP entries. The protocol works in both static and DHCP networks.

Moreover, it can work in large-scale networks without any overhead on the administrator. In addition, the technique doesn't require special hardware to be deployed, as any host can work as ARP server.

The protocol proposed defines three different messages:

A. *Register Message: is a unicast message sent from the client to the server. It contains its IP and MAC address. Also it includes a hashed authentication key.*

B. *Update Message: is a broadcast notification message sent from the server to all users in the network indicating that a new user has entered the network. It also contains the IP and MAC address of that new user.*

C. *Register Response Message: is a unicast message sent from the server to the new user. It contains all static ARP entries of users successfully registered at the server.*

The protocol also defines two different entities:

*a) ARP Client: is a software installed on user's machines. It fulfills the following*

- Automatically get the IP and MAC address of the user and use them to send register message to the server.

- Receive update and register response messages from the server.

- Verify that update or register response messages received are coming from a trusted server.

- Use the IP and MAC pairs received in the update or register response message to add static ARP entries to the user ARP cache.

*b) ARP Server: is a server software that can be installed on any device in the network. It can also be installed on a dedicated server, and has the following functions:*

- Receive register messages from the ARP clients.

- Verify that the message is coming from a trusted user.

- Make use of the IP and MAC pairs encapsulated within the register message to create a list of trusted users in the network.

- Send broadcast update message to notify them that a new user has come to the network.

- Send register response message to the new users.

- Take the proper action regarding users who try to violate the protocol security rules.

The proposed protocol defines two different algorithms for the client and server in order to prevent the ARP spoofing attack

*1) Client Algorithm*
The client algorithm described in Algorithm 1 adds static entry for the server in the client ARP cache to avoid the rogue server threat. Furthermore, it obtains the user IP and MAC address automatically to make the user has no opportunity to send fake information to the server.

---

**Algorithm 1:** *ARP_Client*

*Step 1*: Add static ARP entry for the server.
*Step 2:* Automaitcally get user IP and MAC address
*Step 3:* Formulate the Register Message
*Step 4:* Send the register message to the server.
*Step 5:* Listen to updates from the server
*Step 6:* **if** message received from the server **then**
    Extract the source IP
    **if** source IP = Server IP **then**
        Extract Key, IP, MAC
        **if** received key is correct **then**
            **if** similar MAC in ARP Cache **then**
                Delete this record
            **else**
                Add static ARP entry using extracted IP and MAC address
                Return to step 5
            **end if**
        **else**
            Discard the message
            Return to step 5
        **end if**
    **else**
        Discard the Message
        Return to step 5
    **end if**
**else**
    Return to step 5
**end if**

---

The algorithm checks the source IP address of the received message to be sure that it is coming from the trusted server. It only accepts the IP and MAC addresses encapsulated in the message if the key is correct.

In order to work in DHCP networks where IP and MAC mappings are frequently changing, the algorithm searches for the MAC encapsulated in the message. If matched map is found, it will be changed to the new mapping. Otherwise a new mapping will be added.

Finally if any of the conditions are not met, the algorithm will discard the message and return to listen for another message from the server.

*2) Server Algorithm*
The server algorithm, described in Algorithm 2, listens to incoming register messages from the clients, checks the hash code to be sure that the message is coming from a trusted host. Users are given only three trials to send the correct hash code. If it fails to send the correct hash code within the three trials, the server will block this user.

The blocking action depends on the addressing scheme being used, for DHCP networks, the MAC address of the user will be added to DHCP deny list. Hence, it will not be able to obtain IP configuration from the DHCP server again, for static networks, the server will prevent traffic from this user to reach the server by obstructing its IP address.

---

***Algorithm 2: ARP_Server***

---

***Step 1:*** Add static entry for itself
***Step 2:*** Listen to users Register Messages
***Step 3: if*** register message received from user **then**
      **if** the hash code matched **then**
          **if** wrong trials less than 3 **then**
              **if** similar MAC in ARP cache **then**
                  update this record
              **else**
                  delete this record
              **end if**
              send update message to all users
              send register response message to the new user
              return to step 2
          **else**
              discard the message
              return to step 2
          **end if**
      **else if** it has wrong previous wrong trials **then**
          Increment wrong trials for that MAC
          **if** wrong trials equal 3 **then**
              Add to DHCP deny list Or Block this IP
          **end if**
          Discard the message
          Return to step 2
      **else**
          Add to sucpicious list
          Discard the message
          Return to step 2
      **end if**
    **else**
      Return to step 2

   **end if**

If user wrong trials have reached three or more, its traffic will be discarded even if it gets the correct hash code, and the administrator is the only one who has the ability to remove it from the block list. It is to be noted that we block undesirable users on the network layer instead of the usual blocking criteria using TCP. This enables the protocol proposed to stop DoS attacks completely.

If the key is correct and the number of wrong trials doesn't reach the threshold, the server will search its ARP cache for matching between MAC address encapsulated in the register message received and MAC address in ARP cache. This gives the algorithm the ability to work with DHCP based networks. In turn, it prevents an intruder having the hash code to spoof all ARP cache entries. As a matter of fact, it can only spoof one at a time. If it tries to spoof another one the old spoofed entry will be deleted.

In case a user tries to register with wrong hash code for the first time, it will be inserted in the suspicious users list. When its wrong trials reaches three, it will be moved to the blocked list.

User who has successfully registered at the server will receive a register response message contains the IP and MAC addresses of all successfully registered users to add them as static ARP entries. Moreover, all other users will receive an update message contains IP and MAC address of the new user to add it as a static entry in their ARP cache.

Using the client and server algorithms, every user in the network will have its ARP cache filled with static ARP entries

for all other users in the network. Hence, it will not suffer from the ARP spoofing problem again. And everything is done automatically without any overload on the administrator; this gives the algorithm a greater scalability.

## IV. EXPERIMENTAL RESULTS

Experimental measurement has been chosen to evaluate the performance of the proposed algorithm. The faculty of computers and information, Menofia University (Menofia is one of the districts of Egypt) network, shown in fig.1, is used to conduct the measurements. The network consists of three separate LANs. Each LAN ends with an edge switch. The LANs are connected through a core switch.

LAN 1 consists of 19 users and an ARP server. LAN 2 consists of 13 users and an application server offering web browsing, FTP, and mail services. LAN 3 is a network of 16 users and an Asterisk VOIP server for voice over IP calls between network users.

All PCs are core i5 processor with 4 GB of RAM. The edge switches are cisco catalyst 2960 switch with 24 ports. The core switch is cisco catalyst 4006 switch. Also the wireshark software is used in measuring the values.

The response time metric has been chosen to evaluate the performance of the algorithm. The response time is measured at the different stages of the algorithm and at both sides of the protocol (client and server) taking in mind the different parameters affecting the response time values. These measures speed, reliability, and robustness of the algorithm. Hence, it proves the algorithm efficiency.



Fig. 1. Faculty of computers and Information, Menofia university network

### A. Server Side Measures

#### 1) Authentication time ($T_{auth}$)

It represents the amount of time needed to authenticate a trusted user. Fig.2 shows the authentication time values. The X Axis represents the number of simultaneous attackers trying to authenticate using wrong key. The Y Axis represents the

authentication time values in nanoseconds. The different colors represent the number of users trying to authenticate at the same time.



**Authentication Time**

| | 0 | 1 | 2 | 4 | 8 | 16 |
|---|---|---|---|---|---|---|
| 1 | 5822 | 5813 | 5829 | 5837 | 5841 | 5846 |
| 2 | 6103 | 6111 | 6119 | 6124 | 6129 | 6137 |
| 4 | 6279 | 6283 | 6289 | 6297 | 6306 | 6316 |
| 8 | 6809 | 6815 | 6823 | 6832 | 6842 | 6852 |
| 16 | 7907 | 7917 | 7929 | 7941 | 7958 | 7983 |
| 32 | 8113 | 8124 | 8139 | 8156 | 8169 | 8186 |

Fig. 2. Authentication Time ($T_{auth}$) in nanoseconds versus number of attackers for the different number of users

*2)* *Acceptance time ($T_{acc}$)*

It represents the time spent by the server to accept an authenticated user. Fig.3 shows the acceptance time values. The X Axis represents the total number of ARP cache records. The Y Axis represents the acceptance time values in nanoseconds. The different colors represent the number of users trying to authenticate at the same time.



**Acceptance Time**

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 3631 | 5866 | 8660 | 15085 | 25422 | 46933 |
| 2 | 3703 | 5993 | 8991 | 15423 | 25399 | 48003 |
| 4 | 3957 | 6378 | 9573 | 16234 | 26847 | 49334 |
| 8 | 4208 | 6911 | 10412 | 17411 | 28003 | 50867 |
| 16 | 4421 | 7423 | 11517 | 18534 | 29211 | 51997 |
| 32 | 4987 | 8009 | 12786 | 19765 | 31116 | 54161 |

Fig. 3. Acceptance Time ($T_{acc}$) in nanoseconds versus number of ARP cache records for different number of users

*3)* *Registration time ($T_r$)*

It represents the time taken by the server to add an entry for a new user in the ARP cache. Fig.4 shows the registration time

values. The X-Axis represents the total number of ARP cache records. The Y-Axis represents the registration time values in nanoseconds. The different colors represent the number of users trying to register themselves at the same time.



**Registration Time**

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 5028 | 6704 | 7263 | 11733 | 16761 | 26819 |
| 2 | 5287 | 6998 | 7734 | 12645 | 17432 | 28005 |
| 4 | 5517 | 7305 | 8411 | 14114 | 19634 | 30233 |
| 8 | 5919 | 7913 | 9831 | 15654 | 22344 | 32946 |
| 16 | 6108 | 8726 | 11531 | 16973 | 24987 | 34988 |
| 32 | 6734 | 10113 | 13383 | 18321 | 28166 | 37213 |

Fig. 4. Registration Time ($T_r$) in nanoseconds versus number of ARP cache records for the different number of users

*4)* *Update time ($T_u$)*

It represents the time needed by the server to notify all network users that a new user has entered the network. Fig.5 shows the update time values. The X Axis represents number of simultaneous users trying to communicate with the server. The Y Axis represents the update time values in nanoseconds.



**Update time**

| | 1 | 2 | 4 | 8 | 16 | 32 |
|---|---|---|---|---|---|---|
| | 32126 | 33107 | 34311 | 35876 | 37323 | 39765 |

Fig. 5. Update Time ($T_u$) in nanoseconds versus the number of users

*5)* *Server Convergence time ($T_{Con}$)*

It represents the time taken by the server to send to the new user its ARP cache as a register response message. Fig.6 shows the server convergence time values. The X Axis represents the total number of ARP cache records. The Y Axis represents the server convergence time values in nanoseconds. The different colors represent the number of users trying to communicate with the server at the same time.

## Server convergence time

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 6704 | 11453 | 17879 | 36317 | 74819 | 117333 |
| 2 | 6974 | 12003 | 18973 | 38114 | 76101 | 118454 |
| 4 | 7734 | 13769 | 20541 | 41226 | 78311 | 119621 |
| 8 | 8947 | 16278 | 23001 | 44992 | 81003 | 120872 |
| 16 | 10765 | 19231 | 26742 | 49078 | 83867 | 121877 |
| 32 | 11874 | 21324 | 30991 | 56429 | 86234 | 123117 |

Fig. 6. Server convergence Time $(T_{Con})$ in nanoseconds versus number of ARP cache records for the different number of users

### 6) Detection time ($T_d$)

It represents the time in nanoseconds spent by the server to detect a new attacking user and adds it to the suspicious users list, Fig.7 shows the detection time values.

The X Axis represents the number of suspicious user's list records. The Y Axis represents the detection time values in nanoseconds. The different colors represent the number of users trying to communicate with the server at the same time.

## Detection Time

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 4749 | 7263 | 9777 | 16203 | 28495 | 50006 |
| 2 | 4897 | 7892 | 10223 | 17102 | 29142 | 50873 |
| 4 | 5231 | 8321 | 10983 | 18001 | 30929 | 52147 |
| 8 | 5713 | 9107 | 12314 | 19635 | 32765 | 55725 |
| 16 | 6003 | 10111 | 13867 | 21245 | 34567 | 58836 |
| 32 | 6878 | 11231 | 15311 | 23781 | 37116 | 63139 |

Fig. 7. Detection Time $(T_d)$ in nanoseconds versus number of record in suspicious users list for the different number of users

### 7) Blocking time ($T_b$)

It represents the time spent by the server to add a user to the blocked users list. Fig.8 shows the blocking time values. The X Axis represents the total number of records in the suspicious users list.

The Y Axis represents the detection time values in nanoseconds. The different colors represent the number of users trying to communicate with the server at the same time.
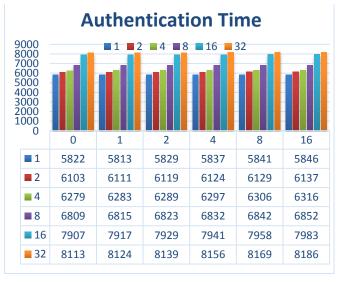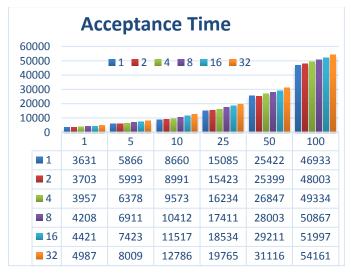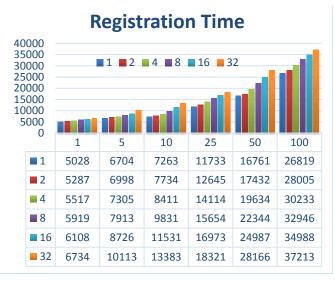
## Blocking Time

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 2793 | 4749 | 6984 | 14526 | 25142 | 47492 |
| 2 | 2987 | 5002 | 7256 | 14987 | 25976 | 48543 |
| 4 | 3347 | 5967 | 7967 | 15789 | 27127 | 50129 |
| 8 | 3983 | 6234 | 8675 | 16430 | 28965 | 52341 |
| 16 | 4673 | 6871 | 9433 | 17987 | 31006 | 54389 |
| 32 | 5987 | 7954 | 10856 | 19768 | 34569 | 57345 |

Fig. 8. Blocking Time $(T_b)$ in nanoseconds versus number of suspicious users list records for different number of users

It can be noted from Fig 2 to 8 that for any measured time $(T_{auth}, T_{acc}, T_r, T_u, T_{con}, T_d, T_b)$, the time needed per user is nearly constant for any number of users for the same number of records.

### B. Client side measures

#### 1) Client Acceptance time ($TC_{acc}$)

It represents the time needed by the client to be sure that the server accepted its register request message. It will be calculated using equation (1):

$$TC_{acc} = RTT + T_{auth} + T_{acc} + T_r + T_u \quad (1)$$

Where $TC_{acc}$ is the client acceptance time, $RTT$ is the round trip time, $T_{auth}$ is the server authentication time, $T_{acc}$ is the sever acceptance time, $T_r$ is the server registration time, and $T_u$ is the server update time.

The $RTT$ value depends of the nature of traffic workload. The users are divided into four groups and every group of users is using one or more of the services: mail, file transfer, VoIP call, or web browsing. The workload depends on the number of services used by every group. The results were taken for 3 different types of workloads: light, normal, and heavy workload. Light workload represents one service usage. Two services are considered for normal workload and three for heavy workload [17].

#### a) Light traffic workload

Fig.9 shows the Client acceptance time values in light traffic workload, The X Axis represents the number of records in Server ARP cache. The Y Axis represents the Client acceptance time in Microseconds, the different colors represent the number of users trying to communicate with the server at the same time.

**Client Acceptance Time**

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 136.607 | 188.509 | 170.878 | 198.781 | 186.15 | 246.724 |
| 2 | 161.2 | 183.209 | 171.951 | 161.299 | 211.067 | 237.252 |
| 4 | 177.064 | 173.277 | 155.584 | 198.956 | 190.098 | 248.194 |
| 8 | 176.812 | 171.515 | 185.942 | 202.773 | 195.065 | 256.541 |
| 16 | 147.759 | 193.389 | 175.3 | 178.771 | 202.479 | 251.291 |
| 32 | 153.599 | 174.011 | 208.073 | 208.007 | 244.216 | 250.325 |

Fig. 9. Client acceptance Time（$TC_{acc}$) in nanoseconds versus number ARP cache records for different number of users in light worload

### b) Normal workload

Fig.11 shows the Client acceptance time values in normal traffic workload, The X Axis represents number of record in Server ARP cache The Y Axis represents the client acceptance time values in microseconds, and the different colors represent the number of users trying to communicate with the server at the same time.

**Client Acceptance Time**

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 424.607 | 382.509 | 445.878 | 419.781 | 393.15 | 421.724 |
| 2 | 438.2 | 353.209 | 378.951 | 411.299 | 466.067 | 471.252 |
| 4 | 410.064 | 400.277 | 458.584 | 386.956 | 442.098 | 500.194 |
| 8 | 433.812 | 384.515 | 370.942 | 427.773 | 459.065 | 467.541 |
| 16 | 392.759 | 438.389 | 398.3 | 462.771 | 416.479 | 527.291 |
| 32 | 441.599 | 392.011 | 444.073 | 392.007 | 476.216 | 478.325 |

Fig. 10. Client acceptance Time（$TC_{acc}$) in nanoseconds versus number ARP cache records for different number of users in normal worload

### c) Heavy traffic workload

Fig.11 shows the Client acceptance time values in Heavy traffic workload. The X Axis represents the number of records in the server ARP cache. The Y Axis represents the client acceptance time values in microseconds. The different colors represent the number of users trying to communicate with the server at the same time.
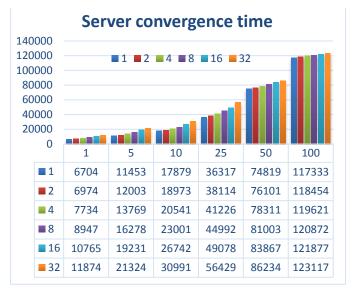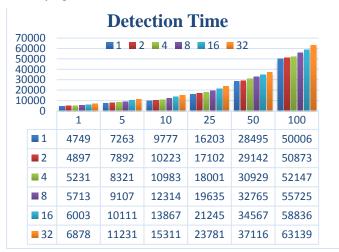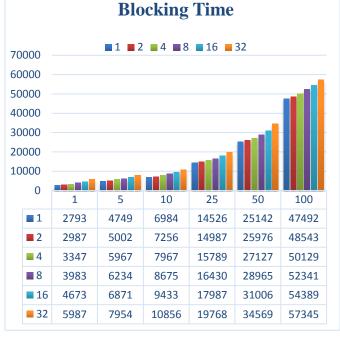
**Client Acceptance Time**

| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 46.607 | 942.509 | 977.878 | 875.781 | 969.15 | 1048.72 |
| 2 | 977.2 | 907.209 | 1010.95 | 920.299 | 926.067 | 1110.25 |
| 4 | 957.064 | 1021.28 | 894.584 | 953.956 | 1075.1 | 1038.19 |
| 8 | 946.812 | 876.515 | 973.942 | 936.773 | 897.065 | 1103.54 |
| 16 | 1001.76 | 923.389 | 890.3 | 1003.77 | 955.479 | 1085.29 |
| 32 | 892.599 | 999.011 | 947.073 | 887.007 | 1010.22 | 964.325 |

Fig. 11. Client acceptance Time（$TC_{acc}$) in nanoseconds versus number ARP cache records for different number of users in heay worload

### 2) Client convergence time ($Tc_{Con}$)

It represents the amount of time needed by the client to process the register response message and add static ARP entries using IP and MAC pairs encapsulated within the message. Fig.12 shows the client convergence time values. The X Axis represents the total number IP and Mac pairs encapsulated in the register response message. The Y Axis represents the client convergence time values in nanoseconds. The different colors represent number of users trying to communicate with the server at the same time.

**Client convergence time**

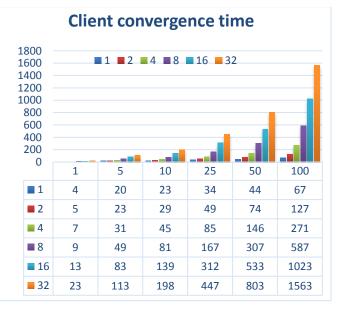| | 1 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|
| 1 | 4 | 20 | 23 | 34 | 44 | 67 |
| 2 | 5 | 23 | 29 | 49 | 74 | 127 |
| 4 | 7 | 31 | 45 | 85 | 146 | 271 |
| 8 | 9 | 49 | 81 | 167 | 307 | 587 |
| 16 | 13 | 83 | 139 | 312 | 533 | 1023 |
| 32 | 23 | 113 | 198 | 447 | 803 | 1563 |

Fig. 12. Client Convergence Time （$Tc_{Con}$) in nanoseconds versus number of ARP cache records for the different number of users

## V. CONCLUSION

In this paper, a solution to the problem of ARP spoofing has been proposed, the solution is an automatic and scalable method of configuring static ARP entries instead of manually configuring. The solution solves the main problems related to this category of solutions Usage of static entries, automation, scalability, manageability, prevention, and cost are the main features of the proposed method. The proposed method has defined two separate algorithms, one for the client, and the other for the server. Experimental evaluation was conducted on the LAN network of the faculty of computers and information, Menofia university of Egypt. The response time metric is used to evaluate the algorithm. The values of the response time were measured at the different stages of the algorithm. Also different types of traffic workloads were used during the measuring the response to show the effect volume of traffic on the response time values. The results prove how fast and accurate the proposed algorithm is since any new user needs less than one millisecond to be safe from ARP problem for heavy workloads.

### REFERENCES

[1] Yafeng Xu and Shuwen Sun , "The study on the college campus network ARP deception defense," 2010 2nd International Conference on Future Computer and Communication (ICFCC), 3(1), pp. 465-467, May 2010.

[2] R. W. Stevens. TCP/IP Illustrated, Volume 1: The Protocols. Addison–Wesley Professional Computing Series, January 1994.

[3] D. Plummer. An Ethernet address resolution protocol, Nov. 1982. RFC 826.

[4] Mohamed Al-Hemairy, Saad Amin, and Zouheir Trabelsi, "Towards More Sophisticated ARP Spoofing Detection/ Prevention Systems in LAN Networks," *2009 International Conference on the Current Trends in Information Technology (CTIT),* pp.1-6, December 2009.

[5] Hu Xiangdong, Gao Zhan, and Li Wei "Research on the Switched LAN Monitor Mechanism and its Implementation Method based on ARP spoofing," *International Conference on Management and Service Science.( MASS '09)*, pp. 1-4, Sept. 2009.

[6] Marco Antônio Carnut and João J. C. Gondim, "ARP spoofing detection on switched ethernet networks: a feasibility study," *5ᵗʰ Symposium on Security in Informatics held at Brazilian Air Force Technology Institute*, November 2003.

[7] Cristina L. Abad and Rafael I. Bonilla, "An Analysis on the Schemes for Detecting and Preventing ARP Cache Poisoning Attacks," *27th International Conference on Distributed Computing Systems Workshops, 2007. (ICDCSW '07)*, page(s): 60, June 2007.

[8] Somnuk Puangpronpitag and Narongrit Masusai, "An Efficient and Feasible Solution to ARP Spoof Problem," *6th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology, 2009. (ECTI-CON 2009)*, 3(1), pp. 910—913, May 2009.

[9] S. Whalen, "An introduction to ARP spoofing," 2600: The Hacker Quarterly, 18(3), 2001, (accessed 13-9-2012). [Online].:http://servv89pn0aj.sn.sourcedns.com/gbpprorg/2600/arp spoofing intro.pdf

[10] http://technet.microsoft.com/en-us/library/cc958841.aspx. ARP Cache, (accessed May 8, 2013).

[11] Zouheir Trabelsi and Wassim El-Hajj, "Preventing ARP Attacks using a Fuzzy-Based Stateful ARP Cache," IEEE International Conference on Communications.( ICC '07), pp. 1355 -1360, June 2007.

[12] Dr. S. G. Bhirud and Vijay Katkar, "Light Weight Approach for IP-ARP Spoofing Detection and Prevention," *2011 Second Asian Himalayas International Conference on Internet (AH-ICI)*, page(s):1-5, November 2011.

[13] Xiangning HOU, Zhiping JIANG, and Xinli TIAN, "The detection and prevention for ARP Spoofing based on Snort," *2010 I*

[14] Andre P. Ortega, Xavier E. Marcos, Luis D. Chiang and Cristina L. Abad, " Preventing ARP cache poisoning attacks: A proof of concept using OpenWrt," Latin American Network Operations and Management Symposium. (LANOMS), pp. 1-9, Oct. 2009.

[15] Ai-zeng Qian, "The Automatic Prevention and Control Research of ARP Deception and Implementation," *2009 WRI World Congress on Computer Science and Information Engineering, ,* 2(1), pp. 555-558, April 2009.

[16] Boughrara, A.; Mammar, S., "Implementation of a SNORT's output Plug-In in reaction to ARP Spoofing's attack," 2012 6th International Conference on Sciences of Electronics Technologies of Information and Telecommunications (SETIT), pp.643,647, 21-24 March 2012

[17] R. K. Jain, "The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling," Prince Hall, April 1991.

# For an Independent Spell-Checking System from the Arabic Language Vocabulary

Bakkali Hamza
Telecom and Embedded Systems Team,
SIME Lab ENSIAS, University of Mohammed V Souissi
Rabat- Morocco

Gueddah Hicham
Telecom and Embedded Systems Team, SIME Lab
ENSIAS,
University of Mohammed V Souissi
Rabat- Morocco

Yousfi Abdellah
Eradiass Team
Faculty of Juridical, Economic
and Social Sciences University of
Mohammed V-Souissi,
Rabat- Morocco

Belkasmi Mostafa
Telecom and Embedded Systems Team, SIME Lab
ENSIAS,
University of Mohammed V Souissi
Rabat- Morocco

*Abstract*—**In this paper, we propose a new approach for spell-checking errors committed in Arabic language.**

**This approach is almost independent of the used dictionary, of the fact that we introduced the concept of morphological analysis in the process of spell-checking. Hence, our new system uses a stems dictionary of reduced size rather than exploiting a large dictionary not covering the all Arabic words.**

**The obtained results are highly positive and satisfactory; this has allowed us to appreciate the validity of our concept and shows the importance of our new approach.**

*Keywords—Arabic language; Lexicon; Misspelled word; Error model; Spell-checking; Edit distance; Morphological analysis; Prefix; Stem, Suffix*

## I. INTRODUCTION

Automatic correction of spelling errors is one of the most important areas in the field of Natural Language Processing (NLP) and it has been a subject of many researches since the 60's [1]. Spell-checking consists in suggesting the closest corrections for a misspelled word, this implies the development of error models and methods allowing the scheduling of plausible corrections and the disposal of a representative lexicon for a given language in order to compare.

Early researches consist in founding a kind of error modeling. Since, appears Damerau's [2] definition which consider a spelling error as a simple combination of elementary edition operations of insertion, deletion, transposition and substitution. Based on Damerau's definition, Levenshtein [3] will define his distance (Levenshtein distance) which is characterized by three of elementary edition operations: insertion, deletion and permutation. Another modeling proposed by Pollock and Zamora [4] consists in associating for each word in the dictionary its alpha-code (consonants of the word), hence the need of having two dictionaries: one for the words and the other for their alpha-codes, and therefore the correction will be done by comparing alpha-codes with the

misspelled word. This method is efficient for permutation errors cases.

We also find among these studies: the decomposition method based on the concept of N-gram language model which is based on decomposing a misspelled word to di-trigrams and compare them to the dictionary di-trigrams in order to produce a similarity index to designate the nearest words to the misspelled word [5]. In 1996, Oflaser [6] defined a new approach called tolerant recognition of spelling errors by using the concept of finite state automaton and a distance called cut-off edit distance. Using this approach, the correction of a misspelled word is done by browsing the dictionary automaton and by calculating the cut-off edit distance for each transition, without exceeding a threshold previously defined in the algorithm. Gueddah, Yousfi and Belkasmi [7] proposed a typical and efficient variant of edit distance by integrating frequency editing errors matrices [8] in the Levenshtein algorithm in order to improve the scheduling of the solutions of an erroneous word in Arabic documents.

Generally in natural languages, and especially in Arabic, existing spelling automatic correction systems do not cover all the misspelled words. Spelling correction was always related to the disposition of a given lexicon covering the totality of misspelled words.

Several studies have been made towards the development of dictionaries adapted for spell-checking systems. Among these studies we cite in particular: the Ayaspell[1] project that aims to generate dictionaries, for example the Arabic lexicon Hunspell-ar Version 3.2 that contains more than 300000 Arabic words designed for free office suite applications of Open Office (writer) and Mozilla Firefox 3, Thunderbird and Google Chrome incorporating the spell-checker Hunspell[2] (originally designed for the Hungarian language).

---

[1] http://ayaspell.sourceforge.net
[2] http://hunspell.sourceforge.net

Despite linguistic resources available to the Arabic language, we note that we do not have yet a robust spell-checker capable of covering all the spelling errors committed. And that raised a major challenge for spell-checking [9].

In order to overcome this limitation raised on spell-checking for the Arabic language, we propose in this paper a new approach which aims to introduce the concept of morphological analysis in the process of spell-checking.

In reality, there are few works that deal with morphological analysis in spell-checking process. Among these works are cited, specially:

- Emirkanian [10] have developed an expert system in the fields of spell-checking, morphological analysis and syntactic analysis for the French language by developing a morpho-syntactic analyzer capable of detecting and correcting spelling, morphological and syntactic errors frequently committed in French documents entries. This system is based on the integration of knowledge-based rules of French for various levels: orthographic (radical's dictionary), morphological (analyzer, suffixes dictionary) and syntactic (syntactic tree, substitution rules, completion rules). This system uses a metrical distance [11] to limit the search space and define the substitution rules.

- In another approach proposed by Bowden and Kiraz [12], they have presented a morpho-graphemic model for spelling and morphological errors correction based on McCarthy morphological analyzer [13]. The advantage of this model is the way it's combines lexical analysis with morphological analysis to determine the correction possibilities.

- Another recent study is presented by Shaalan and his team [14]. This project present a spell-checking system for Arabic language that aims to explore in a first step a huge dictionary of few millions words (13 millions) generated by the AraComLex3 finite state transducer with only 9 millions of valid lexical forms filtered by the AraMorph of Buckwalter morphological analyzer [15]. These ones have explored this dictionary to propose a generic spell-checking model for Arabic by using finite state automaton technology [6] and a specific metrical distance [3] combined with a noisy channel model and also with knowledge-based rules to assign weights to the suggested corrections in order to refine the best solutions.

The major inconvenient of all the spell-checking systems resides in the limitation of the used dictionaries, because they do not cover the totality of the words of a given language. Our idea presented in this paper aims to develop a spell-checking system for Arabic language independently from the vocabulary[4] by introducing morphological analysis in the spell-checking process.

---

[3] http://aracomlex.sourceforge.net
[4] in the broadest sense of vocabulary

## II. THE MORPHOLOGICAL ANALYZER: ARAMORPH

The Buckwalter morphological analyzer [15] developed by LDC (Linguistic Data Consortium), named AraMorph, allows segmenting each word into a sequence of triplet "prefix-stem-suffix". The AraMorph analyzer is formed mainly on three lexicons: prefixes (548 entries), suffixes (906 entries) and stem (78 839 entries). Lexicons are complemented by three compatibility tables used to cover all the possible combinations of prefix-stem (2435 entries), suffix-stem (1612 entries) and prefix-suffix (1138 entries). Thus, the parser will output the stems, prefixes and suffixes associated to the word to be analyzed, and then it checks the validity of these solutions in the lexicon of the system and in the correspondence tables prefix-stem, stem-suffix and prefix-suffix. The stems used in AraMorph are constructed as follows: the stems of root "فعل" are: "فاعل", "فعول", "فعيل" , "فوعل" and "فعال".

## III. THE LEVENSHTEIN DISTANCE

Among the most known metrical methods in the field of spell-checking, we have the unavoidable Levenshtein distance [3], also known as the Edit Distance. The edit distance calculates the minimal number of elementary editing operations required to transform a misspelled word to a dictionary word. Editing operations considered by Levenshtein are: insertion, deletion, and permutation. The procedure of calculating the edit distance between two strings

$P = p_1 p_2 \ldots p_m$ and $Q = q_1 q_2 \ldots q_n$ where the length is respectively $m$ and $n$, consists in calculating recursively step by step in a matrix $O(m, n)$ the edit distance between different substrings of $P$ and $Q$.

The calculation of the cell $(i, j)$ corresponding to the edit distance between the two substrings

$P_1^i = p_1 p_2 .. p_i$ and $Q_1^i = q_1 q_2 .. q_i$, is given by the following recurrent relation:

$$D(i, j) = \text{Min} \begin{cases} D(i - 1, j) + 1, \\ D(i, j - 1) + 1, \\ D(i - 1, j - 1) + \text{cost} \end{cases} \quad (1)$$

with

$$cost = \begin{cases} 0 & if \ p_{i-1} = q_{j-1} \\ 1 & otherwise \end{cases} \quad (2)$$

Admitting these following initializations: $D(i, \emptyset) = i$ and $D(\emptyset, j) = j$, where $\emptyset$ is the empty string.

## IV. INTRODUCING MORPHOLOGICAL ANALYSIS INTO LEVENSHTEIN DISTANCE

Our new idea in this work is to use a dictionary of small size that represents Arabic language stems[5] to correct spelling errors instead of using a large dictionary. In other words, our vision is to invest in a relevant metric method instead of building a dictionary that covers all the words in a given language, which is usually difficult to build.

---

[5] Stems dictionary is the one used by Buckwalter in AraMorph Parser

We note by:

- $T = \{ T_1, T_2, \ldots T_i \}$ The set of Arabic stems.
- $P = \{ P_1, P_2, \ldots P_j \}$ All Arabic prefixes.
- $S = \{ S_1, S_2, \ldots S_k \}$ All Arabic suffixes.
- $L$ means an Arabic lexicon.

Let $W_{err}$ a misspelled word, Levenshtein distance consists in finding the words $W_s$ satisfying the following relation:

$$W_s = \underset{W_i \in L}{ArgMin}(D_{lev}(W_{err}, W_i)) \qquad (3)$$

with $D_{lev}$ presents Levenshtein distance.

According to the morphological analysis approach, there exist $(P_v, T_v, S_v)$ in $P \times T \times S$ such as $W_i = P_v T_v S_v$ respectively for the misspelled word $W_{err} = P_{err} T_{err} S_{err}$, where $P_{err}$ means an erroneous prefix and $T_{err}$ means an erroneous stem and $S_{err}$ means an erroneous suffix.

In order to introduce the morphological analysis concept (used by Buckwalter) in the Levenshtein algorithm, we have defined a new measure noted $M$, as well the measurement between $W_{err}$ and vector $(P_v, T_v, S_v)$ is given by the following formula:

$$M(W_{err}, (P_v, T_v, S_v)) =$$
$$\underset{(P_v, T_v, S_v) \in (P \times T \times S)}{ArgMin}[\underset{W_{err} = P_{err} T_{err} S_{err}}{Min}((D_{lev}(P_{err}, P_v) + D_{lev}(T_{err}, T_v) + D_{lev}(S_{err}, S_v))] \quad (4)$$

The corrections of erroneous word $W_{err}$ are given by:

$$W_s = P_s T_s S_s$$
$$= \underset{(P_v, T_v, S_v) \in P \times T \times S}{ArgMin} M(W_{err}, (P_v, T_v, S_v)) \qquad (5)$$

For all prefixes, stems and suffix respectively belonging to $P, T$ and $S$, we calculate the minimum only on prefixes, stems and suffixes that are compatibles with each other, and that by introducing the three tables of correspondence between prefix-stem, stem-suffix and prefix-suffix already used by Buckwalter.

Example**:**
Let "قغخل" a misspelled word to correct. By applying the formula (4), we get the following solutions in these first orders:

| Min Prefix | Min Stem | Min Suffix |
|---|---|---|
| $D_{lev}$ (Ø, ف)= 1 | $D_{lev}$ (دخل, قغخل)=2 | $D_{lev}$ (Ø, Ø)=0 |
| $D_{lev}$ (ق, ف)= 1 | $D_{lev}$ (غول, غخل)=1 | $D_{lev}$ (Ø, ـ)=1 |
| $D_{lev}$ (ف, قغ)=2 | $D_{lev}$ (دخل, غخل)=1 | |
| $D_{lev}$ (ف, قغخ)=3 | ….. | ….. |
| $D_{lev}$ (ف, قغخل)=4 | …… | …… |

- $M$("قغخل", (ف, "دخل"), Ø )) $=2$ (wich presents the minimal distance in all stems, prefixes and suffixes) → the system suggest the solution "فدخل" as a correction of the misspelled word "قغخل", with distance 2.
- Thus our method returns the solution $M$("قغخل", (ف, "غول"), Ø)) $=2$ → which represents the word "فغول", with distance 2.

## V. TESTS AND RESULTS

To highlight our approach, we have developed a spell-checking program[6] that allows comparing our method to the classical approach of Levenshtein.

The list of words used in this study as reference lexicon for Levenshtein approach contains more than 170000 words extracted from MySpell[7] program of Open Office Writer.

For our approach, we relied on a list of prefixes, suffixes and a list of stems built on Buckwalter approach basis, for example, besides the root "فعل" we find also the stems list of the five forms generated from this root: **"فاعل", "فعال"** and **"فوعل","فعيل","فعول".**

The rectifications suggestions proposed by Levenshtein distance are the word of minimal distance relative to the misspelled word. For our approach, we used the formula (4), explained in the previous paragraph. For our tests, we have used a corpus of 2784 misspelled words. There were three types of errors: addition, deletion and permutation. The table below shows the rate of correction by editing operations:

TABLE I. COMPARATIVE TABLE BETWEEN THE TWO METHODS

| | | Our approach | Levenshtein distance |
|---|---|---|---|
| **Editing operators** | Insertion | **85%** | **44%** |
| | Deletion | **81%** | **61%** |
| | Permutation | **86%** | **46%** |
| Average time / Erroneous word | | **0.10 ms** | **0.19 ms** |

To compare our new approach with Levenshtein's, we have used the following three indicators:

- The correction average time.
- The rate of rectified words.
- The size of each system lexicon.

- We have taken 170000 words as lexicon size for Levenshtein method. For our system, the theorical

---

[6] Developed in Java language under Eclipse platform
[7] http://myspell.sourceforge.net

lexicon size is N words, with: N= Nbre Prefixes x Nbre Suffixes x Nbre Stems ≈ 197x106 words. The real size of our system (lexicon) is much less than this number because the tables of correspondence between suffixes, prefixes and stems reduce this number. Generally, it is a ten of millions of words order.

- Despite the fact that the number of words covered by our method is about 1000 times higher than the size of the lexicon used by Levenshtein method in this study, the average time to correct a word is faster in our method 0.10 ms versus 0.19 ms in Levenshtein's.

- For numeral results regarding correction rate, it is obvious that our system rectify misspelled words 84% more correctly versus an average of 50.3% within Levenshtein distance. This difference is mainly due to the fact that our lexicon contains enough words compared to Levenshtein's. Spelling errors in our system are mainly due to the stems base insufficiency that stays incomplete and do not contains all the stems.

We clearly notice that our system is better than Levenshtein's, either at lexicon level or at runtime level or compared to the correction rate.

## VI. CONCLUSION

We can see clearly, that our system is much better that the one using classical comparison between two lexical forms via Levenshtein distance. The result we have gotten in the previous paragraph shows clearly the interest of our new approach and the facility of integration it has in an automatic spell-checking system.

### REFERENCES

[1] Kukich K.," Techniques for Automatically Correcting Words in Text ", ACM Computing Surveys, Volume 24, No.4, pp, 377-439, December 1992.

[2] Damerau F.J.," A technique for computer detection and correction of spelling errors ", Communications of the Association for Computing Machinery, 1964.

[3] Levenshtein V.," Binary codes capable of correcting deletions, insertions and reversals ", SOL Phys Dokl,pp, 707-710, 1966.

[4] Pollock J. and Zamora A., " Automatic Spelling Correction in Scientific and Scholarly Text ", Communications of the ACM, 27(4), pp, 358-368, 1984.

[5] Ukkonen E., " Approximate string matching with q-grams and maximal matches ", Theoretical Computer Science, 92, pp, 191–211, 1992.

[6] Oflazer K.," Error-tolerant Finite-state Recognition with Applications to Morphological Analysis and Spelling Correction ", Computational Linguistics Archive Volume 22 Issue 1,pp,73-89, March 1996.

[7] Gueddah H., Yousfi A. and Belkasmi M.," Introduction of the Weight Edition Errors in the Levenshtein Distance ", International Journal of Advanced Research in Artificial Intelligence, Volume 1 Issue 5,pp, 30-32, 2012.

[8] Gueddah H. and Yousfi A., " Etude Statistique sur les erreurs d'édition dans la langue Arabe", La 5éme conférence internationale sur les Technologies d'Information et de Communication pour l'Amazighe, IRCAM, Septembre 2012.

[9] Mitton R., " Ordering the suggestions of a spellchecker without using context ", Natural Language Engineering 15 (2), pp, 173-192, 2009.

[10] Emirkanian L. and Bouchard L.H.," La correction des erreurs d'orthographe d'usage dans un analyseur morphosyntaxique du français " dans langue Française N 83, Paris, Larousse, pp, 106-122, 1989.

[11] Romanycia M.H. and Pelletier J.F., "What is an heuristic? " Computational Intelligence, volume 1, pp, 47-58, 1985.

[12] Bowden T. and Kiraz G.A.," A morphographemic model for error correction in nonconcatenative strings ", Proceedings of the 33rd annual meeting on Association for Computational Linguistics, pp, 24-30, 1995.

[13] McCarthy J., "A prosodic theory of non-concatenative morphology ", Linguistic Inquiry 12(3), pp, 373-418, 1981.

[14] Shaalan K., Samih Y., Attia M., Pecina P., Genabith J.V.," Arabic Word Generation and Modelling for Spell Checking ", In the Proceedings of The 8th international conference on Language Resources and Evaluation (LREC'12), pp,719-725, May 2012.

[15] Buckwalter T., " Buckwalter Arabic Morphological Analyzer version 1.0 ", Philadelphia: Linguistic Data Consortium, Catalog No.LDC2002L49, ISBN 1-58563625760, 2002.

# On a New Competitive Measure for Oblivious Routing

Gábor Németh

Dept. of Telecommunications and Media Informatics
Budapest University of Technology and Economics
Budapest, Hungary

*Abstract*—Oblivious routing algorithms use only locally available information at network nodes to forward traffic, and as such, a plausible choice for distributed implementations. It is a natural desire to quantify the performance penalty we pay for this distributedness. Recently, Räcke has shown that for general undirected graphs the *competitive ratio* is only $O(log\ n)$, that is, the maximum congestion caused by the oblivious algorithm is within a logarithmic factor of the best possible congestion. And while the performance penalty is larger for directed networks (Azar gives a $\Omega(\sqrt{n})$ lower bound), experiments on many real-world topologies show that it usually remains under 2. These competitive measures, however, are of worst-case type, and therefore do not always give adequate characterization.

The more different combinations of demands a routing algorithm can accommodate in the network without congestion, the better. Driven by this observation, in this paper we introduce a new competitive measure, the *volumetric competitive ratio*, as the measure of all admissible demands compared to the measure of demands routed without congestion. The main result of the paper is a general lower bound on the volumetric ratio; and we also show a directed graph with $O(1)$ competitive ratio that exhibits $\Omega(n)$ volumetric ratio.

Our numerical evaluations show that the competitivity of oblivious routing in terms of the new measure quickly vanishes even in relatively small common-place topologies.

*Keywords—competitive ratio; oblivious routing; $\ell^p$ norm; $L^p$ norm; throughput polytope; feasible region; probability of congestion; hyper--spherical coordinates*

## I. INTRODUCTION

Routing algorithms are used to drive the process of forwarding traffic from source nodes to remote destination nodes through a network. Often, network links are of limited capacity and it is also the task of the routing algorithm to ensure that no network link gets seriously overloaded with excess traffic. Minimizing congestion, however, may require the global knowledge of the actual traffic demand pattern the users pose to the network, which is very difficult to ensure in a distributed setting. The class of non-adaptive routing algorithms with the property that only knowledge that is locally available at network nodes is used when making forwarding decisions is called *oblivious routing algorithms*. In fact, in oblivious routing a set of paths and corresponding splitting ratios are precomputed off-line, which are then applied in the on-line phase statically to the incoming traffic at network nodes. This scheme is easy to

implement in a distributed fashion, both in virtual-circuit-based as well as in a packet routed environment.

The performance of an oblivious routing algorithm is usually described in terms of the maximum congestion it produces, compared to the congestion produced by an optimal routing algorithm [1]. The maximization is taken over all possible demand combinations. Surprisingly, this *competitive ratio* (or oblivious performance ratio [2]) in undirected networks is only logarithmic [1], [3], [4], and even though no such appealing characterization exists for directed networks, experiment suggests that it rarely surpasses 2 [5]. This makes oblivious routing an attractive choice for implementing distributed routing algorithms.

Being a worst-case measure, the competitive ratio, however, might not always give adequate statistical representation of the performance penalty of distributed routing, as compared to optimal routing. In this paper, we study the competitive ratio arising as the extent to which an oblivious routing algorithm can route without congestion demands that otherwise could be routed in the network by a properly chosen, possibly centralized optimal routing algorithm.

Consider the following formal definition. Given a capacitated graph $G$ and an arbitrary routing algorithm, let the feasible region $D$ be the set of demands that can be routed by the algorithm without congestion in $G$. Let the throughput polytope $T$ of $G$ be the set of demands that can be routed without congestion at all (i.e., the feasible region of an optimal algorithm). Then, the quantity of *all* demands routable in the network is the volume $V(T)$, the quantity of demands routed by the routing algorithm without congestion is $V(D)$, and their, ratio, the so called *volumetric competitive ratio* $\alpha_V = \frac{V(T)}{V(D)}$ represents the fraction of routable demands the algorithm can handle successfully. Easily, $\alpha_V > 1$, and the smaller $\alpha_V$ the better. Sometimes, we use the measure $\alpha_{POC} = 1 - \frac{V(D)}{V(T)}$ as it has appealing statistical interpretation: provided that demands arrive from the set of routable demands T according to a uniform distribution, what is the chance that we encounter congestion at some link in the network. In other words, $\alpha_{POC}$ quanti-

fies the frequency that an operator can expect to find his network in a congested state, in the case when absolutely no information on traffic demands is available *a priori*. Thus, $\alpha_{POC} = 1 - \frac{1}{\alpha_V}$ is called the *probability of congestion (POC)*.

In this paper, we study the competitivity of oblivious routing in terms of the new performance measures.

### A. Related Work

Oblivious routing on hypercubes was first studied by Valiant and Brebner [6]. They give a randomized oblivious routing scheme with $O(\log n)$ competitive ratio. The first oblivious routing algorithm for generic undirected graphs is due to Räcke [1], with the remarkable property that the maximum congestion is within a $O(\log^3 n)$ factor of the lowest possible congestion, attainable by an optimal algorithm (for such an optimal adaptive algorithm, see [7]). This result was non-constructive, as the algorithm's off-line running time was exponential. The competitive ratio was subsequently improved by Harrelson, Hildrum and Rao in [3] to $O(\log^2 n \log \log n)$. They also show a polynomial algorithm for constructing the hierarchical tree decomposition that underlies Räcke's oblivious routing scheme (see also [8]). Finally, it was also Räcke who was able to improve the competitive ratio to $O(\log n)$ in [4]. This bound is asymptotically tight as there are very simple networks (e.g., two dimensional grids) on which no oblivious routing algorithm exists with sub-logarithmic competitive ratio [9], [10].

For directed graphs, unfortunately, no logarithmic congestion guarantee exists. Azar et al. give a $\Omega(\sqrt{n})$ lower bound in [2]. To cope with this intrinsic difficulty, Hajiaghayi et al. [11] present oblivious routing schemes that achieve $O(\log^2 n)$ competitive ratio provided that demands arrive randomly from a known demand distribution. Even though one might expect this assumption to improve the upper bound in undirected networks too, this is not the case: Hajiaghayi et al. give a $\Omega\left(\frac{\log n}{\log \log n}\right)$ lower bound in [12].

The first polynomial algorithm to obtain the best oblivious routing scheme for specific input graphs was introduced by Azar in [2]. A simpler, linear programming-based algorithm was given in [5]. Here, the task is to, given a directed or undirected graph as input, compute the static routing that produces the smallest competitive ratio possible on this graph. This can (and usually is) better than logarithmic. For instance, extensive numerical evaluations suggest that the competitive ratio in most real-world network topologies remains under 2 [5]. Unfortunately, this approach is not suitable to obtain the generic upper bound Räcke could obtain with the use of hierarchical tree decompositions (whose approach, in turn, does not yield optimal oblivious routing schemes for particular graphs).

Demand for more descriptive performance measures for oblivious routing has increased lately [13] [14]. The motivation is not necessarily to quantify the performance of oblivious routing algorithms, but rather to drive the optimization algorithms that compute them. For instance, one would better consider the average network load, or the sum-of-squares of the loads, as the performance measure, in contrast to the maximum load.

Following on the work of Gupta [15], Engler and Räcke in [16] give a universal treatment, able to treat the above case and many more. They define a generic aggregation function that determines how loads at individual links are converted to a congestion measure for the network, and then show a $O(\log n)$-competitive oblivious routing algorithm when the aggregation function is an $\ell^p$ norm. Their development is non-constructive, which was recently remedied by Bhaskara and Vijayaraghavan [17]. Note, however, that these performance measures are still of worst-case nature, meaning that it is the maximum of the congestion measure experienced over all possible demands that determines the outcome.

It seems that Rétvári et al. were the first ones to systematically study the geometric properties of the feasible region D and the throughput polytope T [18], [19]. They showed that, under reasonable regularity conditions, both D and T are compact, down-monotone, *K*-dimensional polyherda (*K* is the number of source-destination pairs). They also showed that no polynomial-size description exists for T even in very small networks. Thus, T is usually given implicitly, in the form of a linear program.

For computing the volumetric ratio, we need to obtain the volume of D and T. Unfortunately, except very low dimensions or special polytopes with high degree of symmetry (e.g., simplices, hyper-cubes), this is a very hard task. In particular, Elekes showed that one cannot construct a general polynomial-time algorithm for calculating the volume of *K*-dimensional bodies [20]. Therefore, randomized algorithms were proposed to break down the complexity and approximate the volume with a prescribed absolute/relative error [21]-[23]. These algorithms rely on Monte--Carlo integration and introduce random walks for sampling. The complexity of the best known randomized method is $O(K^4)$ linear program solver calls [24]. Even though linear programs can be solved in polynomial time, this still can be prohibitive in large networks. What is worse, random-walk-based sampling is another significant source of complexity [25], as a linear program needs to be solved at each step of the random walk, and we need thousands of random samples obtained in possibly thousands of steps.

### B. Our Results

In this paper, we introduce a new competitive measure of non-worst-case type to better characterize the performance of distributed oblivious routing as compared to optimal centralized routing. The measure, called the volumetric competitive ratio, quantifies the fraction of routable demands an algorithm can handle without congestion.

In the first part, we give performance bounds on oblivious routing in terms of the new competitive measure. For this, we develop a geometric model, using which for directed graphs we give a $\Omega(n)$ worst-case lower bound on the volumetric ratio. This behavior is exhibited even in cases when the standard competitive ratio is $O(1)$. Then, we obtain an universal upper bound for the volumentric ratio. At the moment, it is not known whether these bounds are tight.

In the second part of the paper, we conduct brief numerical evaluations on real-world topologies, which indicate that the

measures quickly exhibit the worst case behavior as K increases.

## II. NOTATIONS AND DEFINITION

Let $G = (V, E)$ be a connected directed or undirected graph ($n = |V|$ and $m = |E|$), with positive edge capacities $c_e : e \in E$, and $K$ source-destination pairs $(s_k, d_k) \in V \times V$. A traffic matrix (or, simply, demand) is a column $K$-vector $\theta = [\theta_k : k = 1, \ldots, K]$, where $\theta_k$ represents the request of the $k$-th source-destination pair. Let $S$ be a routing algorithm which, given some demand $\theta$, generates a flow $g_k(e)$ (i.e., a routing) for each $k = 1, \ldots, K$ on each edge $e \in E$. For our purposes, it is enough to know that the output of $S$ is an aggregate flow on each edge e: $S(\theta) = [f(e) = \sum_k g_k(e) : e \in E]$. We denote the relative flow (or load) $f(e)/c(e)$ produced by $S$ on $e$ by $L_S(e)$.

Gupta et al. introduce the notion of aggregation functions agg: $\mathbb{R}^E \to \mathbb{R}$ to aggregate the loads of individual edges into a cost measure [15]. Engler and Räcke in [16] study the case when the aggregation function is an $\ell^p$ norm $\|L_S\|_p = (\sum_{e \in E} |L_S(e)|^p)^{\frac{1}{p}}$. Then, they give oblivious routing algorithms to minimize the *competitive ratio* $\alpha_p$ defined by this $\ell^p$ norm:

$$\alpha_p = \max_{\mathbb{R}_+^K} \frac{\|L_{\mathrm{OBL}}\|_p}{\|L_{\mathrm{OPT}^p}\|_p} = \left\| \frac{\|L_{\mathrm{OBL}}\|_p}{\|L_{\mathrm{OPT}^p}\|_p} \right\|_{L^\infty(\mathbb{R}_+^K)}. \qquad (1)$$

Here, $\mathrm{OPT}^p$ is an optimal routing algorithm that for each demand $\theta$ assigns the flow $L_{\mathrm{OPT}^p}$ that minimizes the $\ell^p$ norm $\|L_{\mathrm{OPT}^p}\|_p$, and $L_{\mathrm{OBL}}$ is the flow produced for $\theta$ by the oblivious routing algorithm that minimizes (1). The second term in (1) comes from substituting max by the $L^\infty$ norm (for a good overview on $\ell^p$ and $L^p$ norms, see [26]). Easily, different settings of $p$ yield different interesting algorithms. For $p = \infty$ specifically we get the well-known "congestion-minimizing" oblivious routing algorithms [1]-[4]. In this paper, we mean by "oblivious routing" this very case (i.e., when $p = \infty$), but throughout the developments we shall often use different settings for $p$. Similarly, the term "competitive ratio" will mean (1) with choosing $p = \infty$, i.e., $\alpha_\infty$.

Our task in this paper is to seek alternatives to the above competitive measure. Our approach is mainly geometric, the main ingredients of which are as follows. Given a routing algorithm $S$, let the *feasible region* $D$ of $S$ be the set of demands $\theta$ to which $S$ orders a routing that does not violate edge capacities:

$$D = \{\theta : f(e) \le c(e) \ \forall e \in E, \text{where } [f(e)] = S(\theta)\}.$$

The throughput polytope $T$ of $G$ is the set of all routable demands, i.e., the feasible region of $\mathrm{OPT}^\infty$. Under the above assumptions, both $D$ and $T$ are $K$-dimensional, compact, convex, down-monotone polytopes [18]. Therefore both sets are measurable in terms of the standard Lebesgue measure, that is, the $K$-dimensional volumes $\mathrm{V}(D)$ and $\mathrm{V}(T)$ exist and are non-zero. Note that we call a set $X$ *down-monotone*, if $x \in X \Rightarrow \forall 0 \le y \le x : y \in X$.

We shall need some definitions from geometry to deal with polytope volumes. Let $P$ denote a convex, compact down-monotone polytope in $\mathbb{R}_+^K$ with $0 \in P$. Let $\phi_P$ denote the gauge functional of $P$, i.e., $\phi_P(x) = \min\{\eta \in \mathbb{R} : x \in \eta P\}$ for all $x \in \mathbb{R}_+^K$. Note that $\phi_P$ is a spherical function, that is, $\forall \alpha \in \mathbb{R} : \phi_P(\alpha x) = \alpha \phi_P(x)$. The reason why we invoke the gauge functional is that it is a natural geometric generalization of the competitive ratio $\alpha_\infty : \alpha_\infty = \max_{\theta \in T} \phi_D(\theta)$. Using this notation, the volume of $P$ is given by [23], [27]

$$\mathrm{V}(P) = \frac{1}{K!} \int_{\mathbb{R}_+^K} e^{-\phi_P(x)} \, d\mu(x), \qquad (2)$$

where $\mu$ is the measure function on $\mathbb{R}^K$. Converting the integration to hyper-spherical coordinates one can rewrite the above from integration on the entire positive orthant to integration on the surface of the unit $K$-ball $\partial B$:

$$\mathrm{V}(P) = \frac{1}{K} \int_{\partial B} \phi^{-K}(x) \, d\mu(x) = \frac{1}{K} \int_{\partial B} d_P^K(x) \, d\mu(x), \qquad (3)$$

where $d_P(x) = \max_{x \in \partial B}\{\eta \in \mathbb{R} : \eta x \in P\}$ denotes the distance of the boundary point of $P$ from the origin in the direction defined by the point $x \in \partial B$. For notational convenience, we shall often omit the dependence on $x$ and simply write $d\mu$ for $d\mu(x)$ and $d_P$ for $d_P(x)$.

## III. THE VOLUMETRIC COMPETITIVE RATIO

As mentioned previously, in this paper we want to characterize the fraction of routable demands a routing algorithm can handle without congestion. Consider the following definition:

*Definition 1:* Given a capacitated network $G$ with $K$ source-destination pairs and a routing algorithm $S$, let $T$ be the throughput polytope of $G$ and let $D$ be the feasible region of $S$ in $G$. Then, the *volumetric competitive ratio* of algorithm $S$ is defined as

$$\alpha_V = \frac{\mathrm{V}(T)}{\mathrm{V}(D)}, \qquad (4)$$

and the Probability of Congestion (POC) is $\alpha_{\mathrm{POC}} = 1 - \frac{\mathrm{V}(D)}{\mathrm{V}(T)} = 1 - \frac{1}{\alpha_V}$.

The reason of why we also define the POC is that it has relevant practical interpretation: $\alpha_{\mathrm{POC}}$ quantifies the chance that we find the network in a congested state, when traffic demands arrive from $T$ according to a uniform distribution.

In the next sections, we search global bounds on the above competitive measures. First, we discuss directed graphs and then we turn to generic bounds on undirected networks.

### A. A Worst-case Upper Bound in Directed Graph

First, we show that there exist directed graphs that exhibit $\Omega(n)$ volumetric competitive ratio, even though the standard competitive ratio $\alpha_\infty$ is $O(1)$.

*Theorem 1:* For any $n \ge 4$, there is a directed graph of $n$ nodes with $\alpha_\infty < 2$ and $\alpha_V > \frac{n}{6}$.

(a)



(b)

Fig. 1.   Directed graph (a) for illustrating the proof of Theorem 1, and the corresponding throughput polytope $T$ and the feasible region $D$ of the oblivious routing for $K = 2$ (b). For this case, $\alpha_\infty = 4/3$ and $\alpha_V = 3/4$.

*Proof:* Consider the directed graph in Fig. 1 for any $K \geq 2$, let all link capacities be 1, and let the source-destination pairs be $(s_k, K + 2) : s_k \in \{1, ..., K\}$. First, we construct the oblivious routing function w.r.t. to the $\ell^p$ norm and calculate $\alpha_\infty(K)$ as a function of $K$, then we obtain an approximation  on $\alpha_V$.

*1. Competitive ratio:* Construct the oblivious routing function as   follows.  The first source-destination pair has only a single path, thus   all its traffic is sent through this path.  The rest of the users $k \in \{2, ..., K\}$ have two paths. Due to the symmetry of the network, it   is enough to consider a general user $k$. Let $\beta$ denote the fraction   of traffic sent at $s_k$ to the path passing through node $K + 1$. We   consider two critical classes of traffic matrices, as these will produce   the largest link load for oblivious routing. One is $[1,1,...,1]^T \in T$, for which the maximum load $1 + \sum_{k=2}^{K} \beta$ occurs on link $(K + 1, K+2)$.  Second, for some general $k \in 2,...,K$ the demand $[0,0,...,2,...,0]^T \in T$, where   there is a single non-zero element in the $k$-th position, causes $2(1 - \beta)$ maximum load on link $(k, K + 2)$. To find $\alpha_\infty(K)$, we seek for $\beta$ so that the maximum load is minimal. This occurs exactly when $1 + \sum_{k=2}^{K} \beta = 2(1 - \beta)$, which yields $\alpha_\infty(K) = \frac{2K}{K+1} \leq 2$, for every finite $K$.

*2. The volumetric ratio:* We need the volume of $T$ and $D$. Instead of calculating the volumes directly, we take lower and upper  approximations. Let $T'$ and $D'$ be such that $T' \subseteq T$ and $D' \supseteq D$. By down-monotonity of $T$, the $K$-hypercube of size 1 resides completely inside $T$.

Moreover, we find $K - 1$ half-hypercubes in $T$ as well, where the $k$-th half-hypercube, $k \in \{2, ..., K\}$, is obtained by placing a $K$-hypercube of size 1 at the point $[0, ..., 1, ..., 0]^T$ (1 is in the $k$-th position) and taking the intersection with the half-

space $\theta_1 + \theta_k \leq 2$. The volume of $T'$ is the sum of the volumes of the above  polytopes: $V(T') = 1 + \frac{1}{2}(K - 1)$.  Second, we give an outer approximation $D'$ for $D$. Simply put, $D$ is enclosed by a hyper-rectangle, whose lower left corner is the origin and whose upper right corner is the point $\left[1, 1 + \frac{1}{K}, ..., 1 + \frac{1}{K}\right]^T$. Let this hyper-rectangle be $D'$ and so $V(D') = \left(1 + \frac{1}{K}\right)^{K-1}$.

Putting all together, we get the desired $\alpha_V = \frac{V(T)}{V(D)} \geq \frac{V(T')}{V(D')} = \frac{1 + \frac{1}{2}(K-1)}{\left(1 + \frac{1}{K}\right)^{K-1}} \geq \frac{n}{6}$ with the substitution $n = K + 2$. ∎

### B. A Worst-case Lower Bound for Undirected Graphs

The result of the previous section is only valid for a class of special directed networks, in which obtaining good approximations on $V(D)$ and $V(T)$ is possible.  For general graphs, this approach is not viable. Therefore, we shall pursue a different approach below: we give a general lower bound on $\alpha_V$ in terms of $\alpha_\infty$.  Consider the following theorem.

*Theorem 2:* $\alpha_V \geq \frac{V(T)}{\sqrt{V(T')}} (\alpha_\infty)^{\frac{K}{2}}$,

where $T' = \{y \in \mathbb{R}_+^K : \exists x \in T, y = x^2\}$.

First, consider the following technical Lemma.

*Lemma 1:* Let $f \leq g$ and $g^{2K}$, $\left(\frac{f}{g}\right)^{2K}$ belong to $L^2(\partial B_r)$ and let $\mu(\partial B_r) = 1$. Then

$$\frac{\int_{\partial B_r} f^K d\mu}{\int_{\partial B_r} g^K d\mu} \leq \frac{\sqrt{\int_{\partial B_r} g^{2K} d\mu}}{\int_{\partial B_r} g^K d\mu} \sqrt{\int_{\partial B_r} \left[\frac{f}{g}\right]^{2K} d\mu}.$$

Lemma 1 is the direct consequence of the Cauchy—Schwarz—Bunyakovsky inequality [26], [28], [29].

*Proof of Theorem 1:* Using (3) and the above Lemma, we get the following upper bound on $\frac{1}{\alpha_V}$.

$$\frac{1}{\alpha_V} = \frac{V(D)}{V(T)} = \frac{\int_{\partial B} d_D^K d\mu}{\int_{\partial B} d_T^K d\mu} = \frac{\int_{\partial B_r} d_D^K d\mu}{\int_{\partial B_r} d_T^K d\mu}$$

$$\leq \frac{\sqrt{\int_{\partial B_r} d_T^{2K} d\mu}}{\int_{\partial B_r} d_T^K d\mu} \sqrt{\int_{\partial B_r} \left[\frac{d_D}{d_T}\right]^{2K} d\mu} \qquad (5)$$

$$= \frac{\sqrt{V(T')}}{V(T)} \sqrt{\left\|\left[\frac{\|L_{OPT^\infty}\|_\infty}{\|L_{OBL}\|_\infty}\right]^K\right\|_{L^{2K}(\partial B_r)}}$$

where the radius $r : \mu(\partial B_r) = 1$.  The last equation comes from observing that the term $\frac{d_D}{d_T}$ is the value the feasible region should be scaled to enclose all the points of the throughput polytope in the selected direction. Thus,

$$\frac{d_T}{d_D} = \frac{\|L_{OBL}\|_\infty}{\|L_{OPT^\infty}\|_\infty}. \qquad (6)$$

(a) NSF (dir)
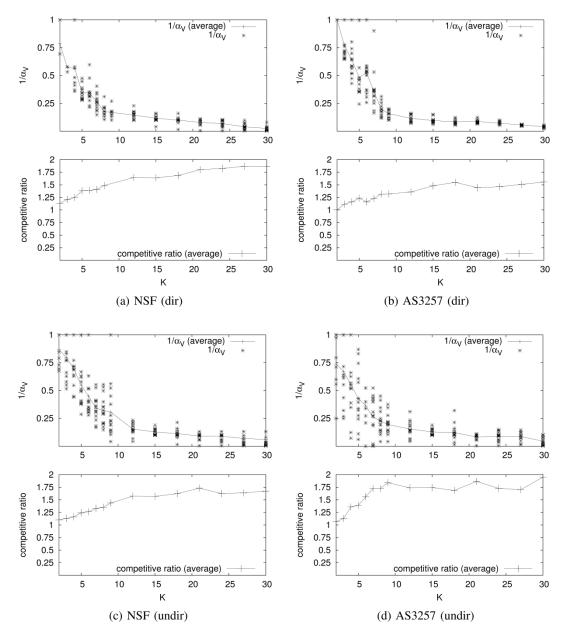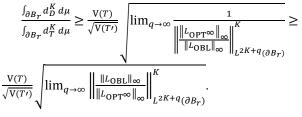
(b) AS3257 (dir)

(c) NSF (undir)

(d) AS3257 (undir)

Fig. 2. The competitive ratio $\alpha_\infty$ and the inverse of the approximated volumetric ratio $\alpha_V$ w.r.t. the number of source-destination pairs $K$ for selected directed and undirected networks.

We rewrite (5) so that the outer norm is a $L^\infty$ norm. First, we use fact that $\lim_{q\to\infty}\|f\|_{L^q(\partial B_r)} = \|f\|_{L^\infty(\partial B_r)}$. Second, choosing let $r$ so that let $\mu(\partial B_r) = 1$ yields $\forall q \in \mathbb{N}: \|f\|_{L^{q+1}(\partial B_r)} \geq \|f\|_{L^q(\partial B_r)}$. Thus, we write (5) as:

$$\alpha_V = \frac{\int_{\partial B_r} d_D^K\, d\mu}{\int_{\partial B_r} d_T^K\, d\mu} \geq \frac{V(T)}{\sqrt{V(T')}}\left(\left\|\left\|\frac{\|L_{\text{OPT}^\infty}\|_\infty}{\|L_{\text{OBL}}\|_\infty}\right\|\right\|_{L^{2K}(\partial B_r)}^K\right)^{-1/2} \geq$$

$$\frac{V(T)}{\sqrt{V(T')}}\left(\lim_{q\to\infty}\left\|\left\|\frac{\|L_{\text{OPT}^\infty}\|_\infty}{\|L_{\text{OBL}}\|_\infty}\right\|\right\|_{L^{2K+q}(\partial B_r)}^K\right)^{-1/2}.$$

Using Jensen's inequality [26] for the concave function $\frac{1}{id_\mathbb{R}}$ yields:

$$\frac{\int_{\partial B_r} d_D^K\, d\mu}{\int_{\partial B_r} d_T^K\, d\mu} \geq \frac{V(T)}{\sqrt{V(T')}}\sqrt{\lim_{q\to\infty}\frac{1}{\left\|\left\|\frac{\|L_{\text{OPT}^\infty}\|_\infty}{\|L_{\text{OBL}}\|_\infty}\right\|\right\|_{L^{2K+q}(\partial B_r)}^K}} \geq$$

$$\frac{V(T)}{\sqrt{V(T')}}\sqrt{\lim_{q\to\infty}\left\|\left\|\frac{\|L_{\text{OBL}}\|_\infty}{\|L_{\text{OPT}^\infty}\|_\infty}\right\|\right\|_{L^{2K+q}(\partial B_r)}^K}.$$

Therefore

$$\frac{\int_{\partial B_r} d_D^K \, d\mu}{\int_{\partial B_r} d_T^K \, d\mu} \geq \frac{V(T)}{\sqrt{V(T')}} \sqrt{\lim_{q\to\infty} \left\| \frac{\|L_{\text{OBL}}\|_\infty}{\|L_{\text{OPT}}\infty\|_\infty} \right\|_{L^\infty(\partial B_r)}^K},$$

which completes the proof.

## IV. NUMERICAL EVALUATIONS

Finally, we seek the answer for the question whether the worst case bounds on the volumetric ratio presented in Section III indeed appear in real networks. Therefore, we conducted some numerical evaluations using the volume approximation algorithm of the previous section. We ran the evaluations on the ISP data maps from the Rocketfuel dataset [30]. We used the same method as in [5] to obtain approximate POP-level topologies: we collapsed the topologies so that nodes correspond to cities, we eliminated leaf-nodes and we set link capacities inversely proportional to the link weights. In this paper, only results for the network AS3257 are shown, as we observed similar results for the rest of the networks as well. Another round of evaluations was conducted on the NSFNET Phase II topology [31]. From these topologies, we generated two series of increasingly more complex networks by adding gradually more source-destination pairs[1]. Recall that the number of source-destination pairs $K$ determines the dimension of the underlying geometric space, and so it has profound impact on $\alpha_V$. In our experiments, $K$ was increased from 2 to 30. For each $K$, fifteen independent samples were generated picking the source and the destination nodes randomly according to a bimodal distribution and then $\alpha_\infty$ and $\alpha_V$ were evaluated for each scenario. The parameters were chosen so that the result has larger than 10% relative error with less than 5% probability. Fig. 2 depicts the results for both networks. Note that we show $\frac{1}{\alpha_V}$ instead of $\alpha_V$ for better visualization.

Our most important observation is that real networks indeed exhibit the worst-case behavior seen in Section III. We observe that with the increase of $K$, $\frac{1}{\alpha_V}$ rapidly approaches 0. This suggests that already in networks with more than a couple of source-destination pairs it is only a very small fraction of all the routable demands the oblivious routing algorithm can handle without congestion. In terms of the probability of congestion, which approaches 1 as $\frac{1}{\alpha_V}$ approaches 0, this basically means that a network adopting oblivious routing will spend most of its time in a congested state, provided that demands arrive uniformly from the set an optimal algorithm could route without any congestion at all. And this is despite of the fact that $\alpha_\infty$ remains low (we see $\alpha_\infty < 2$ in both networks for all $K$).

## V. COCLUSION

Oblivious routing is a promising candidate for minimum-congestion routing in large networks. This is thanks to that, on the one hand, it is a fundamentally distributed scheme and, on the other hand, it comes equipped with a hard performance guarantee, namely, the maximum congestion it causes is within a logarithmic factor of the best possible congestion. This performance characterization, however, is intrinsically of worst-case nature.

In this paper, we introduced an alternative competitive measure, the so called volumetric ratio, which measures the fraction of routable demands an oblivious routing algorithm can route without congestion. We observed that already in very small directed networks (i.e., the ones in Section III-A), oblivious routing algorithms order infeasible routing to $O\left(\frac{1}{n}\right)$ fraction of the, otherwise routable, demands. We showed further worst-case bounds valid for both directed and undirected graphs.

A disadvantage of the new measure is that it is very difficult to numerically evaluate it. In fact, negative results on exact polytope volume computation suggest that we cannot hope for a polynomial time algorithm to compute $\alpha_V$.

Finally, by numerical evaluations we showed that the worst-case behavior we identified clearly manifests itself in real networks.

Easily, this paper is only a first step towards a more thorough performance characterization of oblivious routing. At the moment, it is unclear whether our bounds are tight, and we seriously lack proper upper bounds. It seems though that, with some more work, the geometric model we introduced in this paper will be able to provide these results.

## REFERENCES

[1] H. Räcke, "Minimizing congestion in general networks," in Proceedings of the 43rd Symposium on Foundations of Computer Science, FOCS '02, pp. 43–52, 2002.

[2] Y. Azar, E. Cohen, A. Fiat, H. Kaplan, and H. Racke, "Optimal oblivious routing in polynomial time," in Proceedings of the thirty-fifth annual ACM symposium on Theory of computing, STOC '03, (New York, NY, USA), pp. 383–388, ACM, 2003.

[3] C. Harrelson, K. Hildrum, and S. Rao, "A polynomial-time tree decomposition to minimize congestion," in Proceedings of the fifteenth annual ACM symposium on Parallel algorithms and architectures, SPAA '03, (New York, NY, USA), pp. 34–43, ACM, 2003.

[4] H. Räcke, "Optimal hierarchical decompositions for congestion minimization in networks," in Proceedings of the 40th annual ACM symposium on Theory of computing, STOC '08, pp. 255–264, 2008.

[5] D. Applegate and E. Cohen, "Making intra-domain routing robust to changing and uncertain traffic demands: understanding fundamental tradeoffs," in ACM SIGCOMM, pp. 313–324, 2003.

[6] L. G. Valiant and G. J. Brebner, "Universal schemes for parallel communication," in Proceedings of the thirteenth annual ACM symposium on Theory of computing, STOC '81, pp. 263–277, 1981.

[7] G. Rétvári and G. Németh, "On optimal multipath rate-adaptive routing," in IEEE Symposium on Computers and Communications, ISCC'10, pp. 605–610, IEEE Computer Society, 2010.

[8] M. Bienkowski, M. Korzeniowski, and H. Räcke, "A practical algorithm for constructing oblivious routing schemes," in Proceedings of the fifteenth annual ACM symposium on Parallel algorithms and architectures, SPAA '03, (New York, NY, USA), pp. 24–33, ACM, 2003.

[9] Y. Bartal and S. Leonardi, "On-line routing in all-optical networks," in Proceedings of the 24th International Colloquium on Automata, Languages and Programming, ICALP '97, pp. 516–526, 1997.

[10] B. Maggs, F. Meyer auf der Heide, B. Vocking, and M. Westermann, "Exploiting locality for data management in systems of limited bandwidth," in Foundations of Computer Science, 1997. Proceedings., 38th Annual Symposium on, pp. 284 –293, oct 1997.

[11] M. Hajiaghayi, J. Kim, T. Leighton, and H. Räcke, "Oblivious routing in directed graphs with random demands," in STOC '05, pp. 193–201, 2005.

[12] M. T. Hajiaghayi, R. D. Kleinberg, T. Leighton, and H. Räcke, "New lower bounds for oblivious routing in undirected graphs," in Proceed-

---

[1] A preliminary version of the results were presented as a poster at ACM Sigmetrics 2012 [32].

ings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm, SODA '06, (New York, NY, USA), pp. 918–927, ACM,2006.

[13] P. Harsha, T. Hayes, H. Narayanan, H. Räcke, and J. Radhakrishnan, "Minimizing average latency in oblivious routing," in SODA '08, (Philadelphia, PA, USA), pp. 200–207, 2008.

[14] G. Lawler and H. Narayanan, "Mixing times and $l_p$ bounds for oblivious routing," in Workshop on Analytic Algorithmics and Combinatorics, (ANALCO'09) 4, 2009.

[15] A. G. M. T. Hajiaghayi and H. Räcke, "Oblivious network design," in In Proceedings of the 17th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 970–979, 2006.

[16] M. Englert and H. Räcke, "Oblivious Routing for the $L_p$-norm," IEEE Foundations of Computer Science, pp. 32–40, 2009.

[17] A. Bhaskara and A. Vijayaraghavan, "Computing the matrix p-norm," CoRR, vol. abs/1001.2613, 2010.

[18] G. Rétvári, J. J. Bíró, and T. Cinkler, "Fairness in capacitated networks: A polyhedral approach," in IEEE INFOCOM, vol. 1, pp. 1604–1612, May 2007. G. Rétvári and G. Németh, "Demand-oblivious routing: distributed vs. centralized approaches," in Proceedings of the 29th conference on Information communications, INFOCOM'10, (Piscataway, NJ, USA), pp. 1217–1225, IEEE Press, 2010.

[19] G. Elekes, "A geometric inequality and the complexity of computing volume," Discrete and Computational Geometry, vol. 1, pp. 289–292, 1986.

[20] M. Dyer, A. Frieze, and R. Kannan, "A random polynomial-time algorithm for approximating the volume of convex bodies," J. ACM, vol. 38, pp. 1–17, January 1991.

[21] M. Dyer and A. Frieze, "Computing the Volume of Convex Bodies: A Case where Randomness Provably Helps," in Proceedings of Symposia in Applied Mathematics, vol. 44, 1991.

[22] L. Lovász and M. Simonovits, "Random walks in a convex body and an improved volume algorithm," Random Structures & Algorithms, vol. 4, no. 4, pp. 359–412, 1993.

[23] L. Lovász and S. Vempala, "Simulated annealing in convex bodies and an O*($n^4$) volume algorithm," J. Comput. Syst. Sci., vol. 72, pp. 392–417, March 2006.

[24] R. Kannan, L. Lovász, and M. Simonovits, "Random walks and an O*($n^5$) volume algorithm for convex bodies," Random Structures & Algorithms, vol. 11, no. 1, pp. 1–50, 1997.

[25] G. B. Folland, Real Analysis: Modern Techniques and Their Applications. Pure and Applied Mathematics, John Wiley & Sons, Inc., 2 ed., 1999.

[26] P. Gritzmann and V. Klee, "On the Complexity of Some Basic Problems in Computational Convexity: II. Volume and mixed volumes," in Polytopes: Abstract, Convex and Computational (T. Bisztriczky, P. McMullen, R. Schneider, and A. W. Weiss, eds.), vol. 440 of NATO Adv. Sci. Inst. Ser. C Math. Phys. Sci., pp. 373–466, Dordrecht: Kluwer Acad. Publ., 1994.

[27] H. Royden and P. Fritzpatrick, Real Analysis. Prentice Hall, 4 ed., 2007.

[28] W. Rudin, Real and Complex Analysis. International Series in Pure and Applied Mathematics, McGraw-Hill, 3 ed., 1986.

[29] R. Mahajan, N. Spring, D. Wetherall, and T. Anderson, "Inferring link weights using end-to-end measurements," in ACM IMC, pp. 231–236, 2002.

[30] B. Chinoy and H. W. Braun, "The National Science Foundation network." Tech. Rep., CAIDA, available online: http://www.caida.org/outreach/papers/1992/nsfn/nsfnet-t1-technology.pdf, Sep 1992.

[31] G. Németh and G. Rétvári, "Towards a statistical characterization of the competitiveness of oblivious routing (poster)." ACM Sigmetrics, 2012.

# Design of Reversible Counter

Md. Selim Al Mamun

Dept. of Computer Science and Engineering
Jatiya Kabi Kazi Nazrul Islam University
Trishal, Mymensingh-2220, Bangladesh

B. K. Karmaker

Dept.of Electronics and Communication Engineering
Jatiya Kabi Kazi Nazrul Islam University
Trishal, Mymensingh-2220, Bangladesh

*Abstract*—**This article presents a research work on the design and synthesis of sequential circuits and flip-flops that are available in digital arena; and describes a new synthesis design of reversible counter that is optimized in terms of quantum cost, delay and garbage outputs compared to the existing designs. We proposed a new model of reversible T flip-flop in designing reversible counter.**

*Keywords—Flip-flop; Counter; Garbage Output; Reversible Logic; Quantum Cost*

## I. INTRODUCTION

R. Landauer [1] states that traditional logic operations dissipate heat due to the loss of information bits. It is proved that each bit of information loss generates kTln2 joules of heat energy; where k is Boltzmann's constant and T is the absolute temperature at which computation is performed. C. H. Bennett [2] showed that energy dissipation problem can be avoided if all the gates in the circuits are reversible. This is because reversible logic makes every step of computation to be completely reversible, so that no information is lost at any step of computation.

Research is going on reversible logic and a good amount of research work has been carried out in the area of reversible combinational logic. However, there is not much work in the area of sequential circuit like flip-flops and counters. A counter is a sequential circuit capable of counting the number of clock pulses that have arrived at its clock input. This paper proposes a novel of n bit reversible counter. The efficiency of the proposed design is proved with the help of proper theorems and algorithms.

The rest of the paper is organized as follows: Section 2 presents background on reversible logic. Section 3 describes related works on reversible counter. Section 4 describes the logic synthesis of our proposed reversible counter design and comparisons with other research works. Finally this paper is concluded with the Section 5.

## II. BACKGROUND ON REVERSIBLE LOGIC

This section focuses on the cost metrics used in this paper and describes some popular reversible gates along with their quantum representations.

### A. Cost Metrics

A reversible circuit can be synthesized in several ways, resulting different cost. This section outlines four cost metrics which are generally used to evaluate and compare reversible circuits.

*1) Gate Count:* This refers to the number of gates required to implement the circuit. This is used as a major cost metric in the evaluation of reversible sequential circuit [3]. But gate count is not a good metric for comparison as reversible gates are of different type and have different quantum costs [4].

*2) Garbage Output:* Some outputs are used only to maintain the reversibility of the circuit, but not result the final outputs nor are they used as input to other circuits. These unused outputs are known as garbage outputs.
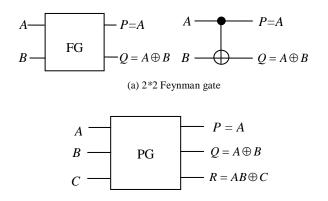
*3) Delay:* The maximum number of gates in a path from any input line to any output line is considered as the delay of the circuit [5]. This paper used the logical depth as the measure of delay for reversible circuit proposed by Mohammadi and Eshghi [6].

*4) Quantum Cost:* The quantum cost of a reversible gate is the number of quantum gates or 1x1 and 2x2 reversible gates required to present the gate. The quantum costs of all reversible 1x1 and 2x2 gates are taken as unity [7].

This paper uses quantum cost, delay and the number of garbage bits as the cost metrics while comparing the proposed design with the existing results.

### B. Quantum Analysis of Popular Reversible Gates

Several reversible logic gates have been designed till now. Some popular reversible gates and their quantum equivalent diagrams are shown in Fig.1. Feynman gate (FG) [8] is the only 2*2 gate which has quantum cost 1. Among 3*3 gate Peres gate (PG) [9] has quantum cost 4, Selim Al Mamun (SAM) [10] gate has quantum cost 4 and Toffoli gate (TG) [11] has quantum cost 5.
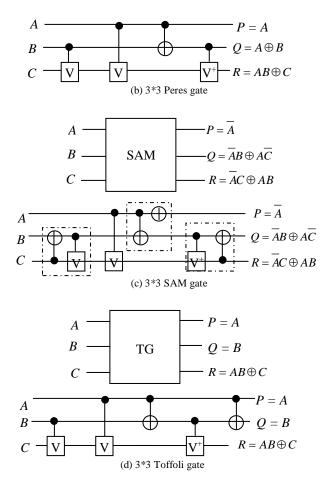


(a) 2*2 Feynman gate

(b) 3*3 Peres gate



(c) 3*3 SAM gate



(d) 3*3 Toffoli gate

Fig. 1.   Quantum analysis of popular reversible gates

### III.   RELATED WORKS ON REVERSIBLE COUNTER

Researchers have worked on many ways on sequential circuit and work is still going on. This section reviews some previous implementation and sequential circuit designs. Researchers [10, 12-15] proposed the implementation of all types of latches, flip-flops and their master-slave design. These works opens a door to the implementation of large sequential circuit like counter.
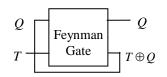
The authors [16-17] carry the above success to the of design reversible counter. Some of the works are implemented by the replacement of latches and gates by their reversible counter parts. Recently Khan [16] used Positive Polarity Reed Muller (PPRM) expression to design synchronous counter. All these works suggest that there is scope for the design and implementation of large sequential circuits like counter.

### IV.   DESIGN AND SYNTHESIS OF REVERSIBLE COUNTER

This section describes our proposed design for n bit counter. Designs for both the asynchronous counter andthe synchronous counter are presented here. While designing counter, this paper also proposed the design of reversible T flip-flop which is the building block of the counter. This paper presents the design of T flip-flop, gated T flip-flop and master slave T flip-flop.

### A.   Proposed T Flip-flop

The characteristic equation of a T flip-flop is $T\overline{Q} \oplus Q\overline{T}$ $= T \oplus Q$. A T Flip-flop can be realized by a single Feynman gate. Our proposed T flip-flop is shown in Fig.2.



Fig. 2.   Design of T Flip-Flop.

Our proposed T flip-flop with Q output has only one gate, quantum cost 1, delay 1 and no garbage outputs.

### B.   Design of clocked T Flip-flop

The characteristic equation of a clocked T flip-flop is $Q = (T \oplus Q).CLK \oplus \overline{CLK}.Q$ . The equation can be simplified as $Q = (T.CLK) \oplus Q$ . This clocked flip-flop is realized by a Peres gate and a Feynman gate. Two designs are proposed here. Fig.3 shows the design of a clocked T flip-flop. Fig.3 (a) is used for synchronous counter and Fig.3 (b) is used for asynchronous counter.



Fig. 3.   Deign of clocked T flip-flop.

Both of our proposed designs have quantum cost 5, delay 5 and garbage output 1. Comparisons of the resources of clocked T flip-flop of our design with the existing design are given in Table I.

TABLE I.      COMPARISONS OF DIFFERENT TYPES OF CLOCKED T FLIP-FLOPS

| Gated T flip-flop design | Cost Comparison | | |
|---|---|---|---|
| | *Quantum Cost* | *Delay* | *Garbage Outputs* |
| Proposed | 5 | 5 | 1 |
| Chuang[12] | 6 | 6 | 2 |
| Thapliyal[14] | 6 | 6 | 2 |

## C. Master Slave T flip-flop

To implement master slave T flip-flop, it needs one flip-flop working as master and another is slave. The same strategy is followed here. For master flip-flop, Peres gate is modified. The input vector $I_V$ and output vector $O_V$ of a 3*3 modified Peres gate, MPG are defined as follows, $I_v = (A, B, C)$ and $O_v =( P = \overline{A}, Q = A \oplus B, R = AB \oplus C)$. The quantum cost of MPG is 4. The block diagram and equivalent quantum representation for 3*3 MPG are shown in Fig.4.



(a)



(b)

Fig. 4.    (a) Block diagram of 3*3 MPG gate and (b) Equivalent quantum representation.

The modification is required to produce negative clocked pulse without generating any additional gate cost. The proposed master slave T flip-flop is shown in Fig.5.
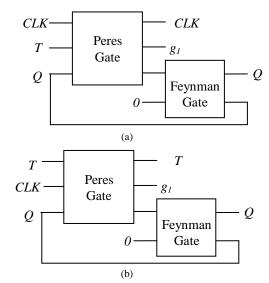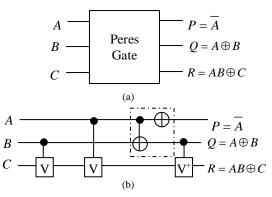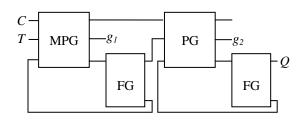


Fig. 5.    Design of master slave T flip-flop.

Our proposed design has quantum cost 10, delay 10 and garbage output 2. Comparisons of resources of master slave T flip-flop of our design with the existing design are given in Table II.

TABLE II.         COMPARISONS OF DIFFERENT TYPES OF MASTER SLAVE T FLIP-FLOPS

| Master-Slave T flip-flop | Cost Comparison | | |
|---|---|---|---|
| | *Quantum Cost* | *Delay* | *Garbage Outputs* |
| Proposed | 10 | 10 | 2 |
| Thapliyal[14] | 11 | 11 | 3 |
| Thapliyal [18] | 17 | 17 | 4 |

## D. Design of Asynchronous Counter

In asynchronous counter, the T flip-flops are arranged in such a way that output of one flip-flop is connected to the clock input of the next higher order flip-flop. The output of a flip-flop triggers the next flip-flop. The flip-flop holding the least significant bit receives the incoming count pulse. Our proposed 4 bit asynchronous counter is shown in Fig.6. The counter is realized by 4 Peres gates and some Feynman gates.
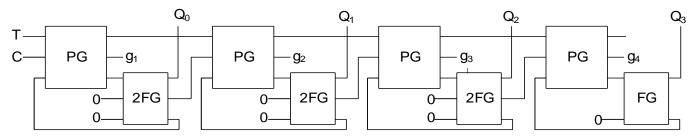


Fig. 6.    Design of 4bit asynchronous counter

Our proposed 4bit asynchronous counter has quantum cost 23, delay 23 and garbage output 4. There is not much good result available about asynchronous counter. Comparisons of our proposed design with existing result are shown in Table III.

TABLE III.         COMPARISONS OF PROPOSED DESIGN OF 4 BIT ASYNCHRONOUS COUNTER WITH EXISTING DESIGN

| 4bit asynchronous counter | Cost Comparison | | |
|---|---|---|---|
| | *Quantum Cost* | *Delay* | *Garbage Outputs* |
| Proposed | 23 | 23 | 4 |
| Rajmohan[17] | 55 | 55 | 12 |

**Theorem 1:** To construct n bit asynchronous counter, if g is the total number of gates required to design the counter producing b number of garbage outputs then g≥2n and b≥n.

**Proof:** Each flip-flop consists of two gates; n bit counter requires n number of flip-flops. No additional gates are required to interconnect each other. So total number of gates required to design the counter is 2n, hence g≥2n. Similarly, every flip-flop produces only one garbage output. No garbage output produced while interconnection among flip-flops. So the total number of garbage out is n, hence b≥n.

**Theorem 2:** The quantum cost of an n bit asynchronous counter is $Q_n \geq 6n-1$.

**Proof:** For n=1, only one Peres gate and one Feynman gate is required to construct the counter. The quantum cost of Peres gate is 4 and quantum cost of Feynman gate is 1. So the total cost 4+1=5.
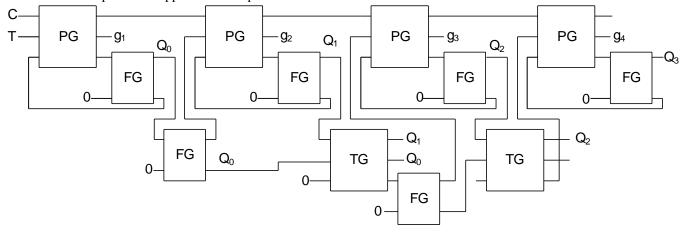
Now for n>1, one Peres gate and one double Feynman gate is required for each flip-flop in the counter except the last one

which requires one Feynman gate instead of double Feynman gate. So for n bit asynchronous counter it needs n Peres gate, (n-1) double Feynman gate and one Feynman gate. The quantum cost of double Feynman gate is 2. So the total quantum cost is 4*n+2(n-1)+1 =6n-1, Hence $Q_n \geq 6n-1$.

*E. Design of Synchronous Counter*

Synchronous counter is different from asynchronous counter in that clock pulses are applied to the inputs of all the

flip-flops at a time. A flip-flop is complemented depending on the input value T and the clock pulse. The flip-flop in least significant position is completed with every clock pulse. A flip-flop in other position is complemented only when all the outputs of preceding flip-flops produces 1. Same strategy is followed here to implement the synchronous counter. Fig.7 shows our proposed 4bit synchronous counter.



Fig. 7. Design of 4bit synchronous counter

Our proposed 4bit synchronous counter has quantum cost 32, delay 32 and garbage output 4. . Comparisons of our proposed design with existing result are shown in Table IV.

TABLE IV. COMPARISONS OF PROPOSED DESIGN OF 4 BIT SYNCHRONOUS COUNTER WITH EXISTING DESIGN

| 4bit synchronous counter | Cost Comparison | | |
|---|---|---|---|
| | *Quantum Cost* | *Delay* | *Garbage Outputs* |
| Proposed | 32 | 32 | 4 |
| khan[16] | 35 | 35 | 4 |

**Theorem 3:** To construct n (≥3) bit synchronous counter, if g is the total number of gates required to design the counter producing b number of garbage outputs then g≥4n-4 and b≥n.

**Proof:** Each flip-flop consists of two gates; n bit counter requires n number of flip-flops. For n=3, one Toffoli and one Feynman is required to carry out all the outputs to the next higher positioned flip-flop. So total number of gates required is 3*2+2=8.

For n>3, 2n number of gates required for the flip-flops and 2(n-2) number of gates are required to carry out all the lower outputs to the next higher outputs. So the total number of gates required is 2n+2(n-2)=4n-4. Every flip-flop produces only one garbage output. No garbage output produced while interconnection among flip-flops and to carry out outputs to next higher flip-flop. So the total number of garbage out is n, hence b≥n.

**Theorem 4:** The quantum cost of an n(≥3) bit synchronous counter is $Q_n \geq 11n-12$.

**Proof:** For n (≥3) bit synchronous counter it requires n flip-flops. Each flip-flop consists of one Peres gate and one Feynman gate. Additional (n-2) Toffoli gates and (n-2)

Feynman gates are required to carry out all the outputs to the next higher flip-flop. So it requires n number of Peres gate, (n-2) number of Toffoli gates and n+(n-2) =2n-2 number of Feynman gates. Quantum cost of Peres gate is 4, Quantum cost of Toffoli gate is 5 and quantum cost of Feynman gate is 1. So total quantum cost = 4*n + 5*(n-2)+1*(2n-2)=11n-12, hence $Q_n \geq 11n-12$.

## V. CONCLUSION

Reducing quantum cost in sequential circuit is always a challenging one. Only a few attempts were made on reversible counter. A novel reversible design for both n bit synchronous and asynchronous counter is proposed. Appropriate algorithms and theorems are presented to clarify the proposed design and to establish its efficiency. As compared to the best reported designs in literature, the proposed designs are better in terms of quantum cost, delay and garbage outputs. The proposed design can have great impact in reversible computing.

### REFERENCES

[1] Rolf Landauer, "Irreversibility and Heat Generation in the Computing Process", IBM Journal of Research and Development, vol. 5, pp. 183-191, 1961.

[2] Charles H.Bennett, "Logical Reversibility of computation", IBM Journal of Research and Development, vol. 17, no. 6, pp. 525-532, 1973.

[3] Perkowski, M., A.Al-Rabadi, P. Kerntopf, A. Buller, M. Chrzanowska-Jeske, A. Mishchenko, M. Azad Khan, A. Coppola, S. Yanushkevich, V. Shmerko and L. Jozwiak, "A general decomposition for reversible logic", Proc. RM'2001, Starkville, pp: 119-138, 2001

[4] J.E Rice, "A New Look at Reversible Memory Elements", Proceedings International Symposium on Circuits and Systems(ISCAS) 2006, Kos, Greece, May 21-24 ,2006, pp. 243-246.

[5] Dmitri Maslov and D. Michael Miller, "Comparison of the cost metrics for reversible and quantum logic synthesis", http://arxiv.org/abs/quant-ph/0511008, 2006

[6] Mohammadi,M. and Mshghi,M, On figures ofmerit in reversible and quantumlogic designs, Quantum Inform. Process. 8, 4, 297–318, 2009.

[7] Md. SelimAl Mamun and Syed Monowar Hossain. "Design of Reversible Random Access Memory." *International Journal of Computer Applications* 56.15 (2012): 18-23.

[8] Richard P.Feynman, "Quantum mechanical computers," Foundations of Physics, vol. 16, no. 6, pp. 507-531, 1986.

[9] A. Peres, "Reversible Logic and Quantum Computers," Physical Review A, vol. 32, pp. 3266-3276, 1985.

[10] Md. Selim Al Mamun and David Menville, Quantum Cost Optimization for Reversible Sequential Circuit, (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 4, No. 12, 2013.

[11] Tommaso Toffoli, "Reversible Computing," Automata, Languages and Programming, 7th Colloquium of Lecture Notes in Computer Science, vol. 85, pp. 632-644, 1980.

[12] M.-L. Chuang and C.-Y. Wang, "Synthesis of reversible sequential elements," ACM journal of Engineering Technologies in Computing Systems (JETC). Vol. 3, No.4, 1–19, 2008.

[13] J. E. Rice, An introduction to reversible latches. The Computer journal,Vol. 51, No.6, 700–709. 2008.

[14] Himanshu Thapliyal and Nagarajan Ranganathan, Design of Reversible Sequential Circuits Optimizing Quantum Cost, Delay, and Garbage Outputs, ACMJournal onEmerging Technologies inComputer Systems,Vol. 6,No. 4,Article 14, Pub. date:December 2010.

[15] Siva Kumar Sastry, Hari Shyam Shroff, Sk.Noor Mahammad, V. Kamakoti", Efficient Building Blocks for Reversible Sequential Circuit Design" 1-4244-0173-9106/$20.00©2006IEEE

[16] Mozammel H A Khan and Marek Perkowski, Synthesis of Reversible Synchronous Counters, 2011 41st IEEE International Symposium on Multiple-Valued Logic, 0195-623X/11 $26.00 © 2011 IEEE

[17] V.Rajmohan, V.Ranganathan,"Design of counter using reversible logic" 978-1-4244-8679-3/11/$26.00 ©2011 IEEE.

[18] H. Thapliyal and A. P. Vinod, "Design of reversible sequential elements with feasibility of transistor implementation" In Proc. the 2007 IEEE Intl. Symp. On Cir.and Sys., pages 625–628, New Orleans, USA, May 2007.

# Modelling and Output Power Evaluation of Series-Parallel Photovoltaic Modules

Dr. Fadi N. Sibai

Process & Control Systems Department

Saudi Aramco

Dhahran 31311, Kingdom of Saudi Arabia

*Abstract*—**Solar energy has received attention in the Middle East given the abundant and free irradiance and extended sunny weather. Although photovoltaic panels were introduced decades ago, they have recently become economical and gained traction. We present a mathematical model for series-parallel photovoltaic modules, evaluate the model, and present the I-V and P-V characteristic plots for various temperatures, irradiance, and diode ideality factors. The power performance results are then analyzed and recommendations are made. Unlike other related work, our evaluation uses standard spreadsheet software avoiding commercial simulation packages and application programming. Results indicate improved power performance with irradiance and parallel connections of cell series branches. Given the hot weather in our region, increasing the number of cells connected in series from 36 to 72 is not recommended.**

*Keywords—solar energy; series-parallel photovoltaic modules; mathematical modeling*

## I. INTRODUCTION

The push for renewable energy solutions is steady to reduce carbon emissions, protect the environment and population health, and effect climate and sea water level changes. Solar energy has received attention in the Middle East given the abundant and free irradiance and sunny weather. Although photovoltaic (PV) panels were introduced decades ago, they have recently become economical and gained traction. Saudi Arabia recently saw the successful installation of a 3.5 megawatt photovoltaic field, the largest solar power plant built in the country, and has plans to install 41 Gigawatts of solar power over the next 20 years. PV fields geographically distributed can be located near loads and cut down on fuel consumed by power generating plants.

PV cells generate DC electricity during the day, which can be immediately consumed or stored in batteries for future consumption. An inverter is used to convert the generated power from DC to AC. The PV panel has a number of solar cells connected in series (typically 36 or 72) with the possibility to connect PV cell series branches in parallel. The solar cell is a p-n junction fabricated in a silicon (or other material) wafer or semiconductor layer. The photovoltaic effect causes the electromagnetic radiation of the solar energy hitting the PV cells to generate electricity. Photons with energy greater than the band gap energy of the wafer's semiconductor material are absorbed, creating electron-hole pairs, which in turn create a photocurrent in the presence of the p-n junction's internal electric field. The photocurrent's magnitude rises with the

number of electron-hole pairs created, which is dependent on the irradiance level. Therefore the magnitude of the current generated by the PV module is proportional to the amount of incident solar radiation. PV module models have evolved to include detailed parameters and even multiple piecewise linear regions or better accuracy. However, most have focused on modeling a number of PV cells in series and stopped short of modeling multiple parallel branches of serially connected PV cells.

In this paper, we present a mathematical model for series-parallel photovoltaic modules, evaluate the model, and present the I-V (module current vs. module voltage) and P-V (module power vs. voltage) characteristic plots for various temperatures, irradiance and diode ideality factors. The power performance results are then analyzed and recommendations are made. Unlike others' work, we model series-parallel PV cells with multiple parallel branches each consisting of a number of PV cells in series, and evaluate the model with a basic spreadsheet application, avoiding application programming and commercial simulators. This is facilitated by the mathematical technique (i.e. Newton-Raphson) used in the model evaluation which makes the model evaluation possible with a spreadsheet application. Thus, a key advantage of our approach is that our model can be reused by a larger population of researchers with no access to commercial simulator packages.

PV modeling and simulations are reviewed in Section II. In Section III, we present our mathematical model for a series-parallel PV module. In Section IV, the nonlinear model is evaluated in Excel and the I-V and P-V characteristics are plotted. Moreover, the results are analyzed for a variety of temperatures, irradiances, diode ideality factors, numbers of PV cells in series, and numbers of PV series branches connected in parallel. The paper concludes in Section V.

## II. PV MODELING AND SIMULATIONS

PV cells and modules have been modeled [1-12] mathematically and/or within commercial simulators. Diode-based circuits have typically represented the PV module. PV models differ in the number of diodes they contain (1 or more). Maximum power tracking methods [2-5] were introduced to determine the best operating parameters for the module and allowing the PV module to generate the maximum power. MATLAB and SIMULINK codes for PV module have been presented in [5-8, 10-11]. Complete system simulation models including MPPT controller, battery, charger, and DC/DC converter were presented in [7]. Our work differs from prior

work in that we model series-parallel PV modules and evaluate the model with a simple spreadsheet application.

Effects of climatic and environmental factors on solar PV performance in arid desert regions were discussed in [9]. Movable non-fixed panels — that can reposition their surfaces toward the sun's incident rays and track the sun's path in the sky and dust removing techniques — are key to optimizing the PV panel's output. The highest solar energy production in the Arabian Gulf region was determined to occur between 11 a.m. and 2 p.m. Cooling techniques were recommended to boost the module's efficiencies in light of high temperature effects. Among relative humidity, temperature and dust, the PV module performance was found to be most affected by dust depositions during long operation times. Smart grid features and renewable energy integration requirements were reviewed in [13] which also introduced a micro inverter and simulated it proving that it meets leakage current standards.

Researchers [10] used various methods for modeling PV panels. Azab [11] presented a piecewise linear model with three diodes that allows simulations of mismatched panels working under different conditions. This piecewise linear model substitutes the single diode by three parallel diodes in the PV model and defines regions where different diodes are turned on. This piecewise linear results in smoother I-V and P-V curves than with the one diode model. Dondi et al [12] presented a PV cell model for solar energy powered sensor networks, suitable over a wide range of irradiance intensity, cell temperature variation and incident angle. Our simple model does not use Azab's more precise piecewise linear model which increases the number of diodes (complexity) in the model, as we are mainly interested in determining the PV module's maximum power output, which our model accurately derives.

Herein, we modify and generalize the PV model in [6] to model series-parallel PV cell modules with multiple parallel branches of PV cells in series.

## III. SERIES-PARALLEL PV MODULE MODEL

The circuit model of the PV cell is a current source $I_L$ in parallel with a diode, and the two parallel branches in series with a series resistance Rs, as shown in Fig. 1.a. A load is placed between the + and − terminals on the right side. The photo current IL is proportional to the amount of light hitting the PV cell. At night or during the absence of light, the PV cell acts as a diode. The diode is responsible for defining the I-V characteristics of the cell. Other models include a third parallel branch with a shunt resistance as a single component but this third branch can be omitted without much loss in accuracy.

The circuit model of Fig. 1.b for $N_S$ PV cells in series is identical to the one depicted in Fig. 1.a, except that the $N_S$ resistances in series add up to a total resistance of $N_S R_S$, and there are now $N_S$ cells in series. As most PV modules include multiple cells in series ($N_S$ = number of cells in series = 36 or 72) and some designs connect multiple cell series branches in parallel ($N_P$ = number of parallel series-connected cell branches), the PV model reduces to the circuit model shown in Fig. 2, where I and V are the module current and module voltage, respectively.



a. One PV cell



b. Ns PV cells in series

Fig. 1. PV cell models



Fig. 2. PV panel model with Np parallel branches, each with Ns cells in series

Note that in Fig. 2, the $N_P$ parallel resistances reduce the total equivalent resistance to $N_S R_S / N_P$. The current source current of the parallel branches combine into the new current source $N_P I_L$. In the following sections, we use the variable "suns" to represent irradiance, where suns is the ratio of the solar irradiance over the solar irradiation under standard operating condition, and is assumed to be 1 (i.e., 1KW/m$^2$) in this paper, unless otherwise indicated.

Nomenclature of the mathematical PV model follows.

$I_L$:          photo current

$I_0$:          saturation current

$I_S$:          photo current

Suns:          irradiance = $G/G_{nom}$ (where 1sun = 1KW/m$^2$)

TaC:          temperature in degrees Celsius

TaK:          temperature in degrees Kelvin

Vg:          band gap energy = 1.12V at temperature T1

T1:          reference start temperature (in deg. K)

T2:          reference end temperature (in deg. K)

T:          junction temperature (in deg. K)

$V_{OC}\_T1$:          open-circuit voltage at temperature T1

$V_{OC}\_T2$:          open-circuit voltage at temperature T2

$I_{SC}\_T1$:          short-circuit current at temperature T1

$I_{SC}\_T2$:          short-circuit current at temperature T2

$R_S$:          series resistance

$N_S$:          number of series-connected cells

$N_P$:          number of parallel-connected cells

q:          electron charge = 1.6E-19 Coul.

k:          Boltzmann's constant = 1.38E-23 J/K

G:          irradiance

n:          diode ideality factor between 1(ideal) and 2

The PV module current I pictured in Fig. 2 is the photo current $I_L$ minus the diode current and is given by

$$I = N_P\, I_L - N_P\, I_0\, (e^{q\ (V/Ns + I\ Rs/Np)/(n\ k\ TaK)} - 1) \tag{1}$$

Note that I appears on both sides of the equation. As under constant radiance and temperature $I_L$ is constant, I is therefore a constant minus the exponential-shaped diode current, and takes the shape of the curve of Fig. 3.
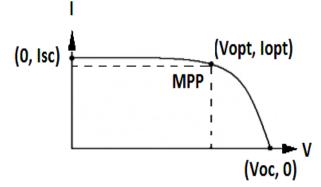


Fig. 3.   Typical I-V characteristic curve of a PV module

The point (Voptimal, Ioptimal) is of interest as it is the point at which the PV module's output power is maximum and is therefore a desired operation point. The I-V curve crosses the horizontal V axis at the open-circuit voltage, Voc, obtained by connecting no load between the + and − terminals of the output and thus occurs when I=0. The I-V curve crosses the vertical I axis at Isc, the short circuit current, obtained by shorting the + and − terminals of the PV output, which occurs when V=0. The line connecting the origin (V, I)=(0, 0) with the maximum

power point (MPP) (V, I)= (Vopt, Iopt) has a slope equal to 1/Ropt, where Ropt=Vopt / Iopt.  At the MPP, the output power Iopt x Vopt is maximum. The ratio of the MPP over the product of the PV cell area and the ambient irradiation gives the maximum efficiency. When the + and − terminals of the PV output are connected to a load, two situations occur:

*1) When the load resistance is small, the PV cell's operating point is on the left side of the Fig. 3 curve and the current is almost constant and very near the short circuit current Isc. The Isc current is the maximum possible current obtained with zero load.*

*2) When the load resistance is large, the PV cell's operating point is on the right side of the Fig. 3 curve and the voltage is almost constant and very near the open-circuit voltage Voc. The Voc voltage is the maximum possible voltage obtained with infinite load (open circuit), and represents the voltage of the PV cell in the dark ($I_D= I_L$ or I=0). It is equal to $Vt\_Ta\ ln(I_L/I_0)$, where Vt_Ta is defined in eq. (6) below.*

Several PV module characteristics are provided by the PV manufacturer.

The photo current is expressed by

$$I_L = I_{SC}\_T1\ suns + (TaK - T1)\ (I_{SC}\_T2 - I_{SC}\_T1)\ /\ (T2 - T1) \tag{2}$$

The reverse saturation current $I_0$ is given by

$$I_0 = (I_{SC}\_T1/(e^{q\ Voc\_T1/(Ns\ n\ k\ T1)} - 1))\ \ \times\ (TaK\ /\ T1)^{3/n}$$
$$\times\ e^{q\ Vg\ (1/TaK - 1/T1)\ /(Ns\ n\ k)} \tag{3}$$

The series resistance is given by

$$R_S = dV/dI\,\big|_{V=Voc} - 1\ /\ Xv\ =\ -1.15\ /\ N_S\ /2 \tag{4}$$

where

$$Xv = I_0\ (q\ /\ (n\ k\ T1))\ e^{q\ Voc\_T1/(Ns\ n\ k\ T1)} \tag{5}$$

Combining the product of n, k, Tak, and 1/q, we obtain

$$Vt\_Ta = n\ k\ TaK\ /\ q \tag{6}$$

A common value for n is 1.2 for Silicon mono and 1.3 for Silicon poly, 1.3 for AsGa, and 1.5 for CdTe.

The existence of Rs renders eq. (1) a nonlinear problem to solve. Typically, the Newton-Raphson method is used to solve such problems, owing to its fast convergence as previously reported. The Newton-Raphson method operates as follows. Given y, a function of x, such that y=f(x)=0, then

$$x = x_0 + [df/dx\,\big|_{x=x_0}]^{-1}\ (y - f(x_0)) \tag{7}$$

and iteratively x can be solved as follows

$$x_{n+1} = x_n - f(x_n)\ /\ f'(x_n) \tag{8}$$

In our case, we reshape eq. (1) in the form "y=f(I)=0" as follows

$$y = f(I) = -I + N_P\ I_L - N_P\ I_0\ (e^{q\,(V/N_S + I\,R_S/N_P)/(n\,k\,TaK)} - 1) = 0$$

(9)

and, according to eq. (8), the iterative Newton-Raphson solution becomes

$$I_{n+1} = I_n - f(I_n)\ /\ f'(I_n)$$

$$= I_n - \{-I_n + N_P\ I_L - N_P\ I_0\ (e^{q\,(V/N_S + I_n\,R_S/N_P)/(n\,k\,TaK)} - 1))\ /$$

$$(-1 - (I_0\ R_S/(n\,k\,TaK))\ e^{q\,(V/N_S + I_n\,R_S/N_P)/(n\,k\,TaK)})\}$$

which can be rewritten, using eq. (6), as

$$I_{n+1} = I_n - \{(-I_n + N_P\ I_L - N_P\ I_0\ (e^{(V/N_S + I_n\,R_S/N_P)/V_{t\_Ta}} - 1))\ /$$

$$(-1 - (I_0\ R_S/(n\,k\,T))\ e^{(V/N_S + I_n\,R_S/N_P)/V_{t\_Ta}})\}$$

(10)

Given input values of the module voltage V, TaC, n, and suns, equations (2), (3), (4) and (6) are solved, and eq. (10) is iteratively solved five times to obtain I. The output PV module power is the product of the final value of I (after five iterations) and V.

## IV. RESULTS

The model equations in the previous section were entered into Microsoft Excel and evaluated. Voltage V was varied in the range 0-25V, temperature was varied in the range -10 deg. C. and 80 deg. C, and n was varied in the range 1.2-1.8. For each value of V, temperature and n, a value of I was computed from which a value of P (I x V) was computed. In this evaluation, the electrical characteristics of the MSX-60 PV module were used: Pmax=17.1V x 3.5A=60W, Isc=3.8A, Voc=21.1V, temperature coefficient of Voc= -(80 +/- 10) mV/degree Celsius, and temperature coefficient of Isc= (0.0065 +/- 0.015)% / degree Celsius. From these values, we derive: T1=25 deg. C, T2=75 deg. C, Voc_T1=21.06 volts/Ns, Voc_T2=17.05 volts/Ns, Isc_T1= 3.8A, Isc_T2= 3.92A, Xv=123.2046 Siemens, and Rs= 0.0078556 ohms.

The results depicted in Figs 4-17 assume that $N_S$ is 36 (i.e. 36 cells in series), $N_P$ is 1, n is 1.2, and suns is 1 (i.e. illumination of $1kW/m^2$), unless otherwise indicated below.

Fig. 4 plots the output current I versus the output voltage V for temperatures between 10-80 degrees Celsius. For a low and fixed output voltage and for output currents near their maximum short circuit current values, the output current I (left side of Fig. 4) generally increases with increasing temperatures.

For lower but fixed I, and for V values near their maximum open circuit values, the output voltage V increases with decreasing temperatures (right side of Fig. 4). Optimum MPP values of I and V occur near the curves' knees.



Fig. 4. ig. 4. I-V plots for various temperatures with Ns = 36, n = 1.2, suns = 1, Np = 1

Fig. 5 plots the output power versus the output voltage for the same temperature range, clearly indicating the $V_{opt}$ values causing the maximum (desired) output power values generated by the PV module. After deriving the $V_{opt}$ value, the optimal I value, $I_{opt}$, can then be easily derived from the plots of Fig. 4. Higher maximum output power values are achieved at decreasing temperatures, and for higher $V_{opt}$ values. Lower temperatures achieve the highest generated powers.



Fig. 5. P-V plots for various temperatures with Ns = 36, n = 1.2, suns = 1, Np = 1

It is desirable to estimate the effect of varying the diode ideality factor, n, and suns. For that purpose, Fig. 6 plots the I-V relationship as n is increased from 1.2 to 1.8, causing only slight changes near the curve's knee. The assumed temperature was 20 deg. C.

Fig. 6. I-V plots for various n values with Ns = 36, suns = 1, Np = 1

Fig. 7 plots the I-V relationship as suns is increased to 2 (from 1). The doubling of suns causes a doubling in current at low voltages (left side) and a slight increase in the voltage at low currents (right side). Significantly, doubling suns also doubles the output power for the same output voltage value, as shown in Fig. 8, and pushes the higher voltage values (right side of the plot) slightly further to the right when n is increased from 1.2 to 1.8, as shown in Fig. 9. As depicted by Fig. 8, the irradiance has a strong effect on the PV module's output power.



Fig. 7. I-V plots for various temperatures with Ns = 36, suns = 2, Np = 1
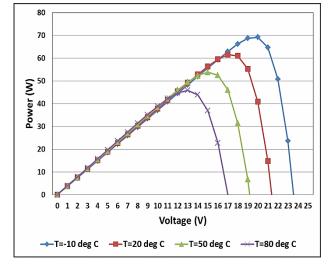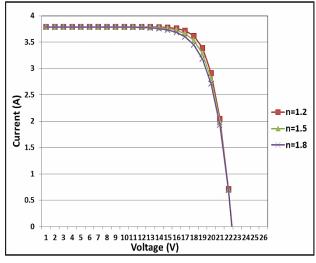


Fig. 8. P-V plots for various temperatures with Ns = 36, n = 1.2, suns = 2, Np = 1



Fig. 9. I-V plots for various n values with Ns = 36, suns = 2, Np = 1

Bringing back suns to 1, and increasing $N_S$ to 72 (from 36) cause the gaps between the I-V curves at the different temperatures to widen and the I values to slightly drop at lower output voltage values, as shown in Fig. 10, and as compared to Fig. 4. These parameter changes also cause the maximum output powers to improve and rise slightly for T = -10 deg. Celsius. Unfortunately, maximum output powers drop for the higher temperatures 20-80 deg. Celsius, as shown in Fig. 11.

Fig. 10. I-V plots for various temperatures with Ns = 72, suns = 1, Np = 1



Fig. 12. I-V plots for various temperatures with Ns = 36, suns = 1, Np = 10



Fig. 11. P-V plots for various temperatures with Ns = 72, suns = 1, Np = 1



Fig. 13. P-V plots for various temperatures with Ns = 36, suns = 1, Np = 10

Bringing back $N_S$ to 36, and increasing $N_P$ to 10 (from 1; i.e. 10 PV panels connected in parallel) cause the I-V curve to drastically move up and the Isc's to be multiplied by 10 as shown in Fig. 12. Luckily, as shown in Fig. 13, the power and peak power values also shot up by 10x. Small changes in I are observed as n is changed from 1.2. to 1.5 and then 1.8, as shown in Fig 14.

When suns is doubled to 2 (with Ns = 36 and Np = 10), the current I and power are further doubled from the case with Ns = 36, Np = 10 and suns = 1, as depicted by Figures 15, 16, and 17. The peak powers are now in the kilowatt range benefitting from the dual effect of enhanced irradiance and cumulative current from multiple parallel modules. When n is varied, slight changes in I are observed but no longer restricted to the optimal I value knee region but also extending to the bottom of the I curve for small I values, as depicted by Fig. 17.



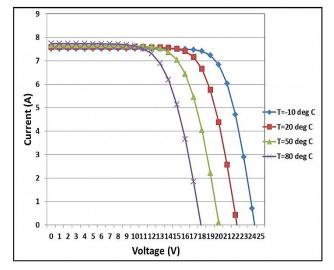Fig. 14. I-V plots for various n values with Ns = 36, suns = 1, Np = 10

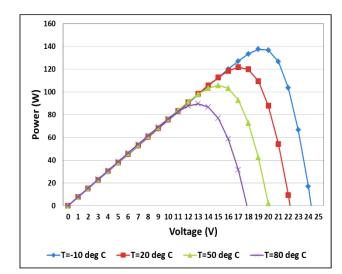Fig. 15. I-V plots for various temperatures with Ns = 36, suns = 2, Np = 10



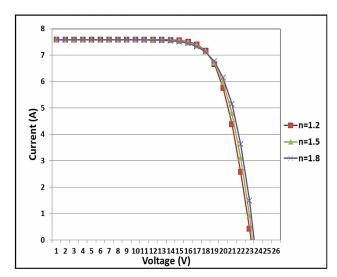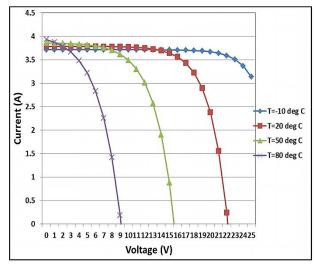Fig. 16. P-V plots for various temperatures with Ns = 36, suns = 2, Np = 10



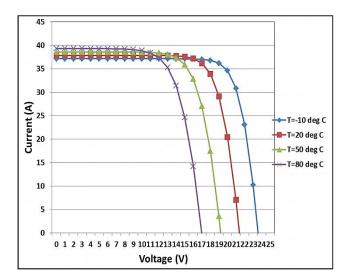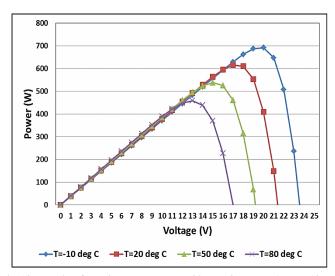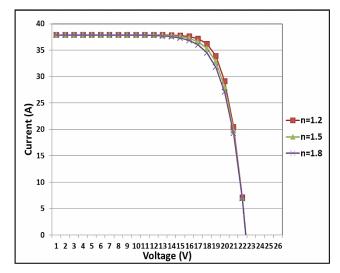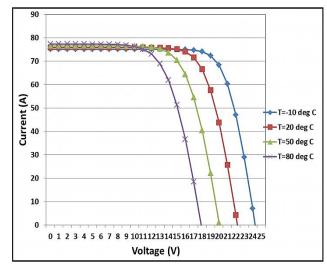Fig. 17. I-V plots for various n values with Ns = 36, suns = 2, Np = 10

TABLE I.      OUTPUT POWER COMPARISONS WITH T = 20 DEGREES CELSIUS AND N = 1.2

| Ns | Suns | Np | Max Power (W) |
|----|------|-----|---------------|
| 36 | 1 | 1 | 61 |
| 36 | 2 | 1 | 120 |
| 72 | 1 | 1 | 58 |
| 36 | 1 | 10 | 610 |
| 36 | 2 | 10 | 1210 |

To compare the effects of varying $N_S$, $N_P$, and suns, on the maximum PV module output power, we assemble these values in Table I. Table I compares the output power values (in W) for various values of Ns, Suns, and Np, and for a temperature of 20 degrees Celsius and an n of 1.2. Given that the power is higher with (Ns = 36, Suns = 2) than with (Ns = 72, Suns = 1), operation investments, for instance in automatic sun tracking and repositioning and in panel glass surface dust cleaning, are more profitable than connecting more PV cells in series.

It should be noted that the Fill factor is the ratio of the MPP over the product Voc x Isc. It has a negative relationship with the temperature. Good quality PV modules have a Fill factor above 0.7.

The results with a series of multiple PV cells correlate well with the commercial data sheet of the MSX-60 PV module whose parameters where used in our evaluation model.

## V. CONCLUSION

Renewable energy including solar sources is planned to consume a larger portion of the GCC regional electric power generation pie chart in the coming years.

This paper presented and evaluated a mathematical model for series-parallel photovoltaic modules based on standard spreadsheet. The inputs to the model were the module voltage V, temperature, irradiance, and diode ideality factor. The model's output are the module current I and module power P. Our evaluation yielded accurate results and was based on Microsoft Excel and avoided Matlab/Simulink or Spice simulation software packages, or application development. This is desirable to make the mathematical model usable to a larger population of researchers with no access to commercial simulators and without the need for coding and developing a software application.

Results indicate improved power performance with higher irradiance (suns) and more parallel connections of PV cell series branches (higher $N_P$). The temperature impact is also significant, rendering higher $N_S$ values detrimental in our high temperature region. In contrast, the impact of the diode ideality factor is rather small.

Tracking mechanisms for repositioning the PV modules toward the sun's rays and dust cleaning methods are therefore recommended to keep the PV module generating power near its MPP point. Given the hot weather in our region, increasing the number of cells connected in series (Ns) from 36 to 72 is not recommended.

As to future work, we may investigate dust cleaning options for PV modules and compare them based on output power performance, water requirements, and initial and long term costs.

REFERENCES

[1] J. A. Gow, and C. D. Manning, "Development of a photovoltaic array model for use in power electronics simulation studies," IEE Proceedings of Electric Power Applications, Vol. 146, No. 2, pp. 193–200, March 1999.

[2] K. H. Hussein, I. Muta, T. Hshino, and M. Osakada, "Maximum photovoltaic power tracking: an algorithm for rapidly changing atmospheric conditions," Proc. Inst. Elect. Eng, vol. 142, no. I, pp. 59-64, Jan. 1995.

[3] J. Youngseok, S. Junghun, Y. Gwonjong and C. Jaeho, "Improved perturbation and observation method (IP&O) of MPPT control for photovoltaic power systems," Proc. 31st Photovoltaic Specialists Conference, Lake Buena Vista, Florida, pp. 1788 – 1791, January 2005.

[4] Nicola Femia, Giovanni Petrone, Giovanni Spagnuolo, and Massimo Vitelli, "Optimization of perturb and observe maximum power point tracking method," IEEE Transactions on Power Electronics, vol. 20(4), pp. 963-973, 2005.

[5] G. Walker, "Evaluating MPPT converter topologies using a MATLAB PV model," Journal of Electrical & Electronics Engineering, Australia, IEAust, 21(1), pp. 49-56, 2001.

[6] Francisco M. Gonzalez-Longatt, "Model of photovoltaic module in Matlab," 2 CIBELEC, 2005.

[7] T. DenHerder, Design and Simulation of Photovoltaic Super System Using SIMULINK, Senior Project, Electrical Engineering Dept. California Polytechnic Univeristy, San Luis Obispo, 2006.

[8] H. Tsai, C. Tu, and Y. Su, "Development of generalized photovoltaic model using MATLAB/SIMULINK," Proc. World Congress on Engineering and Computer Science (WCECS 2008), Oct. 2008.

[9] F. Touati, M. Al-Hitmi, and H. Bouchech, "Towards understanding of the effects of climatic and environmental factors on solar PV performance in arid desert regions (Qatar) for various PV technologies," Proc. First Int. Conference on Renewable Energies and Vehicular Technology, 2012.

[10] Gwinyai Dzimano, Modeling of Photovoltaic Systems, MS Thesis, Electrical and Computer Engineering, Ohio State University, 2008.

[11] Mohamed Azab, "Improved circuit model of photovoltaic array," Int. Journal of Electrical Power and Energy Systems Engineering, vol. 2(3), 2009.

[12] D. Dondi, et al., "Photovoltaic cell modeling for solar energy powered sensor networks," Proc. 2nd IEEE International Workshop Advances in Sensors and Interface (IWASI 2007), 2007.

[13] M. Bouzguenda, A. Gastli, A. H. Al Badi, and T. Salmi, "Solar photovoltaic inverter requirements for smart grid applications," Proc. Innovative Smart Grid Technologies (ISGT)-Middle East, Jeddah, Saudi Arabia, 2011.

# Mining of Web Server Logs in a Distributed Cluster Using Big Data Technologies

Savitha K

Dept. of Computer Science, Research Scholar
PSGR Krishnammal College for Women
Coimbatore, India.

Vijaya MS

Dept. of Computer Science, Associate Professor
GR Govindarajulu School Of Applied Computer Technology
Coimbatore, India.

*Abstract*—**Big Data is an emerging growing dataset beyond the ability of a traditional database tool. Hadoop rides the big data where the massive quantity of information is processed using cluster of commodity hardware. Web server logs are semi-structured files generated by the computer in large volume usually of flat text files. It is utilized efficiently by Mapreduce as it process one line at a time. This paper performs the session identification in log files using Hadoop in a distributed cluster. Apache Hadoop Mapreduce a data processing platform is used in pseudo distributed mode and in fully distributed mode. The framework effectively identifies the session utilized by the web surfer to recognize the unique users and pages accessed by the users. The identified session is analyzed in R to produce a statistical report based on total count of visit per day. The results are compared with non-hadoop approach a java environment, and it results in a better time efficiency, storage and processing speed of the proposed work.**

*Keywords—Big data; Hadoop; Mapreduce; Session identification; Analytics*

## I. INTRODUCTION

A data is a collection of facts from the grids of web servers usually of unorganized form in the digital universe. Around 90% of data present in today's world are generated in last two years [1]. A large amount of the data available in the internet is generated either by individuals, groups or by the organization over a particular period of time. The volume of data becomes larger day by day as the usage of World Wide Web makes an interdisciplinary part of human activities. Rise of these data leads to a new technology such as big data that acts as a tool to process, manipulate and manage very large dataset along with the storage required.

Big Data is a high volume, high velocity and high variety information assets that demand cost-effective [2], innovative forums of information processing for enhanced insight and decision making. Big data, a buzz word in the business intelligence can handle petabytes or terabytes of data in a reasonable amount of time. Big data is distinct from large existing database which uses Hadoop framework for data intensive distributed applications.

Big Data analytics applies advanced analytical techniques of large datasets to discover hidden patterns and other useful information. It is performed using software tools mainly for predictive analysis and data mining. The growing number of technologies is used to aggregate, manipulate, manage and analyze big data. Some of the most prominent technologies are

[3] NoSQL databases that include Cassandra, MongoDb, redis, Hbase and Hadoop framework includes Hadoop, HDFS, Hive, Pig.

Data can come from a variety of sources and in a variety of types that includes not only structured traditional relational data, but also semi-structured and unstructured data. Machine generated data such as click stream logs and email logs of unstructured data are larger in magnitude when compared with human generated data that cannot fit in a traditional data warehouses for further analysis.

Numerous research works are carried out in web log mining, hadoop and some of them are reviewed below. Murat et al., [4] proposed the smart miner framework that extracts the user behavior from web logs. The framework used the smart session construction to trace the frequent user access paths. Sayalee Narkhede et al., [5] introduced the Hadoop-MR log file analysis tool that provides a statistical report on total hits of a web page, user activity, traffic sources. This work was performed in two machines with three instances of hadoop by distributing the log files evenly to all nodes.

Milind Bhandare et al., [6] put forth a generic log analyzer framework for different kinds of log files such as a database or file system. The work was implemented as a distributed query processing to minimize the response time for the users which can be extendable for some format of logs. Parallelization of Genetic Algorithm (PGA) was suggested by Kanchan Sharadchandra Rahate et al., [7]. PGA uses OlexGA package for classifying the document. The train model data is stored in HDFS and the test model categories the text document.

A framework for unstructured data analysis was proposed by Das et al., [8] using big data of public tweets from twitter. The tweets are stored in Hbase using Hadoop cluster through Rest Calls and text mining algorithms are processed for data analysis.

The semi structured log files are large datasets which are challenging to store, search, share, visualize and analyze. Almost 26% of web log types of data require big data technology to perform an analysis [9]. In order to improve the usage of a website and to track the user behavior, in online advertising and E-commerce the web log mining is performed using Hadoop.

The related works so far stated above performs the work with good scalability but fails to experiment the time efficiency between the different modes of hadoop and

necessity of the scalability. The proposed work analyses the working of both hadoop modes and the time efficiency in each mode specifically for semi structure log data, along which a statistical report is made.

This paper effectively utilizes the web logs of NASA website accessed by the user, to identify the session using Apache Hadoop Mapreduce in a distributed cluster. As hadoop does not enforce schema based storage, it processes the semi structured log files easily. The file holds 550MB dataset collected from various time period of the same year, and is stored in HDFS which is scattered in the cluster as blocks through high speed network. The work encompasses the use of both the modes along with non-hadoop approach. The identified session is analyzed based on date and number of times visited using R tool.

## II. HADOOP MAPREDUCE

Hadoop is a flexible infrastructure for large scale computation and data processing on a network of commodity hardware. It allows applications to work with thousands of computational independent computers and petabytes of data. The main principle of hadoop is moving computations on the data rather the moving data for computation. Hadoop is used to breakdown the large number of input data into smaller chunks and each can be processed separately on different machines. To achieve parallel execution, Hadoop implements a MapReduce programming model.

MapReduce a java based distributed programming model consists of two phases: a massively parallel "Map" phase, followed by an aggregating "Reduce" phase. MapReduce is a programming model and an associated implementation for processing and generating large data sets [10]. A map function processes a key/value pair (k1,v1,k2,v2) to generate a set of intermediate key/value pairs, and a reduce function merges all intermediate values [v2] associated with the same intermediate key (k2) as in (1).

Maps are the individual tasks that transform the input records into intermediate records. A MapReduce job usually splits the input data-set into independent chunks which are processed by the map tasks. The framework sorts the output of the map, which are then input to the reduce tasks. Both the input and the output of the processed job are stored in a file-system. Typically just zero or one output value is produced by the reducer. In MapReduce, a mapper and reducer is identified by the following signature,

$$\text{Map } (k1, v1) \rightarrow [(k2, v2)]$$

$$\text{Reduce } (k2, [v2]) \rightarrow [(k3, v3)] \qquad (1)$$

MapReduce suits applications where data is written once and read many times [12]. The data stored in a file system namespace contributes to HDFS which allows master-slave architecture.

The cluster consists of a single NameNode, a master that manages the file system namespace and regulates its access to files by clients. There can be a number of DataNodes usually one per node in the cluster which periodically report to NameNode, the list of blocks it stores.

HDFS replicates files for a configured number of times. It automatically re-replicates the data blocks on nodes that have failed. Using HDFS a file can be created, deleted, copied, but cannot be updated. The file system uses TCP/IP for communication between the clusters. A file is made of several blocks and the target machine holds each block that is chosen randomly on a block-by-block basis. Thus access to a file requires the cooperation of multiple machines, but supports file sizes larger than a single-machine DFS.

## III. WEB LOG MINING USING NON HADOOP APPROACH

The current processing of log files goes through ordinary sequential ways in order to perform preprocessing, session identification and user identification. The developed non hadoop approach loads the log file dataset, to process each line one after another. The log field is then identified by splitting the data and by storing it in an array list as shown in Fig.1. The preprocessed log field is stored in the form of hash table, with key and value pairs, where key is the month and value is the integer representing the month.



Fig. 1. Log file stored as table with unique fields

The developed work is possible to run only on single computer with a single java virtual machine (jvm). A jvm has the ability to handle a dataset based on RAM i.e. if the RAM is of 2GB then a jvm can process dataset of only 1GB. Processing of log files greater than 1GB becomes hectic. The non hadoop approach is performed on java 1.6 with single jvm. Although batch processing can be found in these single-processor programs, there are problems in processing it due to limited capabilities. Therefore, it is necessary to use parallel processing approach to work effectively on massive amount of large datasets.

## IV. LOG MINING USING HADOOP APPROACH

Hadoop framework access a large semi structure data in a parallel computation model. Log files usually generated from the web server comprise of large volume of data that cannot be

handled by a traditional database or other programming languages for computation. The proposed work aims on preprocessing the log file and keeps track on sessions accessed by the user, using Hadoop and the system architecture is shown in Fig.2. The work is divided into phases, where the storage and processing is made in HDFS.
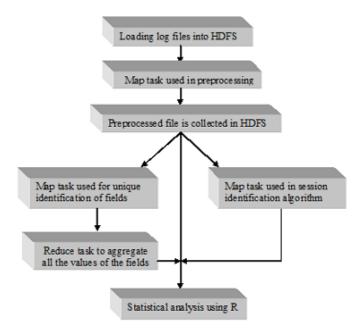


Fig. 2.   System Architecture of the proposed work

Data cleaning is the first phase carried out in the proposed work as a preprocessing step in NASA web server log files. The NASA log file contains a number of records that corresponds to automatic requests originated by web robots, that includes a large amount of erroneous, misleading, and incomplete information. In the proposed work the web log file containing request from robots, spider and web crawlers are removed. These requests are a program that automatically downloads a complete website by following each hyperlink on every page. Request created by web robots are not considered as used data and consequently, it is filtered out from the NASA log data.

The entries that have status of "error" or "failure" have been removed. Also some access records generated by automatic search engine agent is identified and removed from the access log. The foremost important task carried out in data cleaning is the identification of status code. All the log lines even though satisfy the above constraints, only the lines holding the status code value of "200" is identified as correct log. The corresponding fields containing the correct status code are forwarded to the output text file.

The major advantage of this step is to eliminate an error that leads to accurate, error free log data, which produce a quality result and increase in efficiency. After applying data cleaning step, applicable resources are stored in the HDFS as text file, and feed back to the session identification algorithm as input file. The cleaned web log data is used for further

analysis of session identification utilized by the NASA users and also to identify unique user, unique URLs accessed.

Data cleaning, a preprocessing method is applied in the proposed work to filter and minimize the original size of data. The log file that resides in HDFS is given as input to the MapReduce job through FileInputFormat which is an abstract class. Input to map task is given as Long Writable and Reduce task is set to zero as there is no need for aggregation of values.

The TextInputFormat is used for text files to break the files into lines consisting of new line character. The key is the IP address and remaining all the fields are considered as values. Pattern Matching is used to separate the fields, and only log line of successful status count 200 is used for further processing. The processed line is checked to find out 'GET' method and records with "robots.txt" in the requested URL are identified and removed. The preprocessed log value includes the timestamp, date, browser version and total bytes sent for a request.

The session identification plays a major role in web log mining. The algorithm splits the sessions based on unique IP address and time spent by the user from login till logout. If the user exceeds time limit of 30 minutes, a new session number is generated for the similar IP address. The date function is used to calculate the time difference between the same IP address. The map function holds the session number as key and the remaining log fields as values. The reduce function is assigned to zero as no more aggregation is needed. The processed session file also resides in the HDFS that can be downloaded for further analysis.

The preprocessed log file is used to find the user identification, as Mapreduce in general identifies the unique values based on key value pair. The log file consists of different fields and the user is identified based on IP address, which is considered as key and their corresponding count as value. Once all the keys are found, the combiner is used to combine all the values of a specific key, the aggregated result is passed to the reduce task and total counts of IP accessed is listed as text file in HDFS.

Other than the identification of unique user, unique fields of date, url referred, and status code is also identified. These unique values is retrieved and used for further analysis in order to find the total url referred on a particular date or the maximum status code got successes on specific date.

## V.    RESULTS AND INTERPRETATIONS

The web server logs are mined for efficient session identification using Hadoop Mapreduce. This framework is used to compute the log processing in pseudo and fully distributed modes of cluster. The NASA web server logs gathered in four different files are used for processing in hadoop environment. The log data is collected from four different ASCII files of NASA with one line per request. Timestamps have 1 second resolution. The session identification algorithm is performed using the Mapreduce approaches. The log files include HTTP request of NASA web site. The process is analyzed in Ubuntu 12.04 OS with Apache Hadoop-0.20.2 [11].

## A. Pseudo Distributed Mode

Hadoop framework consist of five daemons namely Namenode, Datanode, Jobtracker, Tasktracker, Secondary namenode. In pseudo distributed mode all the daemons run on local machine simulating a cluster. The job tracker and task tracker is set to idle in this mode. Preprocessing is done to eliminate the error that leads to accurate, error free log data, which produce a quality results and increase in efficiency. After applying data cleaning step applicable resources are stored in the HDFS as text file, and feed back to the session identification algorithm as input file. The cleaned web log data as shown in Fig.3 is used to analyze the session identification, unique user and unique URLs. Fig.4 depicts the storage of log data in localhost: 50070 which can be downloaded for further analysis.



Fig. 3. Preprocessed Log Data



Fig. 4. File System storage in localhost

The Map phase maps one set of data to another set of data by a one-to-one rule. It is highly difficult to read the huge file without the use of HDFS where the files are loaded and processed. The total size of 516MB dataset is divided into counters or blocks of same size and processed using a single map task. The working of hadoop daemons and storage is made visible using different web user interfaces.

## B. Fully distributed mode

The fully distributed mode makes use of all five daemons by assigning a job to each one. The cluster is created with two machines of same configuration with 2 GB RAM, Intel® Pentium® Dual CPU T3200 processor and Ubuntu 12.04 OS.

A cluster of two nodes is created with which one acts as slave and other as both master and slave where data is transferred in 100Mb/s speed. The master node contains daemons such as namenode, datanode, jobtracker and tasktracker. The slave node contains tasktracker and datanode. Hadoop clusters are capable to build with inexpensive computers. Even if one node fails, the cluster continues to operate without losing data by redistributing the work to remaining cluster.

The master node contains the IP address of the slave. The Name Node contains all the file system metadata of the developed cluster. The HDFS data is replicated twice as the cluster contains only one slave node. The slave node is identified by the master using its ip address and read the log files which are splitted into block size of 64MB that is default, as in Fig.5 and save it in its hard drive.

The log data is splitted into blocks and processed in two machines in parallel. The code is transmitted from the master to slave to process the work for the given block. The slave node mapper completes the preprocessing task and session identification process and shifts its data to master node. The job tracker holds the metadata of log files, the metadata includes 1 task level with 12 counters.



Fig. 5. Data processing in hadoop file system

The Mapreduce task in cluster is broken into record reader, mapper, combiner, and partitioner. The record reader passes the web log data to the mapper with ip address as key and the remaining fields as value. The map task groups all the data, and the combiner task is set to idle as the proposed work only uses the map task. The IP address 172.24.1.229 acts as a master and the IP 172.24.1.237 acts as master and slave. Both the machine communicates with each other using the IP. The namenode UI for the cluster can be tracked in the master node using the localhost: 54310 as shown in Fig.6 that contains the cluster summary of live and dead nodes.



Fig. 6. Namenode UI with 2 live nodes

The preprocessing work executed in parallel results in 2.15 minutes and 3.04 minutes for session identification of 550MB NASA web log data file. The same work when performed in single node produces 2.48 minutes for preprocessing and 3.02 minutes for session identification. The performance evaluation of Non hadoop approach, Pseudo distributed mode and fully distributed mode is shown in Table-I and Fig.7, which proves that performing the work in distributed mode improves the time in few milliseconds. Expanding the cluster and using terabytes of data would result in better time efficiency.

TABLE I. PERFORMANCE OF DIFFERENT APPROACHES

| NASA Server logs | Milliseconds | Minutes |
|---|---|---|
| Non hadoop approach | 369391 | 6.15 |
| Pseudo distributed mode | 330001 | 5.50 |
| Fully distributed mode | 311400 | 5.21 |



Fig. 7. Performance Evaluation of different approaches in hadoop

The preprocessed data is analyzed using R, a free statistical programming tool. The log files don't contain any specific format, due to which fields could not be separated easily. Data frames are used in the analysis of log files to read the content of the file using read.csv ( ) [14], as comma separated values. Depending upon the format of the log file, the data in R is used to filter, analyze, or manipulate to make it more usable .Using data editor the data frame is created from the preprocessed log file with row ranging till 50, 00,000 as shown in Fig.8.

The package stringr is used to match the retrieved string from the pattern matching of logs. Stringr is used to ensure that function and argument names (and positions) are consistent. Each filed that matches the regular expression is stored as data frame with column names. The individual column is taken back and stored in a separate table to find the unique values as shown in Fig.9.



Fig. 8. Data Frame with preprocessed data

| | row.names | ipaddress | totalcount |
|---|---|---|---|
| 157925 | 129744 | ppp4.cowan.edu.au | 233 |
| 157926 | 156995 | winnie.fit.edu | 233 |
| 157927 | 18584 | 156.26.2.92 | 234 |
| 157928 | 21379 | 163.205.130.2 | 234 |
| 157929 | 24295 | 170.207.35.8 | 234 |
| 157930 | 30934 | 198.77.113.40 | 234 |
| 157931 | 37613 | 204.245.159.7 | 234 |
| 157932 | 55348 | bet.mse.uiuc.edu | 234 |
| 157933 | 56393 | bmixter.clark.net | 234 |
| 157934 | 58025 | butler406b.dorm.tulane.edu | 234 |
| 157935 | 58398 | cab.nysernet.org | 234 |
| 157936 | 60401 | champ.wnet.gov.edmonton.ab.ca | 234 |
| 157937 | 63434 | crl12.crl.com | 234 |
| 157938 | 63549 | crux.izmiran.rssi.ru | 234 |
| 157939 | 64354 | csdanxp4.med.utoronto.ca | 234 |
| 157940 | 64377 | cse.unl.edu | 234 |
| 157941 | 70374 | dial31.lakeheadu.ca | 234 |
| 157942 | 86363 | hiss101.teleserve.ca | 234 |
| 157943 | 99283 | jlyons.speedware.com | 234 |

Fig. 9. Data Frame of unique IP address with count

## VI. CONCLUSION

At present the Big data technology is successfully incorporated for all real domain problems such as web log analysis, fraud detection and text analysis. This paper reveals the importance of one of the Big data technology Hadoop, where the framework handles large amount of data in a cluster for web log mining. Data cleaning, the main part of preprocessing is performed to remove the inconsistent data. The preprocessed data is again manipulated using session identification algorithm to explore the user session. Unique identification of fields is carried out to track the user behavior. Both the algorithms process the NASA web server logs using Mapreduce task. The huge data is stored in HDFS as text files and is accessed around the cluster in blocks. From the results it is observed that the framework developed in hadoop has high performance when compared to the other approaches.

Future work can be extended to number of nodes for prediction analysis to find the next page that will be accessed by the user in a particular website. Also the resulting text files stored in HDFS can be binded with NoSQL databases, and other data mining tools for analysis.

REFERENCES

[1] http://www.web-datamining.net/

[2] Ruchi Verma, Sathyan R Mani, "Use of Big Data Tehnologies in Capital Markets," 2012 Infosys Limited, Bangalore, India.

[3] James Manyika, Brad Brown et.al, "Big Data: The next frontier for Innovation, Competition, and Productivity," McKinsey Global Institute, June 2011.

[4] Murat Ali , Ismail Hakki Toroslu, "Smart Miner: A New Framework for mining Large Scale Web Usage Data," WWW 2009, April 20-24. 2009 Madrid, Spain. ACM 978-1-60558-487-4/09/04.

[5] Sayalee Narkhede and Tripti Baraskar, "HMR Log Analyzer: Analyze Web Application Logs Over Hadoop MapReduce," International Journal of UbiComp (IJU) vol.4, No.3, July 2013.

[6] Milind Bhandare, Vikas Nagare et al., "Generic Log Analyzer Using Hadoop Mapreduce Framework," International Journal of Emerging Technology and Advanced Engineering (IJETAE), vol.3, issue 9, September 2013.

[7] Kanchan Sharadchandra Rahate et al., "A Novel Technique for Parallelization of Genetic Algorithm using Hadoop," International Journal of Engineering Trends and Technology (IJETT), vol.4, issue 8, August 2013.

[8] T. K. Das et al., "BIG Data Analytics: A Framework for Unstructured Data Analysis," International Journal of Engineering and Technology (IJET) vol 5, No 1, Feb-Mar 2013.

[9] Joseph McKendrick, "Big Data, Big Challeneges, Big Opportunities: 2012 IOUG Big Data strategies survey," September 2012.

[10] Jeffrey Dean and Sanjay Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters," Google, Inc

[11] "Hadoop", http://hadoop.apache.org.

[12] Tom White, "Hadoop: The definitive Guide," Third Edition, ISBN: 978-1-449-31152-0-1327616795.

[13] Ian Mitchell, Mark Locke and Aundy Fuller, "The White Book of Big Data"

[14] Http://en.wikipedia.org/wiki/R (programming language).

# Using the Technology Acceptance Model in Understanding Academics' Behavioural Intention to Use Learning Management Systems

Saleh Alharbi

School of ICT, Griffith University
Gold Coast, Australia
Computer Science Department, Shaqra University
Shaqra, Saudi Arabia

Steve Drew

School of ICT, Griffith University
Gold Coast, Australia

*Abstract*—**Although e-learning is in its infancy in Saudi Arabia, most of the public universities in the country show a great interest in the adoption of learning and teaching tools. Determining the significance of a particular tool and predicting the success of implantation is essential prior to its adoption. This paper presents and modifies the technology acceptance model (TAM) in an attempt to assist public universities, particularly in Saudi Arabia, in predicting the behavioural intention to use learning management systems (LMS). This study proposed a theoretical framework that includes the core constructs in TAM: namely, perceived ease of use, perceived usefulness, and attitude toward usage. Additional external variables were also adopted— namely, the lack of LMS availability, prior experience (LMS usage experience), and job relevance. The overall research model suggests that all mentioned variables either directly or indirectly affect the overall behavioural intention to use an LMS. Initial findings suggest the applicability of using TAM to measure the behavioural intention to use an LMS. Further, the results confirm the original TAM's findings.**

*Keywords—Technology Acceptance Model; Higher education; Learning management systems; Saudi Arabia*

## I. INTRODUCTION

The rapid development of information and communications technology (henceforth ICT) makes using ICT imperative. The interest in ICT has drawn substantial research attention[1], and more importantly, ICT contributes directly to the significant changes in teaching and learning that have been occurring in regards to e-learning[2]. The increasing access to ICT creates a new paradigm for education known as e-learning. Therefore, universities around the world have started to revise their strategies in order to adopt technologies that assist in achieving their pedagogical goals. E-learning is commonly defined as the intentional use of ICT in teaching and learning[3].

One of the ICT tools that is incorporated into the education sector is called a learning management system (henceforth LMS). LMS is one of the rapidly-emerging technologies that is widely used in higher education, whether in open-source (Moodle) or commercial LMS such as Blackboard[4]. Paulsen [5] argues that the availability of LMS is considered a critical factor in the success of e-learning. An LMS, alternatively called a learning platform, refers to a wide range of systems that assist teachers and students alike in accessing online learning services [6]. Services provided by an LMS vary from one system to another. However, common services available in an LMS may include access control, performance management, communication facilities, assessments, study schedule documentation, and provision of learning content[7]. Current reports show that more than 95% of all responding universities and colleges in the USA have adopted one or more LMS[8] and that the same adoption rate exists in institutions in the UK[9]. The trend of using LMS in the Middle East is not different. LMS is a promising tool around the globe, including in the Middle East[10]. According to a survey, outcomes about e-learning services provided by 26 Arab universities reveals that 96% adopt LMS as a learning environment to assist in providing blended learning [11]. Blended learning is a term usually used interchangeably with e-learning in the literature involving e-learning in Arab world. In the Gulf Cooperation Council, the education sector has taken care in designing strategic plans to incorporate e-learning[12].

An effective implementation of LMS should highly consider academics who will use such systems for teaching. Therefore, the aim of this research is to develop a theoretical framework based on a well-known technology acceptance model (TAM)[13]. The proposed model contributes to the high volume of research on e-learning in Saudi Arabia, and it will be used to measure academics' behavioural intention for using an LMS. The model is presented in depth in a separate section. The rest of the paper is structured as follows: first, LMS in Saudi Arabia is presented. Then, a brief review on the previous studies of LMS usage in Saudi Arabian higher education is presented, followed by the theoretical framework on which the research model was based. The research context and significance appear next, followed by the methodology section, which provides insight into the research model and hypothesis development. The research methodology section includes a comprehensive structure about the method of validation for the proposed model. The results and discussion are provided prior to the research conclusion and future considerations.

## II. LITERATURE REVIEW

### A. LMS in Saudi Arabia

In Saudi Arabia, although e-learning is in its infancy, most of the Saudi universities keep pace with the development of e-

learning around the world. All governmental universities in Saudi Arabia have a deanship for e-learning and distance learning, created to assist with matching this development and meeting the need to utilise e-learning at universities. The Ministry of Higher Education has initiated an ambitious plan in its establishment of the National Centre for E-learning and Distance Learning (NCeDL). The centre was established to assist in the plan of providing educational tools for local universities[14]. NCeDL contributes to the e-learning industry in the kingdom by providing services and solutions to the local universities. One of the solutions developed locally by the National Centre is an LMS called JUSUR[15]. JUSUR provides academics with features to facilitate their teaching experience, like course and user-management tools, forums, quizzes, and announcements. It also assists in managing the e-learning process by keeping students' data organized, planning courses, making content available to students, tracking students' performance and producing reports about it, facilitating communications with students, and offering testing and assessment tools [16].

JUSUR is hosted and managed by NCeDL. Academics who wish to use the system can register for it by filling out a registration form to create their profiles. Once registered, the NCeDL registration system verifies their academic email. Once verified, the account will be approved and activated. Despite the ease of joining, JUSUR has not been adequately utilised by academics in Saudi Arabia[16]. The highest usage of JUSUR occurs at King Saud University, one of the largest universities in Saudi Arabia. As yet, however, only 55% of the courses at the university are offered through the system. Similar results were reported by [17], who found that the overall utilisation of LMS fell below the satisfactory level. The results of these studies are consistent with another study that showed the harnessing of LMS to accomplish pedagogical benefits in higher education has yet to reach the required level of use[18]. JUSUR LMS has not been the only e-learning system used at Saudi universities. Other commercial LMS such as Blackboard, WebCT, and Design2Learn have been adopted. However, [19] pointed out that only a few faculty members have utilised these systems at each university. Reflecting this issue, a growing number of studies have been conducted on the use of e-learning technologies, whereas research focusing on LMS use receives little attention and remains relatively insufficient. The following section provides the roadmap for these studies.

*B. Previous studies on LMS in Saudi Arabia*

First, it is noteworthy that e-learning is commonly used in place of blended learning in Saudi Arabia. A plethora of studies have examined e-learning in the context of Saudi Arabia; however, little has emerged on LMS usage.

A high percentage of these studies have targeted learner usage of LMS, specifically JUSUR LMS, whereas academics receive only a little attention. Further, most of the studies focus on examining the volume of LMS usage, features used within an LMS, and attitudes towards using such systems. Hence, the previous studies did not target the intentions and behaviours of LMS users. Most importantly, use of the technology acceptance model within LMS in Saudi context is virtually non-existent. Moreover, studies only consider user groups that

have already utilised an LMS. The potential users of LMS, however, are not considered.

Alebaikan and Troudi [20] investigated the use of JUSUR LMS for blended learning in the College of Applied Studies and Community Services at King Saud University. Prior to Alebaikan and Troudi's study [20], an LMS had already been implemented by the faculty to serve the high number of students applying to the college. Their study aimed to interpret students' and academics' perception of a new learning environment with a focus on online discussion features in LMS. From the instructors' point of view, the study concluded that lack of pedagogical and technical experience is an issue in using the Web as a medium of instruction. Further, not all features needed by instructors are available within an LMS. As this study was conducted in one of the largest and most advanced universities in Saudi Arabia, technology integration in teaching within this context could consequently be affected by organizational arrangement[21]. Further, facilitating conditions in which academics would be likely to have more resources and assistance would affect the intention to use the system[22] as they will receive the required support when they need. In addition, Mulkeen [23] suggests that ICT infrastructure should be considered when investigating LMS usage. Finally, it is noted that this study focuses only on online discussion featured within learning management systems that were provided to academics and students prior to the study.

In an attempt to further analyse academics' use of LMS in Saudi Arabia, Asiri, Mahmud, Abu-Bakar, and Ayub [24] suggests a theoretical framework in an attempt to identify factors that influence JUSUR LMS utilisation in public universities in Saudi Arabia. This study is based on the library research approach, and the theoretical framework proposed by the authors was constructed based on well-known theories—namely, the theory of reasoned action[25] and the technology acceptance model[13]. In this study, factors that influence the use of JUSUR LMS are divided into two main categories: internal variables and external variables. First, internal variables consist of three factors that could affect potential users of JUSUR LMS in terms of their attitude, pedagogical beliefs towards e-learning, and level of competency. The authors confirmed that a positive attitude towards JUSUR LMS will likely motivate academics to utilise it. Further, along similar lines with other studies[26, 27], beliefs about e-learning were found important in determining the use of JUSUR LMS. Moreover, the study noted that the use of JUSUR LMS could be predicted by competence level, meaning that having the skills and knowledge to use the system will affect an academic's use of the system. Second, the external variable indicated in this study includes external barriers faced by academics as well as demographic factors. Barriers such as organisational, technological, and social barriers were hypothesised to serve as factors that determine JUSUR LMS usage. Similarly, demographical factors such as gender, computer self-efficacy, and training are also used to predict JUSUR LMS usage.

In a different study, Asiri, Mahmud, Abu-Bakar, and Ayub [28] studied faculty members' utilisation of JUSUR LMS at three public universities in Saudi Arabia and their attitude towards such utilisation. Like the previously-mentioned study,

this study targeted academics who have already utilised LMS to assist them in teaching. The study aimed to determine the volume of JUSUR LMS utilisation that constitutes a moderate level. It is noteworthy that, according to the study, the moderate level is explained as the use of LMS for less than one hour on average twice a month. However, the finding of this study is not consistent with that of other studies mentioned earlier, wherein LMS usage is believed to be below the satisfactory level. Nevertheless, the study confirms that faculty members have a positive attitude towards JUSUR LMS.

In the same way, Hussein [16] studied the attitude of faculty members in Saudi universities towards JUSUR LMS. Similar to other studies, academics in these universities had developed a sufficient awareness and positive attitude towards JUSUR LMS. Despite that, the study confirmed their low level of JUSUR LMS usage, which was not justified within the study.

Similar to the study above, Albalawi and Badawi [29] conducted a study targeting faculty members of the University of Tabuk, which is a public Saudi Arabian university. The study aimed to highlight academics' perception and awareness of e-learning. Surprisingly, the study revealed that almost 63% of faculty members had a negative perception of e-learning. It is worth mentioning that this study was conducted prior to any implementation of e-learning technologies at Tabuk University, making the situation similar to the current study in terms of the absence of LMS.

Other research exists on acceptance of e-learning in general, with a focus on LMS systems. However, this research has used students as subjects [30-33], and students are outside of the current study's research scope.

The main limitation of the previous studies, however, is that they mostly focus on measuring the attitude of faculty members towards already-implemented LMS systems. In other words, most of the existing research focuses on users who have already used an LMS in teaching. Therefore, intention to use an LMS by those who have not used one is not considered. Moreover, although higher education providers in Saudi Arabia are implementing LMS, little has been done to examine the factors that influence academics to use an LMS. Further, the previous studies limited their scope to an examination of the use of JUSUR LMS in Saudi Arabia. However, as stated earlier, JUSUR LMS is not the only LMS employed in Saudi Arabian public universities[11]. In this study, an LMS is defined as any LMS that is either centrally-managed and government-run or privately adopted in a public university.

In response to this gap in literature, this paper develops a research model based on the technology acceptance model (TAM). The following section presents the theoretical framework of the study.

### III. THEORETICAL FRAMEWORK

From the stream of research on information systems (IS), many theories have been proposed to explain the relationship between determinants that would affect technology acceptance. The most common factors are user attitudes, perceptions, beliefs, and actual system use. Frameworks such as the theory of planned behaviour (TPB)[25], diffusion of innovation[34],

the unified theory of acceptance and use of technology (UTAUT)[22, 35], the DeLone and McLean model of IS success[36], and measurement and analysis of computer user satisfaction[37, 38] are popular models used in the context of technology acceptance. Most of these models, however, focus on only technical factors[4].

The technology acceptance model (TAM)[13] is possibly the most widely-used framework in the field of IS for measuring technology acceptance[4, 39-41], and its high validity has been proven empirically in many previous studies[42-44]. Further, Al-Gahtani [45] confirms the validity and reliability of TAM constructs to predict IS adoption in Arab culture, specifically in the Saudi culture. In relation to e-learning and LMS, TAM has also been adopted and tested[46, 47]. Although TAM is a well-known and tested theory in the field of IS, using TAM in predicting and explaining LMS usage has so far received little attention[48].



Fig. 1. The technology acceptance model[13]

TAM was first introduced by Davis [49] around the concept of technology acceptance. As depicted in Figure 1, TAM posits that acceptance of a new IS can be predicted based on users' behavioural intention (BI), attitude towards use (A), and two other internal beliefs: perceived usefulness (U) and perceived ease of use (E). Davis[13] defined perceived usefulness as "the prospective user's subjective probability that using a specific application system will increase his or her job performance within an organizational context" (p. 985) and perceived ease of use as "the degree to which the prospective user expects the target system to be free of effort" (p. 985).

According to TAM, behavioural intention (BI) defines the actual use of a given IS system and therefore determines technology acceptance. Attitude towards use (A) and perceived usefulness (U) jointly influence BI (A). BI is also indirectly affected by perceived ease of use (E). A is directly affected by both U and E, while U is directly influenced by E. Further, TAM theorizes that perceived usefulness and perceived ease of use are affected by external variables. Thus, U and E mediate the effect of external variables on user's attitude and behavioural intention, and therefore the actual system use.

#### A. Shaqra University

Shaqra University is a public university established in 2008, located in Shaqra, Saudi Arabia. In addition to the main campus, there are eight other campuses in geographically

distributed locations that include a total of twenty-one colleges and approximately twenty departments. The latest figures[50] show a total of 761 faculty members and 10,767 enrolled students. Table 1 provides information about the different campuses, colleges, departments, and faculty members' ranks.

TABLE I. Shaqra University Demographics

| Faculty Member Statistics | | | | | |
|---|---|---|---|---|---|
| *Professors* | *Associate Professors* | *Assistant Professors* | *Lecturers* | *Instructors* | *Total* |
| 15 | 35 | 200 | 281 | 230 | 761 |
| Campuses and Faculties | | | | | |
| *Total campuses* | *Faculty/ Colleges* | | | *Departments* | |
| 9 | 21 | | | 20 | |

In line with the national strategic plan, Shaqra University shows interest in incorporating ICT into learning and teaching practices. The university regularly participates in conferences hosted by the ministry of higher education in Saudi Arabia. In addition, soon after the establishment of the university, the deanship of information technology and e-learning was also established. The aim of this initiative is to provide both academics and students with pedagogical and technical support. Moreover, different workshops have been held to raise faculty members' awareness of e-learning. As yet, however, face-to-face teaching is the official medium of instruction at the university.

*B. Study significance*

The significance of the current study stems from various considerations. First, no previous research has sought to investigate faculty members' behavioural intention to use LMS and empirically validate the technology acceptance model at Shaqra University. Moreover, the findings of this study will provide the university with more insight into academics' perception of LMS. Further, this study will pave the way for future research on technology acceptance within the higher education context in Saudi Arabia. Specifically, this study adopted and modified a questionnaire to suit the LMS acceptance context that may be reused in future research.

## IV. RESEARCH MODEL AND HYPOTHESES

The research model is applied to two different groups: academic users and academic non-users. Those in the user group are examined based on their current use of an LMS or their previous experience of usage. Due to their potential to use an LMS, a non-user group is also examined. According to Taylor and Todd [51], TAM has successfully predicted and explained almost equal behavioural intention to adopt a new technology among inexperienced and experienced users. Further Shih [52] noted that TAM can be applied prior to the adoption of a new technology.

In accordance with the research objective and consistent with the related literature, the research model, as shown in Fig. 2, consists of the TAM core constructs and three key moderators. The following section discusses the development of relevant hypotheses.

*A. Hypotheses in relation to TAM variables*

As previously discussed, TAM proposed the following relationship between its constructs: a) Intention to use is positively affected by attitude toward using and perceived usefulness; b) Attitude toward using is positively affected by perceived usefulness and perceived ease of use; and c) perceived usefulness is directly affected by perceived ease of use. In this study, perceived usefulness is defined as the degree to which a faculty member believes that using an LMS would enhance his or her job performance, while perceived ease of use is defined as the degree to which a faculty member believes that learning to use an LMS requires a relatively low degree of effort. The linkage between the different variables has been proven by different studies on e-learning and LMS usage [31, 53-56]. Therefore, the relationships between perceived ease of use, perceived usefulness, attitude toward using, and intention to use an LMS system are hypothesised as the following:

*1) Perceived ease of use positively affects perceived usefulness of an LMS.*

*2) Perceived ease of use positively affects attitudes towards using an LMS.*

*3) Perceived usefulness positively affects attitudes towards using an LMS.*

*4) Perceived usefulness positively affects intention to use LMS.*

*5) Attitude towards using positively affects intention to use LMS.*

Ong, et al. [57] highlights that intention to use e-learning is also effected by perceived ease of use. Therefore, the relationship between perceived ease of use and behavioural intention for use is hypothesised as:

*6) Perceived ease of use positively affects intention to use an LMS.*

*B. Hypotheses in relation to external factors and TAM variables*

The ease of use and usefulness constructs may not be sufficient, and therefore other variables may be needed[58]. Thus, after reviewing the relevant studies, this study suggests three external variables: LMS usage experience, job relevance, and lack of LMS availability. As in figure 2, researchers believe that the suggested external variables moderate the original TAM variables. The following explains the hypotheses on the relationship between external moderators and TAM variables.

Venkatesh and Davis [59]found that experience using technology serves as a critical factor in determining technology acceptance. Thompson, et al. [60] defines usage experience as individual involvement in or exposure to a particular system and the accumulative skills the user gains by using the system. In this study, LMS usage is suggested to moderate TAM variables. LMS usage is defined as academics' previous or current use of an LMS as a medium of instruction within an e-learning environment. Therefore, the following is hypothesised:

*7) LMS usage experience negatively influences the non-user group's intention to use an LMS.*

*8) LMS usage experience negatively influences the non-user group's perceived ease of use of an LMS.*

*9) LMS usage experience negatively influences the non-user group's perceived usefulness of an LMS.*

TAM was extended to incorporate job relevance as a factor that directly affects perceived usefulness [59]. According to Venkatesh and Davis [59], job relevance is "an individual's perception regarding the degree to which the target system is applicable to his or her job" (p.191). Similarly, this study proposes that job relevance affects both perceived ease of use and perceived usefulness. Job relevance in this study is defined as an academic's perception regarding the degree to which an LMS system is relevant to use in managing learning activities at Shaqra University. As found by Venkatesh and Davis [59], job relevance is believed to positively exert a direct effect on perceived usefulness. Consequently, this study argues that job relevance also affects perceived ease of use (PEOU). Therefore, the following are the hypotheses of this study on the relationship between job relevance and TAM variables:

*10) Job relevance positively affects the perceived usefulness of an LMS.*

*11) Job relevance positively affects the perceived ease of use of an LMS.*

### C. Hypotheses on the relationship between lack of LMS availability and TAM variables

At the present, Shaqra University provides no LMS to faculty members. Therefore, as mentioned earlier, academics at Shaqra University can be categorised into the following groups in relation to LMS usage: a) experienced members who have used, and/or are using an LMS in their teaching, or b) inexperienced members who have not utilised an LMS yet. For both categories, the study proposes that a lack of LMS availability has a negative impact on perceived ease of use (E).

Therefore, the relationship between lack of LMS availability and perceived ease of use is hypothesised as following:

*1) H12    Lack of LMS availability negatively affects the perceived ease of use of an LMS.*

### V.    RESEARCH METHOD

The study is quantitative in nature and employs an online survey for data collection. Online surveys provide researchers with various benefits[61], saving researchers time and expenses by overcoming geographic distance. Moreover, they assist in accessing unique subjects. Due to Saudi Arabia's gender-segregated higher education system, the online survey was the appropriate tool to use in order to access both male and female participants. The online survey was developed to examine the relationship between variables proposed in the research model.

### A. Questionnaire

To ensure content validity, the questionnaire used in this study was adapted from the original measurement scales used in TAM[13] and from other literatures[31, 55, 59, 62] with some modifications and the necessary wording changes and validation to fit the context of LMS usage. To avoid issues that can occur in wordings, measurement and ambiguities, the questionnaire was pre-tested by two native English speakers. Sekaran and Bougie [63] highlight that such pre-test is essential because wording problems significantly influence accuracy[64]. The questioner was also translated into Arabic because most of the academics at Shaqra University are native Arabic speakers. For the Arabic version, the back translation method suggested by [65] is used. This method suggests that the questionnaire measurements should be translated by bilingual experts back and forth from the source language to the targeted language. Based on that concept, the English version was sent to two bilingual experts to translate it into Arabic, and the back translation method was followed until the English and Arabic version converged. Finally, the Arabic version was also revised by an expert in the Arabic language for clarity.



Fig. 2.    Research Model

## B. Participants

The participants in this study were 59 faculty members from different colleges and different departments who voluntarily participated in the online survey. All participants in this study were academics working for Shaqra University, who fit well with the aim and context of this study.

## C. Sampling technique

While it is difficult to get responses from a whole population, sampling is an attempt to draw a conclusion based on a small representation in a given population[66]. The sample in this survey is considered a subset of academics at Shaqra University, comprised of some faculty members selected from the institution. The sampling technique used in the present study is non-probability convenience sampling. Convenience sampling is found used in many studies investigating technology acceptance. Further, the technique is used to ensure a better response rate in a short amount of time. Finally, it is against Shaqra University privacy policies to obtain academics' contacts and email addresses from the faculty or its departments. Additionally, using the university's mailing list may have led to the inclusion of other participants who are out of this study scope and therefore distort the findings. Hence, convenience sampling was the optimal technique for the purpose of this study.

## D. Instrumentation

The research instrument consists of two main sections. The first section incorporates a nominal scale to identify respondents' demographic information. The second section uses 7-point Likert response scale where 7: Strongly disagree, 6: Moderately disagree, 5: Slightly disagree, 4: Neutral, 3: Slightly agree, 2: Moderately agree, and 1: Strongly agree. This section includes TAM constructs.

## E. Demographic characteristics

This part of the questioner identifies respondents' basic demographic characteristics. It contains 10 items such as gender, age, teaching experience, academic rank and administration position, academic field, faculty and departmental information, and previous experience with LMS (Table 2).

## F. Measuring TAM constructs

The second section of the survey(Table 3), as discussed in the questionnaire design above, measures TAM constructs used in this study. As shown in table 3, there are 20 items measured in accordance with the current study's research model. The measured items include perceived ease of use (7 items), perceived usefulness (6 items), attitude toward usage (3 items), behavioural intention to use (2 items), and job relevance as an external factor (2 items). It is noteworthy to mention that the seven items used to measure perceived ease of use include one item—lack of LMS availability—that is hypothesised in this research to moderate perceived ease of use.

TABLE II.    Questionnaire – Section I

| Section I: Demographic Characteristics Information |
|---|
| - Gender: |
|   1.   Male |
|   2.   Female |
| - Age |
|   1.   Less than 25 |
|   2.   25-30 |
|   3.   30-40 |
|   4.   40-50 |
|   5.   Above 50 years old |
| - Experience in higher Education(In general, not only at Shaqra University) |
|   1.   Less than 1 year |
|   2.   More than 1 year and less than 3 years |
|   3.   More than 3 years and less than 5 years |
|   4.   More than 5 year and less than 10 years |
|   5.   More than 10 years |
| - Experience at Shaqra University |
|   1.   Less than 1 year |
|   2.   More than 1 year and less than 2 years |
|   3.   More than 2 years and less than 5 years |
| - Academic Rank |
|   1.   Professor |
|   2.   Associate Professor |
|   3.   Assistance Professor |
|   4.   Lecturer |
|   5.   Instructor |
| - Your Academic administrator position |
|   1.   Vice-rector or deputy vice-chancellor |
|   2.   Dean |
|   3.   Associate Dean |
|   4.   Department chairman |
|   5.   Centre director |
|   6.   None |
| - Your academic field |
|   1.   Humanities & Social Sciences |
|   2.   Natural Sciences |
|   3.   Applied Sciences( e.g. engineering, computing& IT) |
|   4.   Medical & Health Sciences |
| - What is your Faculty? |
| - What is your department? |
| - How long have you used, or have been using a Learning |
| - Management System (LMS)? |
|   1.   Have not used a System Management System |
|   2.   Less than a year |
|   3.   1-3 years |
|   4.   3-5 years |
|   5.   More than 5 years |

## G. Data collection

The questionnaire was made available at the beginning of academic year 2013/2014. The survey was distributed online by emailing a convenient sampling of 105 academics with the URL to the survey. Participants had the option to switch between English and Arabic at any time during the survey. At this time, of the 105 questionnaires distributed, 69 responses were recorded (65.71%). Of that, only 59 responses yielded valid responses that were used for analysis. The overall response rate was 56.19%.

TABLE III.　　Questionnaire – Section II

| Section II: Perceived Ease of Use (PEU) | |
| --- | --- |
| I feel that using an LMS would be easy for me | PEU1 |
| I feel that my interaction with LMS would be clear and understandable | PEU2 |
| I feel that it would be easy to become skilful at using LMs | PEU3 |
| I would find LMS to be flexible to interact with | PEU4 |
| Learning to operate LMS would be easy for me | PEU5 |
| it would be easy for me to get LMS to do what I want to do | PEU6 |
| I feel that my ability to determine LMS ease of use is limited by my lack of experience | PEU7 |
| **Section III: Perceived Usefulness (PU)** | |
| Using LMS in my job would enable me to accomplish tasks more quickly | PU1 |
| Using LMS would improve my job performance. | PU2 |
| using LMS in my job would increase my productivity | PU3 |
| Using LMS would enhance my effectiveness on the job. | PU4 |
| Using LMS would make it easier to do my job | PU5 |
| I would find LMS useful in my job | PU6 |
| **Section IV: Attitude Toward Usage (ATU)** | |
| I believe it is a good idea to use a Learning Management System | ATU1 |
| I like the idea of using a Learning Management System | ATU2 |
| Using a Learning Management System is a positive idea | ATU3 |
| **Section V: Behavioural Intention to Use (BIU)** | |
| I plan to use a learning Management System in the future | BIU1 |
| Assuming that I have access to an LMS, I intend to use it | BIU2 |
| **Section IV: Job Relevance (BIU)** | |
| In my job, the usage of a learning Management System is important | JR1 |
| In my job, the usage of a learning Management System is relevant | JR2 |

*H. Ethical issues*

Ethical clearance was obtained prior to the study. Participation in this study was voluntary and data was collected anonymously. This study did not involve personal information about subjects. Prior to commencing the survey in this study, all participants were made aware of the research significance and type of information being collected. The researchers explained that the participation in this research is based on subjects' interest, that they are under no obligation to participate, and that they may decline to participate at any time. Their right to withdraw at any time during the survey was explicitly stated. Further, data confidently was assured. Data collected from this research is to be kept confidential.

*I. Procedure*

In accordance with the hypotheses of the present study, a presentation was designed about LMS. The presentation consists of three main parts. The first provides a general overview about e-learning and presents its basic definition. The second part introduced LMS as a tool used in e-learning. The second part, in turn, includes basic definitions of LMS, describes its main features and components, and provides examples of various examples of LMS. The third part encloses a video presenting a local university experience with e-learning. The presentation was made available in English as well as Arabic. The presentation was designed using Prezi, which is a cloud-based presentation software program that presents ideas in a non-liner manner[67]. It incorporates map-like features to allow users to highlight the concepts that are most important to them and the relationships between these concepts while providing ease of navigation.

As mention earlier, the present study assumes that participants are mainly experienced (referred to as a user-group) and inexperienced (the non-user group) with regard to LMS usage. Within the questionnaire structure described above, all participants had access to the demographic section, which is presented first. At the end of this section, participants were asked whether or not they had previous experience with LMS. Based on their response to this question, participants were directed to the appropriate following section. If they answered that they had no previous experience and proceeded, the following section embedded the presentation described earlier. After the presentation concluded, the non-user could proceed to Section 2, about the perceived ease of use. On the other hand, any other response to the question about LMS usage experience led respondents to skip the presentation and proceed directly to Section 2. The presentation was provided to support this study's hypothesis that there would be no significant differences between the user and non-user groups in their intention to adopt LMS in teaching.

VI.　　DATA ANALYSIS AND RESULTS

*A. Demographics*

The participants were almost equal in terms of gender, with 28 (47.46%) males and 31 (52.54%) females. The majority of participants were between 25 and 40 years, with 28.81% from 25 to 30, 37.29% from 31 to 40, 22.03% from 41 to 50, and 10.17% above 50, with a low minority (1.69%) below 25. Saudi-nationality academics recorded the highest response rate, at 54.24%. The rest of the figures and information are presented in Table 6.

*B. Experience with LMS*

The current study, as previously discussed, investigates the applicability of using TAM on two groups: a user group and a non-user group. As expected, consistent with other figures, almost half of the respondents had not used an LMS (49.15%), while the rest vary in their experience with LMS as follows: Those who had used LMS for less than a year stood at 16.95%; and almost double this number, 28.81%, had used LMS for more than a year but less than 3 years. Only a few respondents had used an LMS for more than three years (3.39% for more than three years and only 1.69% for more than 5 years). In general, the results reveal that 49.15% of respondents were non-users and 50.85% were users. This distribution enables discrimination between the two groups' responses within the study context (Table 4).

TABLE IV.     Academics experience with LMS

| Respondents | | Frequency | Percentage |
|---|---|---|---|
| Have not used a learning management system | | 29 | 49.15% |
| experienced users | | 30 | 50.15% |
| Experience in years | Less than a year | 10 | 16.95% |
| | 1–3 years | 17 | 28.81% |
| | 3–5 years | 2 | 3.39% |
| | More than 5 years | 1 | 1.69% |

## C. Validity and reliability

In addition to the steps mentioned earlier to assess instruments' validity and reliability, a further test was performed. Reliability assessment was done using Cornbach Alpha[68]. Reliability concerns internal consistency between multiple measurements of variables, and Cornbach Alpha is commonly used to measure it[69].

TABLE V.     Instruments reliability Cornbach Alpha

| Scale | Number of Items | Cronbach Alpha |
|---|---|---|
| Perceived ease of use (PEU) | 7 | 0.901 |
| Perceived usefulness (PU) | 6 | 0.924 |
| Attitude towards use (ATU) | 3 | 0.916 |
| Behavioural intention to use (BIU) | 2 | 0.801 |
| Job relevance (JR) | 2 | 0.924 |
| Overall reliability | 20 | 0.958 |

As per many studies (i.e.,[70, 71], constructs are considered to have internal consistency reliability when the Cronbach Alpha value exceeds 0.07.

In this study, the reliability assessment was done using Statistical Package for Social Sciences (SPSS) version 21. All measures in this study show a high level of reliability, ranging from 0.901 to 0.924, with a satisfactory value of 0.801 for behavioural intention to use. All scales exceeded 0.70, and therefore the survey is considered reliable.

## D. Statistical analysis and hypotheses testing

In line with the study objective, correlation analysis was conducted to examine the relationship between the variables used within this study, and therefore to empirically decide whether or not to accept or reject the null hypotheses. The structure of hypothesis testing is as follows. First, hypotheses were tested based on the size of the whole sample. Second, the study investigated the role of prior experience, and therefore hypotheses were tested on the non-user group. Finally, the user-group sample was used for testing the hypothesis. The aim is to provide a comprehensive correlation analysis and then investigate the impact role of prior experience on the correlation significance.

TABLE VI.     Respondents' demographic information

| Respondents | Frequency | Percentage |
|---|---|---|
| *Experience in Higher Education\* in years* | | |
| < 1 | 2 | 3.39 |
| > 2 < 3 | 19 | 32.20 |
| > 3 < 5 | 15 | 25.42 |
| > 5 < 10 | 7 | 11.86 |
| > 10 | 16 | 27.12 |
| *Experience at Shaqra University\* in years* | | |
| < 1 | 6 | 10.17 |
| > 2 < 3 | 22 | 37.29 |
| > 3 | 31 | 52.54 |
| *Academic Rank* | | |
| Associate Professor | 10 | 16.95 |
| Assistant Professor | 17 | 28.81 |
| Lecturer | 16 | 27.12 |
| Instructor | 16 | 27.12 |
| *Administrative Work* | | |
| Vice-rector or deputy vice-chancellor | 1 | 1.69 |
| Associate dean | 6 | 10.17 |
| Department chairman | 11 | 18.64 |
| Centre director | 1 | 1.69 |
| None | 31 | 52.54 |
| Studying in KSA | 2 | 3.39 |
| Studying abroad | 7 | 11.86 |
| *Departments* | | |
| Computer Sciences & Information Systems | 19 | 32.20 |
| Business Administration & Finance Management | 1 | 1.69 |
| English | 10 | 16.95 |
| Physics | 2 | 3.39 |
| Mathematics | 5 | 8.47 |
| Chemistry | 5 | 8.47 |
| Biology | 4 | 6.78 |
| Pharmacology-related | 2 | 3.39 |
| Islamic Studies | 2 | 3.39 |
| Arabic | 5 | 8.47 |
| Home Economics | 2 | 3.39 |
| Other | 2 | 3.39 |
| *Faculty* | | |
| Community College | 2 | 3.39 |
| Education Faculty | 10 | 16.95 |
| Arts & Sciences College | 35 | 59.32 |
| Applied Medical Sciences | 2 | 3.39 |
| College of Sciences & Humanities | 9 | 15.25 |
| Faculty of Pharmacy | 1 | 1.69 |

*E. Hypotheses testing for all participants*

As stated, this section used the sample size as a whole to test the research hypotheses. Hypotheses on the relationship between TAM original variables are presented first.

*1) Hypotheses for TAM variables*

*a) Perceived ease of use positively affects perceived usefulness of an LMS.*

From the correlation analysis result in Table 7, it can be observed that there is a significant positive relationship between the perceived ease of use and perceived usefulness of an LMS. Therefore, H1 is supported.

TABLE VII.    PEOU and PU correlations

| Correlations | | |
|---|---|---|
| *Factors* | | *PU* |
| PEOU | r-value | .576** |
| | p-value | .000 |
| | N | 59 |
| *PEOU: Perceived ease of use; PU: Perceived usefulness* | | |

*b) Perceived ease of use positively affects attitudes towards using an LMS.*

From the correlation analysis result in Table 8, it can be observed that there is a significant positive relationship between the perceived ease of use and attitude towards usage. Therefore, H2 is supported.

TABLE VIII.    PEOU and ATU correlations

| Correlations | | |
|---|---|---|
| *Factors* | | *ATU* |
| | r-value | .513** |
| PEOU | p-value | .000 |
| | N | 59 |
| *PEOU: Perceived ease of use; ATU: Attitude towards usage* | | |

*c) Perceived usefulness positively affects attitudes towards using an LMS.*

From the correlation analysis result in Table 9, it can be observed that there is a significant positive relationship between the perceived usefulness and attitude towards usage. In fact, the relationship between perceived usefulness and attitude towards usage indicates a stronger relationship than the relationship between perceived ease of use and attitude towards usage. In general, H3 is supported.

Table 9.  PU and ATU correlations
TABLE IX.    PU and ATU correlations

| Correlations | | |
|---|---|---|
| *Factors* | | *ATU* |
| | r-value | .691** |
| PU | p-value | .000 |
| | N | 59 |
| *PEOU: Perceived ease of use; ATU: Attitude towards usage* | | |

*d) Perceived usefulness positively affects intention to use LMS.*

From the correlation analysis result in Table 10, it can be observed that there is a significant positive relationship between the perceived usefulness and behavioural intention to use an LMS. Surprisingly, the relationship between perceived usefulness and behavioural intention to use does indicate a strong correlation. However, H4 is supported.

TABLE X.    PEOU and BIU correlations

| Correlations | | |
|---|---|---|
| *Factors* | | *BIU* |
| | r-value | .481** |
| PU | p-value | .000 |
| | N | 59 |
| *PEOU: Perceived ease of use; BIU: Behavioural intention to use* | | |

*e) Attitude towards using positively affects behavioural intention to use an LMS.*

From the correlation analysis result in Table 11, it can be observed that there is a positive relationship between the attitude towards usage and behavioural intention to use an LMS.

However, the relationship does not seem to be significant, and the correlation is not strong. Statistically, H5 is supported.

TABLE XI.    ATU and BIU correlations

| Correlations | | |
|---|---|---|
| *Factors* | | *BIU* |
| | r-value | .265* |
| ATU | p-value | .043 |
| | N | 59 |
| *ATU: Attitude towards usage; BIU: Behavioural intention to use* | | |

*f) Perceived ease of use positively affects intention to use an LMS.*

From the correlation analysis result in Table 12, a significant positive relationship between perceived ease of use and behavioural intention to use can be observed. Therefore, H6 is supported.

TABLE XII.    ATU and BIU correlations

| Correlations | | |
|---|---|---|
| *Factors* | | *BIU* |
| | r-value | .376** |
| PEOU | p-value | .003 |
| | N | 59 |
| *PEOU: Perceived ease of use; BIU: Behavioural intention to use* | | |

*2) The role of prior experience hypotheses*

The current study introduced LMS usage experience as a new moderator believed to affect the original TAM constructs. The users were categorized into two groups based on their previous experience using an LMS. Users who had used, or had been using, an LMS were called the user group. The inexperienced users, named the non-user group, were those who had not utilised an LMS before. The related hypotheses were:

*a) LMS usage experience negatively influences the non-user group's intention to use an LMS.*

*b) LMS usage experience negatively influences the non-user group's perceived ease of use of an LMS.*

*c) LMS usage experience negatively influences the non-user group's perceived usefulness of an LMS.*

TABLE XIII.     The role of prior experience correlations

| Correlations | | | PU | BIU | JR |
|---|---|---|---|---|---|
| *The role of prior experience* | | | **PU** | **BIU** | **JR** |
| **Non-user group** | *PEOU* | r-value | .606** | .476** | .665** |
| | | p-value | .001 | .009 | .001 |
| | | N | 29 | 29 | 29 |
| | *PU* | r-value | | .670** | .863** |
| | | p-value | | .001 | .001 |
| | | N | | 29 | 29 |
| | *BIU* | r-value | | | .720** |
| | | p-value | | | .001 |
| | | N | | | 29 |
| **User group** | *PEU* | r-value | .492** | .247 | .495** |
| | | p-value | .006 | .188 | .005 |
| | | N | 30 | 30 | 30 |
| | *PU* | r-value | | .226 | .627** |
| | | p-value | | .229 | .001 |
| | | N | | 30 | 30 |
| | *BIU* | r-value | | | .172 |
| | | p-value | | | .363 |
| | | N | | | 30 |

As shown in Table 13, there is a positive correlation between TAM variables for both groups, and it is statistically significant in most cases. Interestingly, a stronger correlation between TAM variables occurred for the non-user group. For instance, when comparing the effect of perceived ease of use on behavioural intention to use an LMS, the non-user group showed a significantly stronger and positive correlation. However, with the user group, the correlation was not statically significant. Consequently, it can be concluded that H7 is not supported.

The correlation between the two main constructs—perceived ease of use and perceived usefulness—indicates a significant positive relationship with other variables. In fact, the non-user group shows a higher correlation. In general, both H8 and H9 are not supported.

*3)   The role of the job-relevance hypotheses*
Job relevance was hypothesised to have a positive impact on both perceived ease of use and perceived usefulness, as follows.

*a) Job relevance positively affects the perceived usefulness of an LMS.*

Job relevance correlates strongly with perceived usefulness. Further, there is a significant positive relationship between the two variables. Therefore, H10 is supported (Table 14).

TABLE XIV.      JR and PU correlation

| Correlations | | |
|---|---|---|
| *Factors* | | *PU* |
| JR | r-value | .769 |
| | p-value | .001 |
| | N | 59 |

*JR: Job relevance; PU: Perceived usefulness*

*b) Job relevance positively affects the perceived ease of use of an LMS.*

Job relevance correlates moderately with perceived ease of use. Similar to the correlation for H10, there is a significant positive relationship between the two variables. Therefore, H11 is supported (Table 15).

TABLE XV.       JR and PEOU correlation

| Correlations | | |
|---|---|---|
| *Factors* | | *PEOU* |
| JR | r-value | .592 |
| | p-value | .001 |
| | N | 59 |

*JR: Job relevance; PU: Perceived usefulness*

*4)   The role of the lack of LMS availability hypotheses*
It is hypothesized that the lack of LMS availability would positively affect the perceived ease of use as follows:

*a) Lack of LMS availability positively affects the perceived ease of use.*

Surprisingly, the result shows a negative relationship between lack of LMS availability and perceived ease of use. This result indicates that all participants (in both the user and non-user groups) perceived that lack of LMS availability does not affect ease of use. Therefore, H12 is not supported.

TABLE XVI.      The Lack of LMS availability with PEOU

| Correlations | | |
|---|---|---|
| *Factors* | | *PEOU* |
| Lack of LMS availability | r-value | -.294 |
| | p-value | .024 |
| | N | 59 |

*PEOU: Perceived ease of use*

The table below summarises the hypothesis after the testing was done.

TABLE XVII.     Hypothesis summary

| Hypothesis | Statement | Result |
|---|---|---|
| H1 | Perceived ease of use positively affects perceived usefulness of an LMS. | Supported |
| H2 | Perceived ease of use positively affects attitudes towards using an LMS. | Supported |
| H3 | Perceived usefulness positively affects attitudes towards using an LMS. | Supported |
| H4 | Perceived usefulness positively affects intention to use an LMS. | Supported |
| H5 | Attitude towards using positively affects intention to use an LMS. | Supported |
| H6 | Perceived ease of use positively affects intention to use an LMS. | Supported |
| H7 | LMS usage experience negatively influences the non-user group's intention to use an LMS. | Not supported |
| H8 | LMS usage experience negatively influences the non-user group's perceived ease of use of an LMS. | Not supported |
| H9 | LMS usage experience negatively influences the non-user group's perceived usefulness of an LMS. | Not supported |
| H10 | Job relevance positively affects the perceived usefulness of an LMS. | supported |
| H11 | Job relevance positively affects the perceived ease of use of an LMS. | supported |
| H12 | Lack of LMS availability negatively affects the perceived ease of use of an LMS. | Not supported |

## VII.   Discussion

The current study modified TAM mainly to validate the relationship between the TAM core constructs as well as the effects of moderators proposed with this study. Overall, the statistical analysis shows that the findings of the current study are consistent with the original TAM findings[13]. All TAM-related hypotheses within this study were proven to have positive correlations that are statistically significant. In line with other studies[18, 24, 72-75], academics involved in this study showed a positive attitude towards LMS, and they intent to use an LMS in their work. Further, the study indicates that when users' perceived ease of use increases, the perceived usefulness increases accordingly. As expected, when academics perceived LMS as easy to use, they developed a positive attitude towards utilising it. Similarly, the perceived usefulness increased the degree of positivity toward usage, which subsequently affected the behavioural intention to use. Interestingly, the findings vary between the user-group and the non-user group. The results show that the non-user group shows higher intention towards using LMS. While both groups perceived LMS as easy to use, the statistics show that the non-user group perceived LMS as more useful than the other group.

Job relevance, adapted from[22], showed a strong relationship with perceived usefulness. Academics believed that the use of LMS for teaching is relevant to their job and is a useful matter. On the other hand, lack of LMS availability did not affect academics' perceived ease of use. Based on the presentation provided to them, they believed that LMS would be easy to use.

Gender and academic rank did not correlate significantly with other variables. The findings suggest that gender and academic rank do not reflect a significant correlation with other constructs. However, the effect of age and gender does not fit within the scope of this study objective.

## VIII.   Conclusion and implications

In general, this study modified the original TAM in order to measure academics' behavioural intention to use an LMS. The current study adapts the core constructs used in TAM. Specifically, it validates the relationship between perceived ease of use, perceived usefulness, attitude towards usage, and overall impact on behavioural intention to use. No surprising findings were found regarding the previous constructs. Therefore, this study confirms other empirical evidence and findings based on TAM. Further, the study successfully confirms the applicability of TAM in the Arab world, specifically in Saudi Arabia in higher-education settings.

As suggested by TAM, this study incorporates external variables including lack of LMS availability, job relevance, and experience with LMS usage. First, the unique environment in which data is collected influences the theoretical framework for this study. Subjects' lack of access to LMS during the data collection phase was assumed to exert a moderating effect on the relationship between TAM constructs, specifically the ease of use. However, the findings show that lack of LMS availability does not automatically mean academics believe using an LMS is difficult. The other external variable, job relevance, was also proven to have a strong relationship with TAM constructs. In particular, job relevance within the context of this study positively affected academics' perceived usefulness of an LMS. The role of prior experience with LMS usage was also investigated. The overall results for both experienced and inexperienced users confirm the original TAM findings. Within this study, inexperienced users indicated a higher degree of positivity towards LMS adoption.

The implication of this study can be surmised as follows. First, this study proposed a theatrical framework based on a robust acceptance model (TAM). This framework can be used to predict the behavioural intention to use an IS prior to the actual implementation. Further, the research model is validated with two different groups in the higher-education context. Moreover, this study contributes to the efforts to empirically validate TAM in the Arab world. Most significantly, this study could benefit Shaqra University's management staff in their future plans to adopt e-learning technologies.

### A.  Research limitations and future work

This study is not free of limitations. First and most importantly, this study was limited by time. Further, the findings of this study may not be greatly generalised for various reasons. First, the researchers conducted a size power test, and the result suggests that the sample size should be increased, as a higher sample size would help to make the conclusion more general. In addition, the research framework was designed to be used with LMS. Moreover, the focus on individuals was the main theme of this study. Future studies

could focus on general ICT adoption for teaching and learning. Additionally, collecting data from different groups could be affected by the increase of usage and experience of users[22]. Therefore, longitudinal research may be more suitable to better predicting attitude and behaviour, and hence facilitating comprehensive understanding of the relationships between variables. Moreover, other statistical tests such as factor analysis, multiple regressions, and structural equation modelling could be conducted to confirm variables' validity.

### REFERENCES

[1] W. H. Dutton and B. D. Loader, Digital academe: new media in higher education and learning: Routledge, 2004.

[2] [D. Radcliffe, "Technological and pedagogical convergence between work-based and campus-based learning," *Educational Technology & Society,* vol. 5, pp. p54-59, 2002.

[3] S. Naidu, *E-learning: A guidebook of principles, procedures and practices*: Commonwealth Educational Media Centre for Asia (CEMCA), 2003.

[4] K. A. Al-Busaidi and H. Al-Shihi, "Instructors' Acceptance of Learning Management Systems: A Theoretical Framework," *Communications of the IBIMA,* vol. 2010, p. 2010, 2010.

[5] M. F. Paulsen, "Experiences with Learning Management Systems in 113 European Institutions," *Educational Technology & Society,* vol. 6, pp. 134-148, 2003.

[6] M. F. Paulsen, "Online Education Systems: Discussion and definition of terms," *NKI Distance Education,* 2002.

[7] N. Cavus and A. a. M. Momani, "Computer aided evaluation of learning management systems," *Procedia - Social and Behavioral Sciences,* vol. 1, pp. 426-430, 2009.

[8] P. Arroway, E. Davenport, G. Xu, and D. Updegrove, "EDUCAUSE core data service fiscal year 2009 summary report," *Boulder, CO: EDUCAUSE,* 2010.

[9] T. Browne, M. Jenkins, and R. Walker, "A longitudinal perspective regarding the use of VLEs by higher education institutions in the United Kingdom," *Interactive Learning Environments,* vol. 14, pp. 177-192, 2006.

[10] M. Robinson and M. Ally, "Transition to e-Learning in a Gulf Arab Country," in *The 2nd Annual Forum on e-Learning Excellence in the Middle East, Dubai, UAE,* 2009.

[11] S. Abdallah and F. Albadri, ICT Acceptance, Investment and Organization: Cultural Practices and Values in the Arab World: IGI Global, 2010.

[12] F. Lasrado, "Attitudes towards e-learning: Exploratory evidence from UAE," in *The 2nd Annual Forum on e-Learning Excellence in the Middle East, Dubai, UAE,* 2009.

[13] F. D. Davis, "Perceived usefulness, perceived ease of use, and user acceptance of information technology," *MIS quarterly,* pp. 319-340, 1989.

[14] A. Mirza, "Is E-Learning Finally Gaining Legitimacy in Saudi Arabia?," *Saudi Computer Journal,* vol. 6, 2007.

[15] H. Al-Khalifa, "JUSUR: The Saudi Learning Management System," the 2nd Annual Forum on e-Learning Excellence in the Middle East, 2009, Dubai, UAE., 2009.

[16] H. B. Hussein, "attitudes of Saudi universities faculty members towards using learning management system (jusur)," *TOJET,* vol. 10, 2011.

[17] M. G. Al-Joudi, "Enhancement of the performance of Taif University staff members in the area of information technology: A training need assessment," presented at the the 2nd international conference of e-learing and distance education: unique learning for next generation 21-24 FEB. 2011. Riyadh.

[18] R. Woods, J. D. Baker, and D. Hopper, "Hybrid structures: Faculty use and perception of web-based courseware as a supplement to face-to-face instruction," *The Internet and Higher Education,* vol. 7, pp. 281-297, 2004.

[19] R. Alebaikan and S. Troudi, "Blended learning in Saudi universities: challenges and perspectives," *Research in Learning Technology,* vol. 18, 2010.

[20] R. Alebaikan and S. Troudi, "Online discussion in blended courses at Saudi Universities," *Procedia - Social and Behavioral Sciences,* vol. 2, pp. 507-514, 2010.

[21] Y. Zhao, K. Pugh, S. Sheldon, and J. Byers, "Conditions for classroom technology innovations," *The Teachers College Record,* vol. 104, pp. 482-515, 2002.

[22] V. Venkatesh, M. G. Morris, G. B. Davis, and F. D. Davis, "User acceptance of information technology: Toward a unified view," *MIS quarterly,* pp. 425-478, 2003.

[23] A. Mulkeen, "What can policy makers do to encourage integration of information and communications technology? Evidence from the Irish school system," *Technology, Pedagogy and Education,* vol. 12, pp. 277-293, 2003.

[24] M. S. Asiri, R. Mahmud, K. Abu-Bakar, and A. F. Ayub, "Factors influencing the use of learning management system in Saudi Arabian Higher Education: A theoretical framework," *Higher Education Studies,* vol. 2, p. p125, 2012.

[25] M. Fishbein and I. Ajzen, Belief, attitude, intention and behavior: An introduction to theory and research, 1975.

[26] R. Hermans, J. Tondeur, J. van Braak, and M. Valcke, "The impact of primary school teachers' educational beliefs on the classroom use of computers," *Computers & Education,* vol. 51, pp. 1499-1509, 2008.

[27] C.-P. Kao and C.-C. Tsai, "Teachers' attitudes toward web-based professional development, with relation to Internet self-efficacy and beliefs about web-based learning," *Computers & Education,* vol. 53, pp. 66-73, 2009.

[28] M. J. Asiri, R. Mahmud, K. A. Bakar, and A. F. M. Ayub, "Role of Attitude in Utilization of Jusur LMS in Saudi Arabian Universities," *Procedia - Social and Behavioral Sciences,* vol. 64, pp. 525-534, 2012.

[29] A. Albalawi and M. Badawi, "Teachers' Perception of E-learning at the University of Tabuk," in *World Conference on E-Learning in Corporate, Government, Healthcare, and Higher Education,* 2008, pp. 2434-2448.

[30] I. Alzahrani, "Evaluate Wiki technology as e-Learning tool from the point view of Al-Baha university students: A pilot study with undergraduate students in both faculties of science and education," in *annual meeting of the Education, Learning, Styles, Individual Differences Network (ELSIN). June,* 2012, pp. 26-28.

[31] R. H. Shroff, C. Deneen, and E. M. Ng, "Analysis of the technology acceptance model in examining students' behavioural intention to use an e-portfolio system," *Australasian Journal of Educational Technology,* vol. 27, pp. 600-618, 2011.

[32] J.-H. Wu, R. D. Tennyson, and T.-L. Hsia, "A study of student satisfaction in a blended e-learning system environment," *Computers & Education,* vol. 55, pp. 155-164, 2010.

[33] J. Zouhair, "Surveying Learners' Attitudes Toward a Saudi E-learning System," *International Journal of Information and Electronics Engineering,* vol. 12, September 2012 2012.

[34] M. Rogers Everett, "Diffusion of innovations," *New York,* 1995.

[35] V. Venkatesh, J. Thong, and X. Xu, "Consumer acceptance and use of information technology: Extending the unified theory of acceptance and use of technology," *MIS quarterly,* vol. 36, pp. 157-178, 2012.

[36] W. H. DeLone and E. R. McLean, "The DeLone and McLean model of information systems Success: A ten-year update," *Journal of management information systems,* vol. 19, pp. 9-30, Spring2003 2003.

[37] J. E. Bailey and S. W. Pearson, "Development of a tool for measuring and analyzing computer user satisfaction," *Management science,* vol. 29, pp. 530-545, 1983.

[38] W. J. Doll and G. Torkzadeh, "The measurement of end-user computing satisfaction," *MIS quarterly,* pp. 259-274, 1988.

[39] Q. Ma and L. Liu, "The technology acceptance model: a meta-analysis of empirical findings," *Journal of Organizational and End User Computing (JOEUC),* vol. 16, pp. 59-72, 2004.

[40] D. Kim and H. Chang, "Key functional characteristics in designing and operating health information websites for user satisfaction: An

application of the extended technology acceptance model," *International Journal of Medical Informatics,* vol. 76, pp. 790-800, 2007.

[41] J.-W. Moon and Y.-G. Kim, "Extending the TAM for a World-Wide-Web context," *Information & Management,* vol. 38, pp. 217-230, 2001.

[42] P. Y. Chau, "An empirical assessment of a modified technology acceptance model," *Journal of management information systems,* vol. 13, pp. 185-204, 1996.

[43] K. Mathieson, "Predicting user intentions: comparing the technology acceptance model with the theory of planned behavior," *Information systems research,* vol. 2, pp. 173-191, 1991.

[44] D. A. Adams, R. R. Nelson, and P. A. Todd, "Perceived usefulness, ease of use, and usage of information technology: a replication," *MIS quarterly,* pp. 227-247, 1992.

[45] S. Al-Gahtani, "The applicability of TAM outside North America: an empirical test in the United Kingdom," *Information Resources Management Journal (IRMJ),* vol. 14, pp. 37-46, 2001.

[46] C.-S. Ong and J.-Y. Lai, "Gender differences in perceptions and relationships among dominants of e-learning acceptance," *Computers in Human Behavior,* vol. 22, pp. 816-829, 2006.

[47] J. C. Roca, C.-M. Chiu, and F. J. Martínez, "Understanding e-learning continuance intention: An extension of the Technology Acceptance Model," *International Journal of Human-Computer Studies,* vol. 64, pp. 683-696, 2006.

[48] S. Psycharis, G. Chalatzoglidis, and M. Kalogiannakis, "STUDENTS' ACCEPTANCE OF A LEARNING MANAGEMENT SYSTEM FOR TEACHING SCIENCES IN SECONDARY EDUCATION," *ICT AND OTHER RESOURCES FOR TEACHING/LEARNING SCIENCE,* p. 70.

[49] F. D. Davis Jr, "A technology acceptance model for empirically testing new end-user information systems: Theory and results," Massachusetts Institute of Technology, 1986.

[50] MOHE. (2013, 2/10/2013). *Ministry of Higher Education: Higher Education Statistics Center (Universities Statistics).* Available: http://www.mohe.gov.sa/en/studyinside/universitiesStatistics/Pages/default.aspx

[51] S. Taylor and P. Todd, "Assessing IT usage: The role of prior experience," *MIS quarterly,* vol. 19, pp. 561-570, 1995.

[52] H.-P. Shih, "Extended technology acceptance model of Internet utilization behavior," *Information & Management,* vol. 41, pp. 719-729, 2004.

[53] S. Y. Park, "An Analysis of the Technology Acceptance Model in Understanding University Students' Behavioral Intention to Use e-Learning," *Educational Technology & Society,* vol. 12, pp. 150-162, 2009.

[54] K.-T. Wong, P. S. C. Goh, and M. K. Rahmat, "Understanding Student Teachers' Behavioural Intention to Use Technology: Technology Acceptance Model (TAM) Validation and Testing," *Online Submission,* vol. 6, pp. 89-104, 2013.

[55] R. A. Sánchez and A. D. Hueros, "Motivational factors that influence the acceptance of Moodle using TAM," *Computers in Human Behavior,* vol. 26, pp. 1632-1640, 2010.

[56] N. Kripanont, "Using Technology Acceptance Model to Investigate Academic Acceptance of the Internet," *Journal of Business Systems, Governance, and Ethics,* vol. 1, pp. 13-28, 2006.

[57] C.-S. Ong, J.-Y. Lai, and Y.-S. Wang, "Factors affecting engineers' acceptance of asynchronous e-learning systems in high-tech companies," *Information & Management,* vol. 41, pp. 795-804, 2004.

[58] R. C. King and M. L. Gribbins, "Internet technology adoption as an organizational event: an exploratory study across industries," in *System Sciences, 2002. HICSS. Proceedings of the 35th Annual Hawaii International Conference on*, 2002, pp. 2683-2692.

[59] V. Venkatesh and F. D. Davis, "A theoretical extension of the technology acceptance model: four longitudinal field studies," *Management science,* vol. 46, pp. 186-204, 2000.

[60] R. Thompson, D. Compeau, and C. Higgins, "Intentions to use information technologies: An integrative model," *Journal of Organizational and End User Computing (JOEUC),* vol. 18, pp. 25-46, 2006.

[61] K. B. Wright, "Researching Internet-based populations: Advantages and disadvantages of online survey research, online questionnaire authoring software packages, and web survey services," *Journal of Computer-Mediated Communication,* vol. 10, pp. 00-00, 2005.

[62] I.-L. Wu, J.-Y. Li, and C.-Y. Fu, "The adoption of mobile healthcare by hospital's professionals: An integrative perspective," *Decision Support Systems,* vol. 51, pp. 587-596, 2011.

[63] [63] U. Sekaran and R. Bougie, *Research Methods for Business: A Skill Building Approach*: John Wiley & Sons, 2010.

[64] W. G. Zikmund, J. C. Carr, and M. Griffin, *Business research methods*: CengageBrain. com, 2012.

[65] R. W. Brislin, "The wording and translation of research instruments," 1986.

[66] A. A. Jemain, A. Al-Omari, and K. Ibrahim, "Multistage median ranked set sampling for estimating the population median," *Journal of Mathematics and Statistics,* vol. 3, p. 58, 2007.

[67] A. J. Rockinson-Szapkiw, A. Knight, and J. M. Tucker, "Prezi: Trading Linear Presentations for Conceptual Learning Experiences in Counselor Education," 2011.

[68] L. J. Cronbach, "Coefficient alpha and the internal structure of tests," *Psychometrika,* vol. 16, pp. 297-334, 1951.

[69] J. Hair, W. Black, B. Babin, R. Anderson, and R. Tatham, "Multivariate Data Analysis: Pearson Education," *New Jersey: Hoboken,* 2006.

[70] J. C. Nunnally, *Psychometric Theory*: New York: McGraw-Hill, 1967.

[71] U. Sekaran, *Research methods for business: A skill building approach*: New York, USA: John Wiley & Sons., 2006.

[72] M. Afshari, K. A. Bakar, W. S. Luan, B. A. Samah, and F. S. Fooi, "Factors affecting teachers' use of information and communication technology," *International Journal of Instruction,* vol. 2, pp. 77-104, 2009.

[73] M. S. Albalawi, "Critical factors related to the implementation of web-based instruction by higher-education faculty at three universities in the Kingdom of Saudi Arabia," The University of West Florida, 2007.

[74] D. M. Ball and Y. Levy, "Emerging Educational Technology: Assessing the Factors that Influence Instructors' Acceptance in Information Systems and Other Classrooms," *Journal of Information Systems Education,* vol. 19, pp. 431-444, 2008.

[75] S. S. Al-Gahtani, G. S. Hubona, and J. Wang, "Information technology (IT) in Saudi Arabia: Culture and the acceptance and use of IT," *Information & Management,* vol. 44, pp. 681-691, 2007.

# Spatial Domain Image Steganography based on Security and Randomization

NamitaTiwari
Department of CSE & IT
MANIT
Bhopal, India

Dr. Madhu Sandilya
Department of ECE
MANIT
Bhoapl, India

Dr. Meenu Chawla
Department of CSE & IT
MANIT
Bhopal, India

*Abstract*—In the present digital scenario secure communication is the prime requirement. Commonly, cryptography used for the said purpose. Another method related to cryptography is used for the above objective is Steganography. Steganography is the art of hiding information in some medium. Here we are using image as a means for covering information. Spatial domain image Steganography has been used for the work because of its compatibility to images. Objective of the paper is to increase the capacity of hidden data in a way that security could be maintained. In the current work MSB of the randomly selected pixel have been used as indicator. Result analysis has been performed on the basis of different parameters like PSNR, MSE and capacity.

*Keywords-- Spatial domain; PSNR; MSE*

## I. INTRODUCTION

Steganography can be used to hide or cover the existence of communication of encrypted data. A major drawback to encryption is that the existence of data is not hidden. Data that has been encrypted, although unreadable, still exists as data. A solution to this problem is Steganography [1]. The purpose of both Steganography and Cryptography is to provide secret communication. Cryptography hides the contents of a secret message from an attacker, whereas Steganography even conceals the existence of the message. In cryptography, the system is broken when the attacker can read the secret message [2]. Breaking a steganographic system has two stages: first the attacker can detect that Steganography has been used second he is able to read the embedded message.

In section II requirement and importance of steganography has been described. Section III explains about different techniques of image steganography. Section IV shows LSB technique. Section V describes about different LSB based methods. Section VI defines the objective of proposed work. Section VII explains different parameters for result analysis. Section VIII shows the comparative result analysis. Section IX describes conclusion and future work of the paper.

## II. REQURIMENTS FOR STEGANOGRAPHIC SYSTEM

- Imperceptibility: The stego image and original image should be perceptually identical.

- Undetectable embedded data.

- Security

- Maximizing Capacity of embedded data.

- Robustness: The embedded data should survive against various attacks.

Applications and Importance of a Steganographic system that it is used as Security reinforcement layer to cryptography [1]. It is used in digital watermarks, fingerprinting, defense, business, and education field. Image Steganography is about exploiting the limited powers of the human visual system (HVS)[2]. Within reason, any plain text, cipher text, other images, or anything that can be embedded in a bit stream can be hidden in an image. Image Steganography has come quite far in recent years with the development of fast, powerful graphical computers.

## III. TECHNIQUES FOR IMAGE STEGANOGRAPHY

### A. Spatial Domain based Steganography

It includes LSB (Least Significant Bit) Steganography. The spatial methods are most frequently employed because of fine concealment, great capability of hidden information and easy realization. LSB Steganography includes two schemes: Sequential Embedding and Scattered Embedding. [4]

### B. Transform Domain based Steganography

The method of transform domain Steganography is to embed secret data in the transform coefficients.

### C. Document based Steganography

This method embeds data in documents files by adding tabs or spaces to .txt or .doc files.

### D. File Structure based Steganography

This method inserts secrets data in the redundant bits of cover files, such as the reserved bits in the file header or the marker segments in the file format.

## IV. SPATIAL DOMAIN EMBEDDING

In the LSB technique, the LSB of the pixels is replaced by the message to be sent, this has the effect of distributing bits evenly, thus on average only half of the LSB's will be modified [4, 5].

Least Significant Bit Method
Consider a 24-bit picture
Data to be inserted: character 'A': (10000011)
3 pixels will be used to store one character of 8-bits

Example:

| | | |
|---|---|---|
| 00100111 | 11101001 | 11001000 |
| 00100111 | 11001000 | 11101001 |
| 11001000 | 00100111 | 11101001 |

Embedding 'A'

| | | |
|---|---|---|
| 0010011**1** | 1110100**0** | 1100100**0** |
| 0010011**0** | 1100100**0** | 1110100**0** |
| 1100100**1** | 0010011**1** | 1110100**1** |

## V. LSB BASED STEGANOGRAPHY METHODS

### A. Stego One Bit

This method changes only single LSB of the pixel. Changing the LSB will only change the integer value of the byte by one. This small change is not noticeable. This is the first method to be tested and will involve encoding some of the basic processes required for later Steganographic methods to be tested also. This should have very less effect on the appearance of the image [2].

### B. Stego Two Bit

Using this method two LSBs of one of the colours in the RGB value of the pixels will be used to store message bits in the image.[2] The advantage of this method is that twice as much information can be stored here than in the previous method.

### C. Stego Three Bit

Using this method three LSBs of one of the colours in the RGB value of the pixels will be used to store message bits.The data hiding capacity is three times the storage capacity of Stego One Bit but the image will be even more distorted.

### D. Stego Four Bit

Using this method four LSBs of one of the colours in the RGB value of the pixels will be used to store message bits.The data hiding capacity is 4 times the storage capacity of stego 1 bit, but the image will be more distorted.

### A. Stego Colours Cycle

In order to make the detection of the hidden data more difficult it was decided to cycle through the colours values in each of the pixels in which to store the data [2, 7]. This also means that the same colours were not constantly being

changed. For example the first data bit could be stored in the LSB of the blue value of the pixel, the second data bit in the red value and the third data bit in the green value.SCC technique is an enhancement. This technique is more secure than the LSB. But still it is suffers detecting the cycling pattern that will reveal the secret data. Also it has less capacity than LSB.

### E. Pixel Indicator High Capacity Technique

It uses the least two significant bits of one of the channels to indicate existence of data in the other two channels [5]. Table 1 shows meaning of indicator values for pixel indicator technique.

### F. Triple-A: Based on Randomization

This algorithm can be divided into two major parts: Encryption and Hiding [7].

*1) Encryption:* Part one is related to encrypting the message (M) using AES algorithm, which will produce Enc (M, K). In implementation the key K can be generated from a set of user password.

*2) Hiding:* The RGB Image is used as a cover media. Enc (M, K) is hidden according to triple-A algorithm, which needs to have a pseudorandom number generator (PRNG). The assumption for PRNG is to give two new random numbers in every iteration. The seeds of these PRNGs namely Seed1 (S1) and Seed2 (S2) are formed as a function of the Key (K). S1 is restricted to generate numbers in [0, 6]. S1 random number is used to determine the component of the RGB image, which is going to be used in hiding the encrypted data Enc (M, K). Table 2 shows how (S1) random number selects the RGB components. S2 is restricted to the interval [1, 3]. S2 random number determines the number of the component(s) least significant bits that is used to hide the secret data. Table 3 shows how (S2) random number determines the number of component bits. By combining data from the previous tables, we can see that the minimum number of bits used in each pixel is 1. If we use only one bit of one chosen components of the RGB image.

TABLE I.     MEANING OF INDICATOR VALUES FOR PIXEL INDICATOR TECHNIQUE

| Indicator Channel | Channel-1 | Channel-2 |
|---|---|---|
| 00 | No Hidden data | No Hidden data |
| 01 | No Hidden data | 2 Bits of Hidden Data |
| 10 | 2 Bits of Hidden Data | No Hidden data |
| 11 | 2 Bits of Hidden Data | 2 Bits of Hidden Data |

TABLE II.     SEED 1 RANDOM NUMBER USAGE

| | Random Number | Meaning to the algorithm |
|---|---|---|
| **1st PRNG** | 0 | Use R |
| | 1 | Use G |
| | 2 | Use B |
| | 3 | *Use RG* |
| | 4 | *Use RB* |
| | 5 | *Use GB* |
| | 6 | *Use RGB* |

TABLE III.     SEED 2 RANDOM NUMBER USAGES

| | Random Number | Meaning to the algorithm |
|---|---|---|
| **2st PRNG** | 1 | Use 1 bit of the component(s) |
| | 2 | Use 2 bit of the component(s) |
| | 3 | Use 3 bit of the component(s) |

TABLE IV. PROPSED ALGORITHM

| 3 MSB of channel (R/G/B) | | | Channel (R/G/B) |
|---|---|---|---|
| 0 | 0 | 0 | 1 bit Hidden Data |
| 1 | 0 | 0 | 2 bit Hidden Data |
| 1 | 1 | 0 | 3 bit Hidden Data |
| 1 | 1 | 1 | 4 bit Hidden Data |

The maximum is 9 bits if we used all the three components with three bits.

Capacity factor = number bits used inside a pixel to hide part of the secret message/ the number of bits in the pixels itself

Capacity factor can be in the range from 1/24 to 9/24.

## VI. OBJECTIVE OF THE STUDY

Objective of the work is to increase the capacity of hidden data in the way that image should be less distorted and unauthorized person cannot detect that Steganography is going on. Randomization has been used to select the pixel to overcome the sequential pattern. It is proposed to use MSB (Most Significant Bit) of the pixel as indicator for data hiding. Message is hiding in all three channels according to the indicator bits in MSB.

In Proposed Method there are two phases.

### A. First Phase

The indicator channel has been decided by the user, suppose red is an indicator channel and three MSB of red channel are 101, and then MSB of indicator channel will decide that which channel is used for data hiding. (101 used for RGB respectively).

For example:

R        G        B

1        0        1

Here 1 indicates to hide the data and 0 indicate for not hiding the data. In above case, channel R and B will use for hiding data and channel G will not use for hiding data.

### B. Second Phase

In first phase we have decided the channels for data hiding, now in second phase we will decide the no. of bits to be hidden. Three MSB of selected channel (R/G/B) will decide that how many bits of message will be hidden in that particular channel. Table 4 shows the method of second phase.

To optimize the proposed algorithm it has been checked on different parameters.

## VII. PARAMETERS FOR RESULT ANALYSIS

### A. Capacity

This is the term refers to the amount of data that can be hidden in the medium. It is defined as the maximum size that can be hidden in the medium. It is defined as the maximum size that can be embedded subject to certain constraints [10].

Capacity factor = number bits used inside a pixel to hide part of the secret message/ the number of bits in the pixels itself

### B. MSE

Mean squared error is the average squared difference between a reference image and a distorted image [10].

$$MSE = \sum_{i=1}^{M} \sum_{j=1}^{N} (X_{i,j} - Y_{i,j})^2 \qquad (1)$$

### C. PSNR

Peak signal to noise ratio is the ratio between the reference signal and the distorted signal in an image [9, 10].

$$PSNR = 10 \log_{10} \left[ \frac{I_{max}^2}{MSE} \right] dB \qquad (2)$$

### D. Histogram Analysis

Histogram analysis has been performed on original image and stego image to differentiate between both images. For more accurate analysis it is proposed to perform Histogram analysis on each channel separately.

## VIII. RESULT ANALYSIS

Table 5 shows the comparison of proposed algorithm with existing 7 algorithms. 1 bit and 2 bits are better in MSE and PSNR then proposed algorithm but capacity of proposed algorithm is better than all existing algorithms.

5KB of data has been taken for hiding in 512*512 size image. All Steganographic algorithms have been performed on same data and on same image size. It is observed that proposed algorithm is using less amount of image for hiding the same data. For hiding 5 KB data the algorithm takes only 2.50 % pixels of whole image and 1bit method takes 15.04 % pixels of the image. If requirement for Steganographic algorithm is high capacity and good MSE and PSNR then proposed algorithm could be used for this purpose.

## IX. CONCLUSION AND FUTURE WORK

Steganography has its own place in the field of security. Steganography used in an open-systems environment such as the Internet and Far-fetched applications, privacy protection, authentication, data integrity, intellectual property rights protection.

Proposed method is achieving highest capacity among all existing methods without any distortion in image. When proposed method has been performed on different images, it has given constant result but other existing methods gave different results on different images. Proposed method is secure and undetectable because of randomness.

In future, techniques to improve security for data hiding by using randomization will be used to extend this work.

TABLE V.        RESULT ANALYSIS

| | 1 bit | 2 bit | 3 bit | 4 bit | SCC | PIT | AAA | Proposed Algorithm |
|---|---|---|---|---|---|---|---|---|
| MSE | .0261 | .062 | .1646 | .6216 | .0259 | .0611 | .1060 | .1606 |
| PSNR | 63.96 | 60.20 | 55.96 | 50.19 | 63.99 | 60.27 | 57.87 | 56.07 |
| Capacity in percentage | 15.04 | 7.52 | 5.01 | 3.76 | 5.01 | 7.52 | 3.00 | 2.50 |

REFERENCES

[1] S. Venkatraman, A. Abraham, M. Paprzycki,"Significance of steganography on Data Security", *International conference on Information Technology: Coding and Computing(ITCC'04)*, Las Vegas, 5-7 April 2004.

[2] K. Baily, K. Curran, "An Evaluation of Image Based Steganography Methods using visual inspection and automated detection techniques", Multimedia Tools & Applications, Vol. 30, No.1, pp. 55-88, July 2006, Springer.

[3] V. Lokeswara Reddy, A. Subramanyam, P. Chenna Reddy, "Implementation of LSB Steganography and Its Evaluation for Various File Formats" ,International Journal of Advanced Networking and Applications, Vol.02, Issue:05, pp.868-872, 2011

[4] Chen Ming, Zhang Ru, Niu Xinxin, Yang Yixian,"Analysis of Current Steganography Tools: Classification & Features", *IEEE International conference on Intelligent Information Hiding and Multimedia signal Processing (IIH-MSP' 06)* ,Pasadena, CA, USA, 2006.

[5] Adnan Abdul Aziz Gutub, "Pixel Indicator Technique For RGB Image Steganography", *Journal of Emerging Technologies in Web Intelligence, Vol.2 No.1, February 2010, Academy Publisher.*

[6] Mohammad Tanvir Parvez and Adnan Gutub, "RGB Intensity Based Variable-Bits Image Steganography", *APSCC 2008 – 3rd IEEE Asia Pacific Services Computing Conference, Yilan, Taiwan, 9-12* December 2008.

[7] Adnan Gutub, Ayed Al-Qahtani, AbdulazizTabakh," Triple-A: Secure RGB Image     Steganography Based on Randomization*" IEEE International Conference on Computer System and Application, Rabat, May-June 2009*, pp.400-403

[8] Mamta Juneja and Parvinder Singh Sandhu, " Designing a Robust Image Steganography  Technique Based on LSB Insertion and Encryption" , *IEEE International Conference on Advances in recent technologies in Communication and Computing , Kottayam, Kerela, 2009, pp.302-305*

[9] Amirtharajan, R., RambhatlaSubrahmanyam, Pakalapati J S Prabhakar, Kavitha, R, and Balaguru, "MSB over hides LSB - A dark communication with integrity", IEEE International Conference on Internet Multimedia Services Architecture and Applications, IMSAA 2011.

[10] Rengarajan Amirtharajan, K. Ramkrishnan, M. Vivek Krishna, Nandhini. J, John Bosco and Balaguru Rayappan, "Who decides hiding capacity? I, the Pixel Intensity, *IEEE International Conference on Recent Advances in Computing and Software System"(RACSS'12), Chennai  2012, pp-71-76.*

[11] Mehdi Kharrazi, Husrev T. Sencar, and NasirMemon" Image Steganography: Concepts and    Practice*" WSPC/Lecture Notes Series: 9in x 6in, Institute of mathematical sciences, Singapore* April 22,2004

[12] M Amin, M. Salleh, *S. Ibrahim, M.R.K atmin, and M.Z.I. Shamsuddin,"* Information Hiding using  Steganography", *4th National Conference on Telecommunication Technology Proceedings, Shah Alam, Malaysia*, 2003, pp.21-25.

[13] Mehdi Kharrazi, Husrev T. Sencar, Nasir Memon, "Performance Study Of Common Image   Steganography And Steganalysis Techniques", *Journal of Electronic Imaging Vol.15 No.4, Oct–Dec 2006, pp. 041104-1-15.*

[14] Saeed Sarreshtedari, Mohsen Ghotbi and Shahrokh Ghaemmmaghami, "On The Effect Of Spatial To Compressed Domain Transformation In LSB-based Image Steganography", *IEEE International Conference on Computer Systems and Applications* , 2009, pp.260-264.

[15] Hassan Mathkour, Batool Al-Sadoon, Ameur Touir, "A New Image Steganography Technique", *IEEE 4th International Conference on Wireless Communications, Networking and Mobile Computing*, 2008, pp. 1-4.

[16] Farhan Khan and Adnan Abdul-Aziz Gutub, "Message Concealment Techniques using Image based Steganography*",    Department of Computer Engineering, King Fahd University of Petroleum & Minerals, Dhahran*,31261, Kingdom of Saudi Arabia

[17] HediehSajedi, Mansour Jamzad,"Cover Selection Steganography Method Based on Similarity of Image *Blocks", 8th International Conference on Computer and Information Technology Workshops*, *IEEE,* 2008, pp. 379-384.

[18] Yanming Di, Huan Liu, Avinash Ramineni, and Arunabha Sen," Detecting Hidden Information in Images : A Comparative Study", *Department of Computer Science and Engineering Arizona State University*, Tempe, AZ 85287,2005

[19] Li Zhi, Sui Ai Fen, "Detection of Random LSB Image Steganography", IEEE 60th *Vehicular technology Conference*, 2004, pp.2113-2117.

[20] Wien Hong and Tung-Shou Chen, "A Novel Data Embedding Method Using Adaptive Pixel Pair Matching", *IEEE Transactions on Information Forensics and Security, Vol.7,No.1,February 2012.*

[21] WeiqiLuo,Fangjun Huang and Jiwu Huang, "Edge Adaptive Image Steganography Based on LSB Matching Revisited" , *IEEE Transactions on Information Forensics and Security, Vol.5, No.2, June  2010.*

[22] M. Khodaei and K. Faez, "New adaptive Steganographic method using least significant bit substitution and pixel value differencing" , IET Image processing , Vol.6, iss.6, pp 677-686, 2012

# Mining Interesting Positive and Negative Association Rule Based on Improved Genetic Algorithm (MIPNAR_GA)

Nikky Suryawanshi Rai
PG Research Scholar (CSE),
RITS, Bhopal (M.P.) India

Susheel Jain
Asst Prof(CSE)
RITS, Bhopal (M.P.) India

Anurag Jain
HOD (CSE)
RITS, Bhopal (M.P.) India

*Abstract*—**Association Rule mining is very efficient technique for finding strong relation between correlated data. The correlation of data gives meaning full extraction process. For the mining of positive and negative rules, a variety of algorithms are used such as Apriori algorithm and tree based algorithm. A number of algorithms are wonder performance but produce large number of negative association rule and also suffered from multi-scan problem. The idea of this paper is to eliminate these problems and reduce large number of negative rules. Hence we proposed an improved approach to mine interesting positive and negative rules based on genetic and MLMS algorithm. In this method we used a multi-level multiple support of data table as 0 and 1. The divided process reduces the scanning time of database. The proposed algorithm is a combination of MLMS and genetic algorithm. This paper proposed a new algorithm (MIPNAR_GA) for mining interesting positive and negative rule from frequent and infrequent pattern sets. The algorithm is accomplished in to three phases: a).Extract frequent and infrequent pattern sets by using apriori method b).Efficiently generate positive and negative rule. c).Prune redundant rule by applying interesting measures. The process of rule optimization is performed by genetic algorithm and for evaluation of algorithm conducted the real world dataset such as heart disease data and some standard data used from UCI machine learning repository.**

*Keywords*—*Association rule mining; negative rule and positive rules; frequent and infrequent pattern set; genetic algorithm*

## I. INTRODUCTION

Association rule mining is a method to identify the hidden facts in large instances database and draw interferences on how subsets of items influence the existence of other subsets. Association rule mining aims to discover strong or interesting relation between attributes. All generalized frequent pattern sets are not very efficient because a segment of the frequent pattern sets are redundant in the association rule mining. This is why, traditional mining algorithm produces some uninteresting rules or redundant rules along with the interesting rule. This problem can be overcome with the help of genetic algorithm. Most of the data mining approaches use the greedy algorithm in place of genetic algorithm. Genetic algorithm is produced by optimized result as compare to the greedy algorithm because it performs a comprehensive search and better attributes interaction [1]. In genetic algorithm population evolution is simulated. Genetic algorithm is an organic technique which uses gene as an element on which solutions (individuals) are manipulated. Generally association

rule is used to finding positive relationship between the data set. Negative association rule is also vital in analysis of intelligent data. Negative association rule mining is adopted where a domain has too many factors and large number of infrequent pattern sets in transaction database. Negative association rule mining works in reverse manner and it define decision making capability, whether which one is important instead of checking all rules. However problem with the negative association rule is it uses huge space and can take more time to produce the rules as compare to the conventional mining association rule. In the generalized association rule database is scanned once and transaction is transformed into space reduced structure. The association rule mining problem can be decomposed in statistical and unconditional attributes in a database. The application of association rule mining is used to analyze various situation like market basket analysis, banks, whether prediction, pattern reorganization, multimedia data etc.

The process of optimization of interesting association rule mining used genetic algorithm. Genetic algorithm works in multiple levels of constraints for minimum support value and individual confidence value of frequent and infrequent patterns. The proposed method enhances the process of rule optimization for large datasets. The rest of paper is organized as follows. In Section II describes about related work of association rule mining. Section III describes about proposed method. Section IV describes about experimental result algorithm followed by a conclusion in Section V.

## II. RELATED WORK

This section describes some related work to negative and positive association rule mining.

An Improved apriori algorithm is used minimum supporting degree and degree of confidence ,for extracting association rules .But it has suffered from "frequent pattern sets explodes "and "rare item dilemma " [2].

Improved multiple minimum support (MSapriori) based on notion of support difference and define how to deal with the problem caused by frequent pattern sets explodes ,but still suffer from rare item dilemma[3].

Primary stage of association rules, all algorithms based on single minimum support and those algorithms suffer from "rule missed" and "rule explosion" problem. An efficient

method to extract rare association rules .In this method the probability and introduces multiple minsupp value to discover rare association rules. One obstacle of this algorithm is that, it produces large number of uninteresting pattern sets [4].

PNAR_MDB on PS measures is introduced to discover PNAR in multi-databases. PNAR_MDB on PS extract interesting association rules by weighting the database (the weight of database must be determined) and used the correlation coefficient to remove the confliction of rules [5].

Reveal knowledge hidden in the massive database and proposed an approach for Evaluation of exam paper. This paper introduces a new direction, applies interesting rules mining to evolution of completive exam and finds out some useful knowledge. But this algorithm need repeatedly database scan and takes more time to perform I/O operation [6].

Some algorithm uses comparison support and comparison confidence (comsup, comconf) for extracting interesting relationship between pattern sets [7].

According to correlation and dual confidence measures association rules are classified in to positive and negative association rules ,but one drawback of dual confidence, is if less confidence would be a lot of rules even produce large number of contradict rules ($\neg C \rightarrow \neg D$), if greater confidence may missed useful positive association rules[8].

Generalized Negative Association Rules (GNAR) is produced interesting negative rules ,this approach could speed up execution time efficiently through the domain taxonomy tree and extract interesting rules easily, advantage of taxonomy tree is to eliminate large number of useless transaction [9].

Another approach to solve key factors of interesting rules is PNAR algorithm, this algorithm efficiently define frequent pattern sets for interesting rules, NAR based on correlation coefficient and modified pruning strategy[10].

PNAR_IMLMS produces valid association rules based on correlation coefficient but one demerits of this algorithm, negative rules extract from uninteresting pattern sets which is useless [11]. Optimized association rule mining with genetic algorithm produces more reliable interesting rules compare to previous method.

Mining association rules using multiple support confidence values and several studies have been addressed the issue of mining association rules using Multiple Level Minimum Supports [12].

## III. PROPOSED ALGORITHM

This paper proposed a novel algorithm for optimization of association rule mining, the proposed algorithm resolves the problem of negative rule generation and also optimized the process of rule generation. Interesting association rule mining is a great challenge for large dataset. In the generation of interesting rules association existing algorithm or method generate a series of negative rules, which generate rules which affect performance of association rule mining. In the process of rule generation various multi objective associations rule mining algorithm is proposed but all these are not solve.

This paper proposed an improved approach to mine association rule In this algorithm we used a MLMS for multi level minimum support for constraints validation. The scanning of database divided into multiple levels as frequent level and infrequent level of data according to MLMS. The frequent data logically assigned 1 and infrequent data logically assigned 0 for MLMS process. The divided process reduces the uninteresting item in given database.

The proposed algorithm is a combination of MLMS and genetic algorithm along this used level weight for the separation of frequent and infrequent item. The multiple support value passes for finding a near level between MLMS candidates key. After finding a MLMS candidate key the nearest level divide into two levels, one level take a higher odder value and another level gain infrequent minimum support value for rule generation process. The process of selection of level also reduces the passes of data set. After finding a level of lower and higher of given support value, compare the both values of level by vector function. Here level weight vector function work as a fitness function to define the selection process of genetic algorithm

Here we implemented the combinatorial method of MLMS and genetic algorithm for the mining of positive and negative item sets. The key idea is to generate frequent and infrequent item sets and with these item sets positive and negative association rules are generated. MLMS algorithm is used for the generation of rules [12], since the association rule mining seems to be better when the association rules are less, hence the minimization of these positive and negative rules can be done using genetic algorithm. The proposed technique can be described as follows:

*1) Take an input dataset which contains number of attributes and instance values with single or multiple classes.*

*2) Initialize the data with length of the item sets k=2, 3, 4 and pass support and confidence (Para b).*

*3) Generate all the frequent and Infrequent item sets from MLMS algorithm for an item set of length k=2, 3, 4.*

*4) Generate positive association rules from frequent items sets and negative association rules from infrequent item sets.*

*5) Initialize all the general parameters involved in genetic algorithm.*

*6) Generate the child chromosomes of the positive and negative association rules and calculate the fitness value of each individual child chromosomes. Compare the individual fitness value of each child with the average fitness value and regenerate positive and negative association rules.*

*7) Crossover and mutate the remaining child chromosomes and reinitialize the fitness value and recalculate and regenerate final positive and negative rules.*

### A. Load Datasets

The association rules generated from the proposed algorithm needs datasets containing a number of transaction values. Here we use a number of datasets i.e. small and large dataset, a dataset with single and multiple classes. So the performance of the proposed methodology is tested for each datasets.
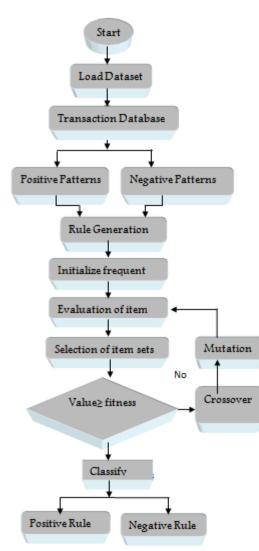
Fig. 1. Shows that proposed block model of algorithm.

### B. Support and Confidence

Here the association rules can be generated on the basis of item set length, support and confidence. Suppose sup and cf are the support and confidence respectively. Let k be the length of the item set. For an item set $A \subseteq I$, the support is A.count / |TD|, where A.count is the number of transactions in TD that contain the itemset A. The support of a rule $A \Rightarrow B$ is denoted as sup $(A \cup B)$, where A, $B \subseteq I$, and $A \cap B = \Phi$ while the confidence of the rule $A \Rightarrow B$ is defined as the proportion of s $(A \cup B)$ above s $(A)$, i.e., cf $(A \Rightarrow B)$ = s $(A \cup B)$ /s $(A)$.

### C. Generate Frequent and infrequent item sets

Here use MLMS algorithm for the generation of frequent and infrequent item sets. Form these frequent and infrequent item sets positive and negative association rules are generated.

A frequent itemset I: sup $(I) \geq$ minsupp
An infrequent itemset J: sup $(J) \leq$ minsup*p*

### D. Correlation factor

For the generation of positive and negative association rules from these item sets, first of all correlation coefficient between the items sets is computed using:

$$corr_{AB} = \frac{cov(A, B)}{\sigma_A \sigma_B}$$

Where cov(A, B) represents the covariance of two variables and σ represents the standard deviation. Then compare the correlation coefficient with the correlation strength. Generate all the rules of the form

Positive association rules:

$$A \cap B = \phi$$
Supp (AU B) ≥ minsupp
Supp (A U B) / supp (A) ≥ minconf

Negative association rules:

$$A \cap B = \phi$$
Sup (A) >= minsupp, Sup (B) > minsupp,
and sup (A U ~B) >= minsupp
Sup (A U ~B)/sup (A) >= minconf

If the correlation coefficient is greater than or equal to α and if they meet the conditions VARCC(A,B,α,mc)=1 and VARCC(¬A,¬B,α,mc)=1. if the correlation coefficient is lower than or equal to -α and if they meet the conditions VARCC(A,¬B,α,mc)=1 and VARCC(¬A,B,α,mc)=1.

### E. Initialization of Parameters

The genetic algorithm when applied should be initialized by certain parameters such as selection, crossover and mutation as well the number of iterations it will performed during working. There are various solutions that must be chosen randomly to form an initial population. The size of the population will depends on the problem

### F. Fitness Function

The population selection for Genetic Algorithm is based on Fitness Function:

$$m(S) = \frac{Ai}{wi} + \frac{Bi}{L \times (1 - wi)}$$

Ai = {frequent item support}

Wi= {level of Wight value of MLMS}

Bi = {those value or Data infrequent}

The selection policy based on the foundation of individual fitness and concentration p(i) is the selection of individual whose fitness value is greater than one and m(s) is a value whose fitness is less than one but close to the value of 1.

The genetic operators find out the search capability and convergence of the algorithm.

## G. Reproduction Operators

The child chromosomes that are not used in the sets will now be crossover and mutate so that the new fitness value is generated and again from parent, child chromosomes are generated. The process repeats until the rules generation finishes: Example:

$$1\ 0\ 1\ 0\ 0\ 1\ 0$$
$$\downarrow$$
$$1\ 0\ 1\ 0\ 1\ 1\ 0$$

Mutation operator has been chosen to insure high levels of diversity in the population. We adopted PCA-mutation in (Munteanu 1999b), and shown that it has very good capabilities in maintaining higher levels of diversity in the population. We briefly summarize the PCA-mutation operator, as follows: The population **X** of the GA can be viewed as a set of $N$ points in a $l$-dimensional space, where $N$ is the size of the population and $l$ is the length of the chromosome. It can be shown (Munteanu 1999b) that a GA converging has the effect of decreasing the number of Principal Components (PCs) as calculated with the Principal Components Analysis (PCA) method on data **X**.

  a) *Select a random point on the two parents.*

  b) *Split parents at this crossover point.*

  c) *Produce children's by exchanging trails.*

  d) *Mutation typically in range (0.6, 0.9).*
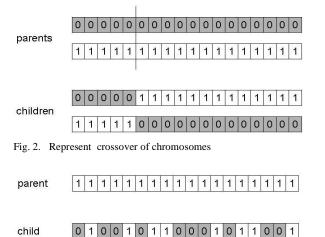


Fig. 2. Represent crossover of chromosomes



Fig. 3. Represent mutation of chromosomes.

## IV. SIMULATION RESULT

This section shows the performance of MIPNAR_GA algorithm for mining both interesting positive and negative rules. Experiments are performed on a computer Intel Pentium dual core processor with 2.10 GHZ of CPU, running on a Windows 7 ,64-bit operating system and 4 GB of memory .All codes are implemented under the Java Compiler (JDK 1.6 and Weka 3.6.9) and Net Beans IDE version 6.9. Test the performance of proposed algorithm on 4 datasets from UCI machine learning website, which involve, Heart diseases, Breast Cancer, Wine and Iris. All information related to datasets are shown in Table 1.

TABLE I.　CHARACTERISTICS OF DATASETS

| Dataset | No of Attributes | No of Instances | Classes |
|---|---|---|---|
| Heart Disease | 14 | 303 | 2 |
| Breast Cancer | 10 | 286 | 2 |
| Iris | 14 | 178 | 3 |
| Wine | 5 | 150 | 3 |

Because MIPNAR_GA is designed to mine positive and negative rules from positive (frequent) and negative (infrequent) patterns with different input parameter (support, confidence, itemset length), it will be compared with the base algorithm PNAR_IMLMS for mining interesting positive and negative rules. The results are representing in table 2 to 7 where the number of interesting positive (A→B) and negative rules are represent as (A→¬B, ¬A→B, ¬A→¬B).

TABLE II.　SHOW THAT GIVEN VALUE OF SUPPORT (65%) CONFIDENCE (55%) AND ITEM **LENGTH 2** ALGORITHM PNAR_IMLMS GENERATED TOTAL NUMBER OF INTERESTING POSITIVE AND NEGATIVE RULES FOR UCI DATA SET

| Datasets | | PNAR_IMLMS | | | |
|---|---|---|---|---|---|
| | | A→B | A→¬B | ¬A→ B | ¬A→¬B |
| Heart Disease | FIS | 8 | 3 | 1 | 8 |
| | inFIS | 0 | 16 | 20 | 14 |
| Breast Cancer | FIS | 33 | 0 | 0 | 42 |
| | inFIS | 0 | 17 | 17 | 23 |
| Iris | FIS | 7 | 0 | 0 | 8 |
| | inFIS | 0 | 12 | 12 | 16 |
| Wine | FIS | 19 | 5 | 4 | 16 |
| | inFIS | 0 | 12 | 14 | 43 |
| **Total** | | **67** | **65** | **68** | **170** |

TABLE III.　SHOW VALUE OF SUPPORT (75%) CONFIDENCE (65%) AND ITEM **LENGTH 3** ALGORITHM PNAR_IMLMS GENERATED NUMBER OF INTERESTING POSITIVE AND NEGATIVE RULES FOR UCI DATA SET

| Datasets | | PNAR_IMLMS | | | |
|---|---|---|---|---|---|
| | | A→B | A→¬B | ¬A→ B | ¬A→¬B |
| Heart Disease | FIS | 47 | 1 | 1 | 52 |
| | inFIS | 0 | 20 | 22 | 45 |
| Breast Cancer | FIS | 104 | 0 | 0 | 117 |
| | inFIS | 0 | 6 | 8 | 37 |
| Iris | FIS | 18 | 0 | 0 | 16 |
| | inFIS | 0 | 11 | 12 | 6 |
| Wine | FIS | 148 | 6 | 4 | 143 |
| | inFIS | 0 | 24 | 27 | 146 |
| **Total** | | **317** | **68** | **74** | **562** |

TABLE IV.    SHOW THAT GIVEN VALUE OF SUPPORT (55%) CONFIDENCE (45%) AND ITEM **LENGTH 4** ALGORITHM MIPNAR_GA GENERATED NUMBER OF INTERESTING POSITIVE AND NEGATIVE RULES FOR UCI DATA SET

| Datasets | | PNAR_IMLMS | | | |
|---|---|---|---|---|---|
| | | A→B | A→¬B | ¬A→B | ¬A→¬B |
| Heart Disease | FIS | 141 | 3 | 3 | 144 |
| | inFIS | 0 | 0 | 0 | 0 |
| Breast Cancer | FIS | 265 | 0 | 0 | 294 |
| | inFIS | 0 | 0 | 0 | 10 |
| Iris | FIS | 10 | 0 | 0 | 11 |
| | inFIS | 0 | 5 | 5 | 0 |
| Wine | FIS | 656 | 18 | 14 | 667 |
| | inFIS | 0 | 0 | 0 | 0 |
| Total | | 1072 | 26 | 22 | 1126 |

TABLE V.    SHOW THAT GIVEN VALUE OF SUPPORT (65%) CONFIDENCE (55%) AND ITEM LENGTH 2 ALGORITHM MIPNAR_GA GENERATED NUMBER OF INTERESTING POSITIVE AND NEGATIVE RULES FOR UCI DATA SET

| Datasets | | MIPNAR_GA | | | |
|---|---|---|---|---|---|
| | | A→B | A→¬B | ¬A→B | ¬A→¬B |
| Heart Disease | FIS | 3 | 0 | 0 | 2 |
| | inFIS | 0 | 7 | 9 | 10 |
| Breast Cancer | FIS | 12 | 0 | 0 | 16 |
| | inFIS | 0 | 6 | 6 | 8 |
| Iris | FIS | 2 | 0 | 0 | 3 |
| | inFIS | 0 | 5 | 5 | 8 |
| Wine | FIS | 10 | 1 | 0 | 9 |
| | inFIS | 0 | 8 | 8 | 16 |
| Total | | 27 | 26 | 28 | 72 |

TABLE VI.    SHOW THAT GIVEN VALUE OF SUPPORT (75%) CONFIDENCE (65%) AND ITEM LENGTH 3 ALGORITHM MIPNAR_GA GENERATED NUMBER OF INTERESTING POSITIVE AND NEGATIVE RULES FOR UCI DATA SET

| Datasets | | MIPNAR_GA | | | |
|---|---|---|---|---|---|
| | | A→B | A→¬B | ¬A→B | ¬ A→¬B |
| Heart Disease | FIS | 31 | 0 | 0 | 35 |
| | inFIS | 0 | 10 | 12 | 18 |
| Breast Cancer | FIS | 75 | 0 | 0 | 80 |
| | inFIS | 0 | 1 | 2 | 9 |
| Iris | FIS | 5 | 0 | 0 | 8 |
| | inFIS | 0 | 4 | 5 | 1 |
| Wine | FIS | 130 | 2 | 2 | 112 |
| | inFIS | 0 | 16 | 16 | 97 |
| Total | | 241 | 33 | 37 | 360 |

TABLE VII.    SHOW THAT GIVEN VALUE OF SUPPORT (55%) CONFIDENCE (45%) AND ITEM LENGTH 4 ALGORITHM MIPNAR_GA GENERATED NUMBER OF INTERESTING POSITIVE AND NEGATIVE RULES FOR UCI DATA SET

| Datasets | | MIPNAR_GA | | | |
|---|---|---|---|---|---|
| | | A→B | A→¬B | ¬A→B | ¬A→¬B |
| Heart Disease | FIS | 97 | 1 | 0 | 98 |
| | inFIS | 0 | 0 | 0 | 0 |
| Breast Cancer | FIS | 203 | 0 | 0 | 215 |
| | inFIS | 0 | 0 | 0 | 2 |
| Iris | FIS | 3 | 0 | 0 | 4 |
| | inFIS | 0 | 1 | 1 | 0 |
| Wine | FIS | 598 | 7 | 5 | 602 |
| | inFIS | 0 | 0 | 0 | 0 |
| Total | | 898 | 9 | 6 | 921 |

Table 2-7 shows the number of interesting positive and negative rules generated from useful positive and negative patterns with different input parameter. These rules are mined with two algorithms, the PNAR_IMLMS algorithm [12] and the MIPNAR_GA. For example, in Table 2 to 4 the number of interesting positive and negative rules mined by PNAR_IMLMS are 67 to 303 and 317 to 704 and 1072 to 1174, whereas in table 5 to 7 represent the total number of interesting positive and negative rules mined by MIPNAR_GA are 27 to 126 and 241 to 430 and 898 to 936 respectively .We can say that the algorithm MIPNAR_GA can successfully produce fewer rules than PNAR_IMLMS. In figure 3 to 5, P represent positive rule X→Y, $N_1$ represent A→¬B, $N_2$ represent ¬A→B, and $N_3$ represent ¬A→¬B.



Fig. 4.   Shows the comparative value of no of rules and reduces rules of two algorithms by optimization process from Tabel 3 to table 6.

Fig. 5. Shows the comparative value of no of rules and reduces rules of two algorithms by optimization process from Tabel 2 to table 5.
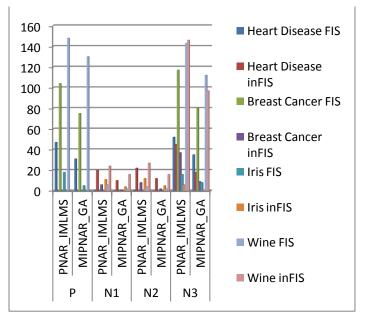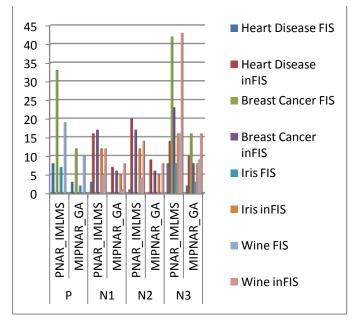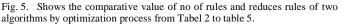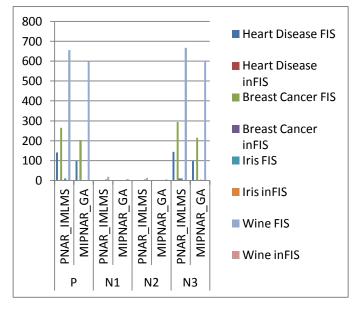


Fig. 6. Shows the comparative value of no of rules and reduces rules of two algorithms by optimization process from Table 4 to table 7

We theoretically proofed a relation between locally large and globally large patterns that is used for pruning at each level to reduce the searched candidates. We derived a locally large threshold using a globally set minimum recall threshold. Pruning achieves a reduction in the number of searched candidates and this reduction has a proportional impact on the reduction of large number of negative rules. In future, some revision might take place to achieve two goals.

*a) Various measures are added to this method for working with grid computing environment.*

*b) To improve the efficiency of the algorithm*

## V. CONCLUSION AND FUTURE WORK

This paper proposed a novel method for optimization of interesting positive and negative association rule. The defined algorithm is combination of MLMS and genetic algorithm. The observation is that when modify the scan process of transaction, generation of rule is fast. With more rules emerging it implies there should be a mechanism for managing their large numbers. The large generated rule is optimized with genetic algorithm.

REFERENCES

[1] By Pengfei Guo Xuezhi Wang Yingshi Han: "The Enhanced Genetic Algorithms for the Optimization Design" 978-1-4244-6498-2/10 © IEEE (2010).

[2] By WEI Yong-Qing, YANG Ren-hua, LIU Pei-yu: "An Improved Apriori Algorithm for Association Rules of Mining" 978-1-4244-3930-0/09/$25.00 © IEEE (2010).

[3] By R. Uday Kiran and P. Krishna Reddy:" An Improved Multiple Minimum Support Based Approach to Mine Rare Association Rules "978-1-4244-2765-9/09/$25.00 © IEEE (2009).

[4] By Sandeep Singh Rawat and Lakshmi Rajamani: "Probability Apriori based Approach to Mine Rare Association Rules".In 3rd Conference on Data Mining and Optimization (DMO), © IEEE (2011).

[5] By Shi-ju SHANG, Xiang-jun DONG, Jie LI, Yuan-yuan ZHAO: "Mining Positive and Negative Association Rules in Multi-database Based on Minimum Interestingness" 978-0-7695-3357-5/08 $25.00 © IEEE (2008).

[6] By XING Xue CHEN Yao WANG Yan-en:"Study on Mining Theories of Association Rules and Its Application" .International Conference on Innovative computing and communication Asia –Pacific Conference on Information Technology and Ocean Engineering 978-0-7695-3942-3/10 $26.00 IEEE (2010).

[7] By LI Tong-yan, LI Xing-ming: "New Criterion for Mining Strong Association Rules in Unbalanced Events" .Intelligent Information Hiding and Multimedia Signal Processing 978-0-7695-3278-3/08 $25.00 © IEEE (2008).

[8] By Xiufend Piao, Zhan long Wang, Gang Liu: "Research on mining positive and negative association rules based on dual confidence" Fifth International Conference on Internet Computing for Science and Engineering. 978-1-4244-9954-0/11 $31 © IEEE (2011).

[9] By Li-Min Tsai, Shu-Jing Lin and Don-Lin Yang: "Efficient Mining of Generalized Negative Association Rules" in International Conference on Granular Computing 978-0-7695-4161-7/10 $ 26 © IEEE (2010).

[10] By CH.Sandeep Kumar, K.Shrinivas, Peddi Kishor T.Bhaskar: "An Alternative Approach to Mine Association Rules" 978-1-4244-8679-3/11 $26.00 © IEEE (2011).

[11] By Dong, X., Niu, Z., Shi, X., Zhang, X., Zhu, D.: Mining both Positive and Negative Association Rules from Frequent and Infrequent Itemsets. ADMA, LNAI 4632, Springer-Verlag Berlin Heidelberg (2007)

[12] By Dong, X., Niu, Z., Zhu, D., Zheng, Z., Jia, Q:" Mining Interesting Infrequent and Frequent Itemsets Based on MLMS Model". The Fourth International Conference on advanced Data Mining and Applications, ADMA (2008).

# De Jong's Sphere Model Test for a Human Community Based Genetic Algorithm Model (HCBGA)

Nagham Azmi AL-Madi

School of Computer Sciences
AL-Zaytoonah Private University
Amman, Jordan

*Abstract*—A new structured population approach for genetic algorithm, based on the custom, behavior and pattern of human community is provided. This model is named the Human Community Based Genetic Algorithm (HCBGA) model. It includes gender, age, generation, marriage, birth and death. Using the De Jong's first function 1, "The Sphere Model" comparisons between values and results concerning the averages and best fits of both, the Simple Standard Genetic Algorithm (SGA), and the Human Community Based Genetic Algorithm (HCBGA) model are obtained. These results are encouraging in that the Human Community Based Genetic Algorithm (HCBGA) model performs better in finding best fit solutions of generations in different populations than the Simple Standard Genetic Algorithm. The HCBGA model is an evolution of the simple Genetic Algorithm (SGA).

The result of this paper is an extended of the result concerning algorithm in [6].

*Keywords—Genetic Algorithms (GAs); Evolutionary Algorithms (EA); Simple Standard Genetic Algorithms (SGA), Human Community Based Genetic Algorithm (HCBGA) model; De Jongs' functions, the Sphere model*

## I. INTRODUCTION

Genetic Algorithms (GAs) were early proposed in the 1960s and 1970s. These search algorithms were initially proposed by Holland, his colleagues and his students at the University of Michigan. GA's are based on nature and mimic the mechanism of natural selection [1, 3, 5, 6, 7, 8, 9].

In his book "Adaptation in Natural and Artificial Systems" [1] Holland initiated GA's as a new area of study. Theoretical foundations besides exploring applications were also presented.

The solution to the problem is represented as a genome (or chromosome) [1, 3, 4, 5, 6] in GAs. The operators such as the crossover and mutation of GA are applied to initialize the population [1, 3, 4, 5, 6]. And with their natural selection they have an iterative procedure usually used to optimize and select the best chromosome (solution) in the population. This population consists of various solutions to hard complex problems and is usually generated randomly [5, 14]. Fig. 1 below represents the Simple Standard GA evolution flow.



Fig. 1. Evolution flow of genetic algorithm [5]

GAs attracted many researchers to search and optimize complex problems. In addition, they proved to be efficient in solving different combinatorial optimization problems. They are considered heuristic search algorithms that solve unconstrained and constrained problems [3]. GAs plays a main role in designing complex devices such as aircraft turbines, integrated circuits and many others [3].

GAs has many advantages in terms of global optimization. On the other hand, from these advantages; potential disadvantages appear [3].

## II. RELATED WORK

John Holland, his colleagues and his students have designed some kind of artificial system software to explain adaptive processes of natural systems [3, 6].

In a certain problem, GA is unaware of the problem itself. It only needs the input parameters. After that, GA represents these inputs in a chromosome format. It differs from other search algorithms in that it has this unique characteristic [3]. This is the reason why GAs can be applied to many types of complex problems [1, 3].

Researches began using GAs to solve some academic problems such as the traveling salesman problem and the 8 Queens problem [3, 5, 6, 9]. Years later, GAs grew rapidly. Applications of GAs were increased to optimize complex scheduling problems and many other types of problems that are hard to efficiently maximize [7].

In the Simple Standard Genetic algorithm parents are selected randomly. There are no constraints in choosing two individuals to mate together [36]. Researchers in this field tried to tackle this problem. They tried to design structured population and control the interaction of the individuals in this population [36].

From many researches on GAs different types and models of GAs appeared such as Cellular GA [36], Island GA [37], Patchwork GA [38, 39], Terrain-Based [40], and religion-Based GA [41]. Below we will discuss some of them briefly.

### A. Cellular GAs (CGA)

By Gorges-Schleuter, 1989 [36]. It is called a diffusion model. A two-dimensional Grid world is used here to arrange the individuals where these individuals interact with each other by the direct neighborhood of each individual [42]. These individuals will be distributed on a graph which is connected together; each individual connects with its neighborhood by a genetic operator. This type of GAs is designed as a probabilistic cellular automation. A self-organizing schedule is added to reproduce an operator [43]. The individual which can interact with its immediate neighbors can only be held in the cell.

### B. Terrain-based GA (TBGA)

TBGA showed better performance than the CGA with less parameter tuning [40]. This was discussed in a previous study [36]. At every generation each individual should be processed, and the mating will be selected from the best of four strings, located above, below, left, right.

It is a more self-tuning model compared to cellular genetic algorithm [40]. In which many combination parameter values will be located in different physical locations.

### C. Island Models (IGA)

According to the increasing complex problems which appear in evolutionary computation, more advanced models of evolutionary algorithms (EAs) appear. Island models are considered a family of such models [45]. Here the individuals are divided into sections. We call each section a subpopulation which is referred to as an island. These island models are able to solve problems in a better performance than standard models [46, 43]. There is a specific relation between islands through some exchange of some individuals between islands. This process is called migration; this is what island models are famous of, and without these migrations, each island is considered as a set of separate run. Therefore migration is very important [47, 45].

### D. Patchwork Model

This type was introduced by Krink et al., (1999). A combination of ideas from cellular evolutionary algorithms, island models, and traditional evolutionary algorithms where used in this model [38, 39]. Here the grid is a two dimensional grid of fields, each field can have a fixed number of individuals. The patchwork model is considered a self-organized, spatial population structure [44]. In a GA population, in order to allow self-adaptation, patchwork model is used as a base. It contains a grid world and some interesting agents. In modeling biological systems the patchwork model is considered as a general approach.

### E. Religion-Based EA Model (RBEA)

It was introduced by Rene Thomsen et al. [44]. The religion-based EA model is based on a part of religious concept which is attracting believers. It attracts new believers to a religion which puts more control than other models such as cellular EA and the patchwork models [41].

### III. HUMAN COMMUNITY BASED GENETIC ALGORITHM (HCBGA) MODEL

AL-Madi and Khader [6] presented a new approach for structured population of GAs so-called Social-Based Genetic Algorithm (SBGA). They applied some constraints on the Simple Standard Genetic Algorithm (SGA) in order to control its randomness in selecting parents. This paper provides a new structured population approach for genetic algorithm, based on the custom, behavior and pattern of human community. This includes gender, age, generation, marriage, birth and death. As such, this model is named the Human Community Based Genetic Algorithm (HCBGA) model. This model is an evolution of the simple Genetic Algorithm (SGA). It is considered an extension to results given in [6].

### A. HCBGA Chromosome Representation

In the HCBGA, the chromosome represents the genome information and additional attributes that would help in simulating human community behavior. In addition, being in the same society- as the population is divided into subgroups or islands- is a dependable constraint for recombination. The problem of age is considered also by adding an attribute for the age. The age attribute takes three values: youth, parent, and grandparent. This chromosome representation (the presence of father and mother pointers) will keep all family relations which divides the subgroups into a Directed Acyclic Graph (DAG).

All the standard operations in the SGA will be changed in order to add restrictions on each operation including: Social constraints such as the Male/Female 'operator', this will be added in the selection part which will restrict choosing two different couples. In addition the Birth *operator* which is generating a new population, and the Death *operator* which will discard the worse individuals.

The development of this new model was carried out in a series of steps. This was done in gradual steps to enable the

measurement of the enhancements to be carried out. The major steps are as illustrated in Fig. 2.

### B. HCBGA Method

Initially, the first individual is selected randomly from the population - this will be the first parent. Based on the first parent's type (whether a male or a female), the second parent will be chosen such that it is the opposite type of the first parent. This process is repeated for a number of individuals creating the initial population. Next come the stages of selection and crossover, bringing up two new children or *offspring's*. Repeating this for a number of couples a second population will be generated.

Again, the previous process is repeated until the maximum number of generations is reached. (The next main important thing is that the *two individuals* must *not* share the same parents).
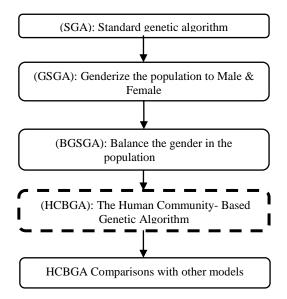


Fig. 2. Development of the HCBGA model

### IV. DE JONG'S FUNCTIONS

De Jong's functions were initially introduced in his thesis entitled "An analysis of the behavior of a class of genetic adaptive systems" [8, 11]. These different functions were used as evaluation functions for the genetic algorithm structure. Many different optimization problems were explained in a novel way using these kinds of functions. This made them the most widely used functions for experimenting Genetic Algorithms functionality and allowing direct comparisons with existing available results [8, 12].

### A. De Jong's function (1): (The Sphere Model)

De Jong's function no. (1) is considered the easiest and simplest test function among De Jong's other functions [10]. It is also called "The Sphere Model". It is a good example of a continuous, strong convex, unimodel function [9, 10].

The structure of the first functions of De Jong functions is defined as follows:

Function definition:

$$f_1(x) = \sum_{i=1}^{n} x_i^2 \qquad -5.12 \le x_i \le 5.12$$

$f_1(x) = $ sum(x(i)^2), i =1:n, 5.12<=x(i)<=5.12.

Global minimum:

f(x)=0, x(i)=0, i=1:n.

The Sphere model serves as a test case for convergence velocity and is well known and widely used in all fields of evolutionary algorithms occurring in the test sets of Schwefel, De Jong, and Fogel [9, 10]. The three-dimensional topology of the Sphere model which shows the Visualization of De Jong's function (1) is shown in Fig. 3 below.



Fig. 3. The Sphere model in a very large area from -500 to 500, [10]

### V. EXPERIMENTAL RESULTS

In this paper we have used the first of De Jong's functions - "The Sphere model" to test the Human Community Based Genetic Algorithm (HCBGA) model. We also used it as a test on the Simple Standard Genetic Algorithm (SGA) in order to compare between both algorithms.

A population size of 350 and a randomly selected one-point crossover are used in a process that is both standard and simple [34]. A random integer (crossover point) and a crossover rate of 50% are chosen according to the maximum length of the chromosome in the model. This is the place in the chromosome at which, with probability, the crossover will occur. If the crossover does occur, then the bits up to the random integer of the two chromosomes are swapped. The mutation of a solution is a random change to a gene value [34, 35]. After several experiments of different mutation rates, the most suitable mutation rate is 0.04. The selection method used is the roulette wheel. The number of generations is 100. The implementation part was programmed in C# (C Sharp) Language Version (5.0) on a Pentium 4, HP-Compaq laptop. This function generates values randomly, whereby the value is restricted to between (-5.12 and 5.12). As mentioned earlier, this was defined in the De Jong's first function.

By applying the Sphere model on both the Simple Standard Genetic Algorithm (SGA) and on the Human Community Based Genetic Algorithm (HCBGA) model we can compare the performance of both algorithms. The comparisons in Figures 4 and 5 below show that the constraints put on the new Human Community Based Genetic Algorithm (HCBGA) model has results in better performance to HCBGA than the Simple Standard Genetic Algorithm (SGA) which depends mainly on its randomness in finding the best fit solution.

It is shown that in the Human Community Based Genetic Algorithm (HCBGA) model the average converge toward the optimal solution better than the Simple Standard Genetic Algorithm (SGA), and the best fit values in the Human Community Based Genetic Algorithm (HCBGA) model also show better findings of best fit values in comparison to the Simple Standard Genetic Algorithm (SGA).

*1) Diversity measurement*

A pair-wise Hamming distance is used in this paper as a measurement to the diversity of the six models SGA, GSGA, BGSGA, HCBGA, CGA and IGA using the Dejong's first test function (f1) problem. This is shown in Figures (4) and (5) respectively.

Fig. 4 illustrates a convergence which occurred in the SGA model, where the curve go down towards the zero x-axis at the second generation. This occurence is considered a fast convergence. This fast convergence indicates a loss of diversity. As a reason this happened due to the existance of similar or identical individuals as all individuals are of same gender in the SGA, in addition there is no constraints when selecting partners to mate any individual could mate with any individual as long as they have high fitness values. This causes a wide possibility of identical individuals to mate producing similar individuals in the next generation and by repeating this process over the generations it leads the search to get stuck around the same solution which causes the algorithm to find a local optimum and fall in a premature convergence.
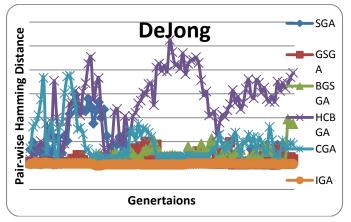


Fig. 4.  Pair-wise Hamming Distance for six models SGA, GSGA, BGSGA, HCBGA, CGA and IGA over the 100th generations.

In the GSGA model, it is seen that the curve converges at the 15th generation, then a small improvment in the pair-wise Hamming distance which gave higher values is indicated from the 16th through the 100th generations. But still the GSGA's

Pair-wise hamming distance has a low value which is near to the zero x-axis. On the other hand, GSGA shows a slower convergence towards the zero compared to the SGA model. This slower convergence occured due to the division to male and females between individuals giving a better oportunity to the GSGA model to search the search space for better solutions than in the SGA. This division applied a better diversity in the GSGA population better than the SGA . But until now there are no restrictions when choosing partners to mate and there is no balance between the number of males and females in the GSGA population which  indicates a loss of diversity as shown in Fig. 4.

The BGSGA model with both the division to male and female and the balance of the individuals in the population to 50% males and 50% females caused a much slower convergence towards the zero x-axis as the values of the pair-wise Hamming distance are getting higher over the generations denoting by this a better diversity in the BGSGA population against the GSGA and the SGA models.

Similarly to the SGA model the CGA and the IGA models have a fast convergence towards the zero x-axis. This is due to both CGA and IGA have common features as in the SGA model whereas there is no existance of sexual gender between the CGA and IGA's population. By this a loss of diversity exists due to the random selection between mates to mate which could be similar or identical and they produce new similar or identical individuals in the next generations.

Relating to the pair-wise hamming distance the HCBGA model shows no convergence towards the zero x-zxis over the generations and gives higher pair-wise hamming distance than the other models which indicates that this model has a better diversity. This is due to the different constraints put on the individuals as the existance of the balanced genderized population made a kind of balanced divers population. Besides that the humman community rules raised in the rules of marriage which restricts mating to a prohibited female gave the HCBGA a huminizing population with a balanced diversity in the population. All together gave the HCBGA algorithm the oppurtunity to search for better solutions in the space of potential solutions maintaining by this the diversity and avoiding falling into a premature convergence.
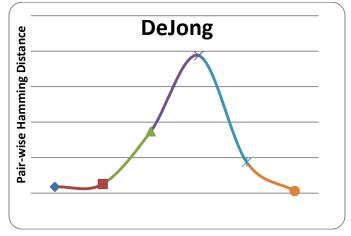


Fig. 5.  Pair-wise Hamming Distance for six models SGA, GSGA, BGSGA, HCBGA, CGA and IGA at the 100th generation.

From Fig. 5, the HCBGA's population diversifies better than the other models whereas, the later shows that the individuals are centralized near zero of x axis due to thier low pair-wise Hamming distance. This means that the distance between the individuals is very short. It indicates that the individuals are almost similar or identical which leads to premature convergence. However, in the HCBGA model, the individuals spread over the search space in a tuning way far away from the x axis as such there are differences between individuals which gives a more divers population and avoids the model to fall in a premature convergence. This is indicated from the high pair-wise Hamming distance between individuals of the HCBGA model as in Fig. 5.

### 2) Statistical Analysis for the models

A statistical analysis has been conducted on the results of the experiments of the paper using the SPSS version 5.0. Table (I) summarizes the results of 20 experiments which compares between 6 models SGA, GSGA, BGSGA, HCBGA, CGA and IGA using the DeJong's first function test problem. In a minimization problem the lower the mean value is the best the model is. Table (I) shows that the HCBGA is the best model based on the lowest mean, standard deviation and variance where its mean value = 1.06 and its standard deviation value = 0.089. If the standard deviation is a high value it means that the individuals don't spread towards the minimum, else the low value of the standard deviation explains the spread of individuals towards the minimization. It is found in Table (I) that the HCBGA model has the lowest standard deviation. This indicates that individuals in the population are spreading around the mean in a balanced distribution. In addition, the variance value of the HCBGA = .008 which is also lower compared to the other models, this indicates a variation in the data, so the HCBGA model has achieved more diversity between its individuals than the other models. By this, the HCBGA model could achieve a better fitness value which means better performance than other models.

TABLE I.     MEAN, STANDARD DEVIATION AND VARIANCE OF THE POPULATION FOR SGA, GSGA, BGSGA, HCBGA, CGA AND IGA MODELS USING THE DEJONG'S FIRST FUNCTION (F1) PROBLEM AFTER 100 GENERATIONS

| Models | No. Generations Statistic | Mean Statistic | Std. Error | Std. Deviation Statistic | Variance Statistic |
|--------|--------|--------|--------|--------|--------|
| SGA | 100 | 3.354 | .0617 | .6167 | .380 |
| GSGA | 100 | 2.03 | .035 | .354 | .126 |
| BGSGA | 100 | 1.371 | .0269 | .2687 | .072 |
| HCBGA | 100 | 1.06 | .009 | .089 | .008 |
| CGA | 100 | 3.975 | .0605 | .6053 | .366 |
| IGA | 100 | 3.630 | .0637 | .6370 | .406 |

TABLE II.     FRIEDMAN TEST SHOWS RANKS BETWEEN SGA, GSGA, BGSGA, HCBGA, CGA AND IGA MODELS

| Models | Mean Rank |
|--------|-----------|
| CGA | 5.35 |
| IGA | 5.20 |
| SGA | 4.45 |
| GSGA | 2.94 |
| BGSGA | 1.92 |
| HCBGA | 1.15 |

The lowest rank in a minimization problem is considered the best. Table (II) shows the mean ranks of the six models; HCBGA model clearly outperforms the other models as it achieved the lowest rank. Since HCBGA yields the best rank against the other models this means that this model has achieved better fitness values in its populations along the 100 generations towards the optimal minimum.

TABLE III.     KENDALL'S W TEST SHOWS SIGNIFICANT DIFFERENCES BETWEEN SGA, GSGA, BGSGA, HCBGA, CGA AND IGA MODELS

| N | 100 |
|---|-----|
| Kendall's W(a) | .898 |
| Chi-Square | 448.932 |
| Df | 5 |
| Asymp. Sig. | .000 |
| Monte Carlo Sig. | .000 |

In Table (III), N is the number of generations the chi-square indicates a test of independence, whereas its value is very high in Table (III) meaning that the HCBGA model is independent from other models. The Df is the degree of freedom its value is k-1 where k is the number of models tested where in this test there is 6 models so the Df value is 5. In addition, the Kendall's W value is .898 which is a high value near to 1, this indicates a full agreement that the HCBGA model performs significantly better in exploring the search space for best solutions than other models. Finally, in Table (III) it shows a Monte Carlo significant value of .000 which means the HCBGA model has a 100% effect and it has a high significant difference over the other models with a level of confidence of 99% due to .000 is less than 5%.

## VI. CONCLUSION

In this paper, a test function of the De Jong's function 1 which is also called "The Sphere Model" is used to evaluate and compare results between the Simple Standard Genetic Algorithm (SGA) and a new approach for structured population of GA called the Human Community Based Genetic Algorithm (HCBGA) model.

It is concluded based on the analysis results that the Human Community Based Genetic Algorithm (HCBGA) model is better in terms of best finding as shown in our given results than the Simple Standard Genetic Algorithm (SGA) and other enhanced models (CGA and IGA).

The Average of the Human Community Based Genetic Algorithm (HCBGA) model is trying to converge towards the minimum despite its restricted constraints to the best values. In addition, the findings of the best solutions of best fit values are better in the Human Community Based Genetic Algorithm (HCBGA) model than in the Simple Standard Genetic Algorithm (SGA).

This model could be considered as a new enhanced model of the SGA. The HCBGA performed better and produced better results in terms of the average of the individuals' fitness as well as the best fit value of individuals in the population, which lead to global optima.

References

[1] J. H. Holland, "Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control and Artificial Intelligence" MIT Press, Cambridge, MA, USA, 1992.

[2] http://ai-depot.com

[3] Y.Zheng, S. Kiyooka,"Genetic Algorithm Applications". www.me.uvic.ca/~zdong/courses/mech620/GA_App.PDF/. 1999.

[4] www.doc.ic.ac.uk/~nd/surprise_96/journal/vol1/hmw/article1.html

[5] Y. Liao, and C.T. Sun, "An Educational Genetic Algorithms Learning Tool" ©2001 IEEE.

[6] N. A. AL-Madi, A. T. Khader, "A SOCIAL-BASED MODEL FOR GENETIC ALGORITHMS". Proceedings of the third International Conference on Information Technology (ICIT), AL-Zaytoonah Univ., Amman, Jordan, 2007.

[7] J. DALTON, "GENETIC ALGORITHMS" NEWCASTLE ENGINEERING DESIGN CENTRE, HTTP://WWW.EDC.NCL.AC.UK. 2007.

[8] E. Bagheri, H. Deldari, 2006. "Dejong Function Optimization by means of a Parallel Approach to Fuzzified Genetic Algorithm" Proceedings of the 11th IEEE Symposium on Computers and Communications (ISCC'06) 0-7695-2588-1/06 $20.00 © 2006 IEEE

[9] T. Back. "Evolutionary Algorithms in theory and practice Evolution Strategies, Evolutionary Programming, Genetic Algorithms". Accessed 2008.

[10] http://www.geatbx.com/docu/fcnindex.html #P74_1604.

[11] K. Dejong, "an analysis of the behavior of a class of genetic adaptive systems" PhD thesis, University of Michigan, 1975.

[12] L. Randy, Haupt S., Practical Genetic Algorithms, 2003, Wiley-IEEE Publication.

[13] T. Back. Evolutionary Algorithms in Theory and Practice, Oxford, New York, 1996.

[14] B. Bhanu, S. Lee, and J. Ming. Self-optimizing image segmentation system using a genetic algorithm. In R. K. Belew and L. B. Booker, editors, Proceedings of the Fourth International Conference on Genetic Algorithms and Their Applications, pages 362-369, San Mateo, CA, July 1991. Morgan Kaufmann.

[15] G. A. Cleveland and S. F. Smith. Using genetic algorithms to schedule flow shop releases. In J. D. Schaffer, editor, Proceedings of the Third International Conference on Genetic Algorithms and Their Applications, pages 160-169, San Mateo, CA, June 1989. Morgan Kaufmann.

[16] N. L. Cramer, A representation for the adaptive generation of simple sequential programs. In J. J, Grefenstette, ed., Proceedings of the First International Conference on Genetic Algorithms and Their Applications. Erlbaum, 1985.

[17] M. Dorigo and U. Schnepf. Genetic-based machine learning and behavior-based robotics: a new syndissertation. IEEE Transactions on System, Man, and Cybernetics, SMC-23(1):144-154, 1993

[18] L. J. Fogel, A. J. Owens, and M. J. Walsh. Artificial Intelligence through Simulated Evolution. John Wiley & Sons, New York, 1966.

[19] C. Fujiki, and J. Dickinson, Using the genetic algorithm to generate Lisp source code to solve the Prisoner's dilemma. In J. J. Grefenstette, ed., Proceedings of the First International Conference on Genetic Algorithms and Their Applications. Erlbaum, 1987.

[20] D. E. Goldberg. Genetic Algorithms in Search, Optimization and Machine Learning. Addison-Weley, Reading, MA, 1989.

[21] D. E. Goldberg and J. Richardson. Genetic algorithms with sharing for multimodal function optimization. In J. J. Grefenstette, editor, Proceedings of the Second International Conference on Genetic Algorithms and Their Applications, pages 41-49, Hillsdale, NJ, July 1987. Lawrence Erlbaum Associates.

[22] J. J. Grefenstette. Credit assignment in rule discovery systems based on genetic algorithms. Machine Learning, 3(2/3):225-245, 1988.

[23] J. J. Grefenstttte, R. Gopal, B. J. Rosmaita, and D. V. Gucht. Genetic algorithm for the traveling salesman problem. In J. J. Grefenstette, editor, Proceeding of the First International Conference on Genetic Algorithms and Their Applications, pages 160-168, Hillsdale, NJ, July 1985. Lawrence Erlbaum Associates.

[24] S. A. Harp, T. Samad, and A. Guha. Towards the genetic syndissertation of neural networks. In J. D. Schaffer, editor, Proceedings of the Third International Conference on Genetic Algorithms and Their Applicatins, pages 360-369, San Mateo, CA , June 1989. Morgan Kaufmann.

[25] K. A. De Jong. Learning with genetic algorithms: an overview. Machine Learning, 3(2/3:121-138, 1988.

[26] C. L. Karr. Design of an adaptive fuzzy logic controller using a genetic algorithm. In R. K. Belew and L. B. Booker, editors, Proceedings of the Fourth International Conference on Genetic Algorithms and Their Applications, pages 450-457, San Mateo, CA, July 1991. Morgan Kaufmann.

[27] K. Kristinsson and G. A. Dumont. System identification and control using genetic algorithms. IEEE Transactions on System, Man, and Cybernetics, SMC-22(5):1033-1046, 1992.

[28] J. R. Koza, Genetic Programming: On the Programming of Computers by Means of Natural Selection. MIT Press, 1992.

[29] J. R. Koza, Genetic Programming II: Automatic Discovery of Reusable Programs. MIT Press, 1994.

[30] G. G. Miller, P. M. Todd, and S. U. Hegde. Designing neural networks using genetic algorithms. In J. Schaffer, editor, Proceedings of the Third International Conference on Genetic Algorithms and Their Applications, pages 379-384, San Mateo, CA, June 1989. Morgan Kaufmann.

[31] N. H. Packard, A genetic learning algorithm for the analysis of complex data Complex Systems 4, no. 5:543-572, 1990.

[32] G. Syswerda and J. Palmucci. The application of genetic algorithms to resource scheduling. In R. K. Belew and L. B. Booker, editors, Proceedings of the Fourth International Conference on Genetic Algorithms and Their Applications, pages 502-508, San Mateo, CA, July 1991. Morgan Kaufmann.

[33] M. Mitchell, An Introduction to Genetic Algorithms. MIT Press, 1997.

[34] E. K. Prebys,The Genetic Algorithm in Computer Science, MIT Undergraduate Journal of Mathematics, ee.sharif.ir/~poshtkoohi, accessed 2007.-

[35] www.geocities.com, "A Genetic Knapsack Problem Solver". Accessed 2007.

[36] M. Gorges-Schleuter. "ASPARAGOS an asynchronous parallel genetic optimization strategy". In J. D. Schaffer, editor, Proceedings of the Third International Conference on Genetic Algorithms, 1989, pp. 422–427.

[37] Back, T., Fogel, D. B., Michalewicz, Z., Et Al., Eds. Handbook on "Evolutionary Computation". IOP Publishing Ltd and Oxford University Press, 1997.

[38] Krink, T., Mayoh, B. H., and Michalewicz, Z. "A Patchwork model for evolutionary algorithms with structured and variable size populations". In Proceedings of the Genetic and Evolutionary Computation Conference, Vol. 2, 1999, pp. 1321-1328.

[39] Krink, T., and Ursem, R. K. "Parameter control using the agent based patchwork model". In Proceedings of the Second Congress on Evolutionary Computation (CEC-2000). 77-83.

[40] Gorden, V. S., Pirie, R., Wachter, A., and Sharp, S. "Terrain-based genetic algorithm (TBGA): Modeling parameter spaceasterrain".In Proceedings of the Genetic and Evolutionary Computation Conference, vol. 1, 1999, pp. 299- 235.

[41] Rene Thomsen, Peter Rickers and Thiemo Krink. "A Religion-Based Spatial Model For Evolutionary Algorithms". EvAlife Group, Dept. of Computer science, University of Aarhus, Denmark [Online] [Accessed 2006].

[42] Whitley, D. "Cellular genetic algorithms". In Proceedings of the Fifth International Conference on Genetic Algorithms, Forrest S. (ed.). Morgan Kaufmann 1993, pp. 658.

[43] Enrique Alba, Mario Giacobini and M. Tomassini,: "Comparing Synchronous and Asynchronous Cellular Genetic Algorithms". In J. J. Merelo et al., editor, Parallel Problem Solving from Nature – PPSN V11, Springer-Verlag, Heidelberg, 2002, 2439: 601-610.

[44] Ursem, R. K. "Multinational evolutionary algorithms". In Proceedings of the Congress of Evolutionary Computation, Vol. 3, 1999, pp. 1633-1640.

[45] Zbigniew Skolicki, and K. De Jong. "The influence of migration sizes and intervals on island models". In Proceedings of Genetic and Evolutionary Computation Conference – GECCO-2005. ACM Press, 2005, 1295-1302.

[46] Mohamed A. Belal, Iraq H. Khalifa, "A Comparative Study between Swarm Intelligence and Genetic Algorithms", Egyptian Computer Science Journal, Vol. 24, No. 1, 2002.

[47] In: 2002 Water Resources Planning & Management Conference, Roanoke, VA. American Society of Civil Engineers (ASCE) Environmental & Water Resources Institute (EWRI), 2002.

# DES: Dynamic and Elastic Scalability in Cloud Computing Database Architecture

Dr.K.Chitra

Dept of Computer Science

Govt Arts College

Melur, Madurai Dt, TN, India

B.Jeeva Rani

Dept of Computer Science

Research Scholar, Bharathiar University

Coimbatore, TN, India

*Abstract*—**Nowadays, companies are becoming global organizations. Such organizations do not limit themselves in conducting business in one country. They need dynamic, elastic, scalable cloud computing platform that operates around-the-clock. Full functionality, adaptability, non-stop availability and reduced cost are the major requirements that are expected from cloud computing services. Planned or unplanned system outages are the enemies of the successful business in cloud computing environment. Hence it requires highly available, elastic scalable systems. In this paper, we analyze the benefits of cloud computing and evaluates the database architectures namely shared-disk database architecture and shared-nothing database architecture for high availability, Dynamic and Elastic Scalability.**

*Keywords*—*Cloud computing; Availability; Scalability; Shared-disk; Shared-nothing*

## I. INTRODUCTION

Accessing the remotely available computing resources like hardware and software, over a network is termed as Cloud computing. Cloud computing is a collection of integrated and networked hardware, software and Internet infrastructure called a platform. Using the Internet for communication and transport, it provides hardware, software and networking services to the clients. Users need utilities like computing power and storage on demand. Cloud - the name comes from the cloud-shaped symbol used in system diagrams. Cloud computing systems provide services with a user's data, platform for software development and computation.

Cloud computing environment requires platform which supports the key design principles of the cloud architecture [1]. One of the core design principles is dynamic scalability i.e. the ability to add and remove servers on demand. Unfortunately, the majority of available database servers are unable to satisfy this key requirement [2]. This paper analyzes the benefits of cloud computing and evaluates the database architectures namely shared-disk database architecture and shared-nothing database architecture for their applicability with cloud computing. Section 2 lists out the benefits of cloud computing. Section 3 explains the inconsistency in cloud database. Section 4 analyzes the shared-nothing and shared-disk database architectures of cloud databases.

## II. BENEFITS OF CLOUD COMPUTING

Cloud computing is driven by tangible and powerful benefits. Cloud computing is also known as "Elastic Computing". It should be noted that the underlying database is not very elastic and scalable. The following are the features of cloud environment:

*1) Computing power is elastic only if workload is parallelizable*

*2) Data is stored at an untrusted host*

*3) Data is replicated across large geographic distances*

*4) Hard to maintain ACID while data is replicated over large geographic distances*

In a traditional IT model, each development effort needs expertise on staff. But Cloud computing model enables development to be staffed by experts and these services are accessed by large number of customers [2].

The ideal platform needs servers, networking equipment, data storage/backup, power, redundant high-speed connectivity etc. But they are very expensive and can result in a huge start-up cost to develop a single project and this effort may also fail [1]. But with Cloud computing the same investment can be utilized over a large number of projects [4].

In an IT company, separate person is needed to manage and schedule the backup process. But with cloud, backup is highly automated for thousands of customers.

### B. Key Benefits of Cloud Computing

*1) Reduced IT operating costs (25%)*

*2) Shifting Capital Expenses to Operating Expenses*

*3) Increased efficiency (55%)*

*4) Agility*

*5) Dynamic scalability*

*6) Enabled us to offer new products/services (24%)*

*7) Simplified maintenance*

*8) Large scale prototyping/load testing*

*9) Increased ability to innovate (32%)*

*10) Diverse platform support*

*11) Faster management approval*

*12) Improved employee mobility (49%)*

*13) Freed current IT staff for other projects (31%)*

*14) Faster development*

With the combination of the above listed benefits of cloud computing, we are seeing an explosion of cloud options. Most importantly the Database-as-a-Service (DaaS) - provisioning of database services in the form of cloud databases. The following sections of this paper will focus on the requirements of cloud databases and the various database architectures.

### III. CLOUD DATABASE

Business applications demand that the cloud database be ACID (Atomicity, Consistency, Isolation and Durability) compliant. But the cloud database architecture provides Scalability, Availability, but weak form of consistency. The following are the two example scenarios that explain the major problems faced by the companies due to inconsistency.

#### A. Example I: Customer Cloud Access

Consider a customer-centric music instruments website with DB access. The user makes a search for a violin and the user must be given instant access to the data orelse he/she may jump to other instrument sites. If the data given to the user said that the chosen violin is in inventory and sale is done. It happened because of data inconsistency, but violin is not really in inventory. Hence, the customer is notified that it is on backorder and the violin will be shipped soon.

#### B. Example II: DB in Corporate Cloud

Consider a company that sells raw materials to manufacturers. A huge company makes a purchase of a load of raw materials which are required to keep its production line running. If the inventory database is incorrect because of inconsistent data, the shipment is delayed. So the company who purchased the raw materials is forced to shut down a production line at a cost.

### IV. DYNAMIC SCALABILITY IN CLOUD DATABASE ARCHITECTURE

Cloud computing databases are dynamically and elastically scalable. Hence they can provide around-the-clock data services to the clients. Now we discuss about the two types of database architectures that can be used in cloud computing platform.

1) *Shared-Nothing Database Architecture*
2) *Shared-Disk Database Architecture*

#### B. Shared-Nothing Database Architecture

Dynamic scalability is one of the basic principles of cloud computing. In the shared-nothing DB architecture, data is partitioned and kept one partition per server [3] i.e. if there are two servers, 50% of the data is kept in each. When we add a third server, 33% of the data is maintained per server as shown in Figure 1.

It becomes a tedious and time consuming process.

- When many users request the same data, one server will be fully loaded.



Fig. 1. Shared-Nothing Architecture

- When user requests related information, we need to check for the data in many servers. It is very big time consuming process and degrades the system.
- Automating the partitioning process remains an indefinite goal.
- Repartitioning make the performance of the system poor.

Hence the performance is decreased with the addition of more servers [3]. These are the current problems faced while deploying a shared-nothing database in the cloud [5]. Amazon, Facebook and Google use a persistence engine that uses replicated tables that store and retrieve information. This approach supports elastic or dynamic scalability.

#### C. Shared-Disk Database Architecture

The shared-disk database architecture is illustrated in Figure 2. It eliminates the need to partition data. It uses group of low-cost servers with single collection of data. Without partitioning, the whole data is placed in all of the servers. It is a master-master configuration. If a node fails, any other node provides service [3]. Hence it supports elastic and dynamic scalability, high-availability, reduced maintenance.



Fig. 2. Shared-Disk Architecture

In addition to the vital feature 'elastic scalability', the shared-disk database architecture has other important advantages also that make it very much suitable for deployment in the cloud computing environment.

The following are some of the advantages of shared-disk architecture:
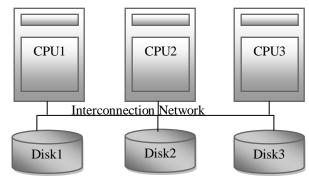
1) *Fewer servers required*
2) *Lower cost servers*
3) *Scale-in*
4) *Simplified maintenance/upgrade process*
5) *High-availability*
6) *Reduced partitioning and tuning services*
7) *Reduced Support Costs*

### D. Shared-Nothing Vs Shared-Disk Database Architecture

The following Table 1 compares the Shared-Disk and Shared-Nothing Database Architectures in using the parameters like configuration, number of servers required, Utilization of CPU, the cost of servers, Scalability, Maintenance, Availability, Partitioning, Support cost, downtime, Adaptability and performance.

## V. CONCLUSION

Cloud computing platform needs cloud compatible database architecture. This paper has analyzed the database architectures that can be used with cloud computing environment. Shared-Disk Architecture is write-limited because the related locks must coordinate across the cluster. Shared-Nothing architecture is useful for systems that need high-throughput writes when you shard your data and must be clear about the transactions that span different shards. If you require scalability over consistency, you can use shared-nothing architecture. For a business system with two or three servers, Shared-Disk is the best option. Dynamic scalability is the core principle of cloud computing. But, Shared-nothing databases architecture does not support dynamic scalability. On the other hand, Shared-disk database architecture supports elastic scalability and provides high-availability.

TABLE I.  COMPARISON OF SHARED-DISK AND SHARED-NOTHING DB ARCHITECTURE

|  | Shared-Disk Architecture | Shared-Nothing Architecture |
|---|---|---|
| **Configuration** | Master-master configuration | Master-slave configuration |
| **Servers Required** | Less no of servers | More no of servers |
| **CPU Utilization** | High | Low |
| **Cost of Servers** | Lower-cost servers | Expensive servers |
| **Scalability** | High | Almost unlimited scalability |
| **Maintenance** | Upgraded individually, cluster remains online | Entire cluster must be shut down |
| **Availability** | High | Low |
| **Partitioning** | No partitioning | Data partitioned across cluster |
| **Support Cost** | Reduced | High |
| **Downtime** | Less but scheduled only | More scheduled and unscheduled downtime |
| **Adaptability** | Quick adaptability to changing workloads | Less |
| **Performance** | Performs best in a heavy read environment (for example, a data warehouse) | Works well in a high-volume, read/write environment |

REFERENCES

[1] Dr.K.Chitra, B.Jeevarani, "Study on Basically Available, Scalable and Eventually Consistent NOSQL Databases", International Journal of Advanced Research in Computer Science and Software Engineering, Vol 3 (4), July - 2013, ISSN: ISSN: 2277–128x, pp. 1-5.

[2] Dr.K.Chitra, B.Jeevarani, "Transaction Processing on Cloud Computing", International Journal of Scientific and Research in Computer Science Applications and Management Studies, Volume2 Issue 5, ISSN: 2319-1953, November 2013.

[3] Craig S. Mullins, Architectures for Clustering: Shared Nothing and Shared Disk, 2003.

[4] Divyakant Agrawal, Amr El Abbadi, Sudipto Das, and Aaron J. Elmore, Database Scalability, Elasticity, and Autonomy in the Cloud.

[5] Mike Hogan, Scaling-up, Scaling-out & Scaling-in, 2009.

[6] Ben, Distributed Data Storage, www.benstopford.com, Tuesday, November 24th, 2009.

# Wireless LAN Security Threats & Vulnerabilities:

## A Literature Review

Md. Waliullah

Dept. of Computer Science & Engineering,
Daffodil International University,
Dhaka, Bangladesh

Diane Gan

School of Computing and Mathematical Science
University of Greenwich,
London, UK

*Abstract*—**Wireless LANs are everywhere these days from home to large enterprise corporate networks due to the ease of installation, employee convenience, avoiding wiring cost and constant mobility support. However, the greater availability of wireless LANs means increased danger from attacks and increased challenges to an organization, IT staff and IT security professionals. This paper discusses the various security issues and vulnerabilities related to the IEEE 802.11 Wireless LAN encryption standard and common threats/attacks pertaining to the home and enterprise Wireless LAN system and provide overall guidelines and recommendation to the home users and organizations.**

*Keywords—WLAN; IEEE 802.11i; WIDS/WIPS; MITM; DoS; SSID; AP; WEP; WPA/WPA2*

## I.    INTRODUCTION

Over the last twelve years, 802.11 Wireless LAN's have matured and really reshaped the network landscape. 802.11n is now rapidly replacing Ethernet as the method of network access. The rapid proliferations of mobile devices has led to a tremendous need for wireless local area networks (WLAN), deployed in various types of locations, including homes, educational institutions, airports, business offices, government buildings, military facilities, coffee shops, book stores and many other venues. Besides, the facilities of flexibility and mobility of wireless devices has been attracted by most organizations and consumers all over the world.  Low cost of hardware and user friendly installation procedures allow anyone to set up their own wireless network without any specialist knowledge of computer networks.

However, the increased development of Wireless LAN has increased the potential threats to the home user, small businesses and the corporate world. Unlike a wired network, a WLAN uses radio frequency transmission as the medium for communication. This necessarily exposes layer 1 and layer 2 to whoever can listen into the RF ranges on the network. Wireless insecurity has been a critical issue since Wired Equivalent Privacy (WEP), an IEEE standard security algorithm for wireless networks, was compromised [1]. To address the significant security flaws in the WEP standard, the Wi-Fi alliance developed the 802.11i standard, called Wi-Fi Protected Access (WPA) and WPA2 [1]. However, many researchers have shown that the IEEE 802.11i standard cannot prevent eavesdropping, various denial of service attacks including de-authentication and disassociation attacks. Moreover, 802.11i's pre-shared key mode of WEP for flexibility and backward compatibility has made it easier for most hackers to perform a Dictionary and Brute force attack [2].

Recently, a scanning experiment based on London conducted by the security firm Sophos has revealed that more than one in four Wi-Fi networks in London are poorly secured or not secured at all [3]. Of 100,000 Wi-Fi networks detected on a 90 Km route, 8% of the Wi-Fi networks detected used no encryption at all. This figure excludes intentionally open networks such as coffee shops, hotels and Wi-Fi hotspots. Approximately, 9% of Wi-Fi networks detected were using default network names such as "default" or a supplier name enabling the hacker to break passwords more easily. More importantly, the experiment revealed that 19% of the Wi-Fi networks detected used obsolete WEP as the encryption standard which has already proved to be easily cracked within a second, using readily available hacking tools [3]. So, the security of a wireless LAN still remains the top concern in the home and corporate network.

This paper discusses the vulnerabilities and security issues pertaining to the IEEE 802.11 security standard and describes major well known attack/threats to the home and enterprise wireless LAN system. The remainder of the paper is organised as follows. A brief overview of WLANs are outlined in section II. Related work is presented in section III. The common vulnerabilities and security issues pertaining to the IEEE 802.11 security standard and WLAN are discussed in section IV. This is followed by an over view of the common threats/attacks on WLAN technology. Common guidelines and an overall recommendation is presented in section VI, and a conclusion is outlined in section VII.

## II.    OVERVIEW OF WLAN

An access point (AP) and a network interface card (NIC) are the two basic components of a WLAN. An AP typically connects the wireless clients or stations to each other by means of antenna and then connects to the wired backbone through a standard Ethernet cable. A NIC normally connects a wireless station to the AP in the Wireless LAN [4].  Any devices that have the ability to communicate with 802.11 networks are called a station i.e. laptops, printers, media servers, smartphones, e.g. IPhones, Windows mobile handsets, VoIP phones etc. All 802.11 stations operate in two ways, either in ad-hoc mode, where stations are connected to each other, or in infrastructure mode, where stations are communicating with each other via the access points to reach some other network [5].

Companies install as many access points as it takes to cover an entire building or even a campus. The whole network is configured with the same network name to act as one huge wireless network, which is called an extended service set identifier (ESSID) or a standard service set identifier (SSID). For example, if an IPhone wants to connect to a WLAN, it starts by scanning all channels; sends a probe request and listens for beacon frames that are sent by access points to advertise themselves. Then, it compares all those beacons and probe responses to the desired SSID and selects the best available access point.

Finally, the IPhone will send an authenticating packet and will associate the request to that access point by establishing a 4-way handshake mechanism. Once the IPhone is associated and authenticated it can send and receive data using that wireless network [5].

## III. RELATED WORK

Sheldon, Weber, Yoo and Pan [1] described how the wireless LAN encryption standards such as WEP, WPA/WPA2 are vulnerable to attack. They presented some of the attacks on encryption standards such as Chop-chop attack, Brute force, Beck-Tews, Halvorsen-Haugen and the hole 196 attacks etc. Wang, Srinivasan, and Bhattacharjee [2] proposed a 3-way handshake model instead of the usual 4 way handshake method for the 802.11i protocol. They suggested how their alternative method can effectively prevent denial of service (DoS) attacks including de-authentication, disassociation and memory/CPU DoS attacks. Souppaya and Scarfone [6] discussed the need for security concerns and these should be applied from the configuration design stage to implementation and evaluation through to the maintenance stage of the WLAN. They provided some general guidelines and recommendations in order to reduce the vulnerabilities and prevent the most common threats.

Pan Feng [4] suggested that more than 70% of the WLAN security issues are due to human factors, such as data theft by acquaintances or colleagues. He addressed that remaining 30% of security threats are technology related. Reddy, Rijutha, Ramani, Ali, and Reddy [7] demonstrated how WEP can be cracked by freely available open source software tools such as Netstumbler, Ministubler, Airopeek, Kismat, Cain etc. They have mainly focused on securing WLANs by realizing miscellaneous threats and vulnerabilities associated with 802.11 WLAN standards and have used ethical hacking to try to make these more secure.

Li and Garuba [8] and Deng Shiyang [9] discuss various encryption standards relating to 802.11 WLAN, their vulnerabilities and security flaws. Stimpson et al [10] describes war driving techniques as a useful tool for assessing security and vulnerabilities of home wireless networks.

However, none of the above researchers has elaborately presented WLAN security vulnerabilities, threats and general guidelines/recommendations for securing them. Realizing the vulnerabilities, understanding the most common threats and providing general guidelines and recommendation in order to protect WLAN network and make them more secure for the home user and for enterprise networks is the aim of this paper.

## IV. WLAN VULNERABILITIES

Wireless LANs have gained much more popularity than wired networks because of their flexibility, cost-effectiveness and ease of installation. However, the increasing deployment of WLANs presents the hacker or cracker with more opportunities. Unlike wired networks, WLANs transmit data through the air using radio frequency transmission or infrared. Current wireless technology in use enables an attacker to monitor a wireless network and in the worst case may affect the integrity of the data. There are a number of security issues that presents the IT security practitioner, system administrator securing the WLAN with difficulties [11].

As the name implies Wired Equivalent Privacy (WEP) was intended to provide users with the same level of privacy as that of a wired LAN. However, when this protocol was first developed by the IEEE 802.11b Task Force in 1999, it quickly proved to be less secure than its wired equivalent. WEP comes as 64 bit or 128 bit but the actual transmission keys are 40 bits and 104 bits long. In each case the other 24 bits is an Initialization Vector (IV). Before transmission, the packets are encrypted with a symmetric encryption algorithm (RC4) using a session key which is made up of the IV and the default transmit key. The IV is randomly generated for each session but the default transmit key is fixed. The IV is sent in the packet along with the data. Once the encrypted packet reaches the receiving end, it decrypts the packet using the same session key [12].

However, WEP has some serious security problems. It fails to meet the fundamental security goals of confidentiality, integrity and authentication. The main problem with WEP is that the 40 or 104 bit keys are static and common to all users in the WLAN. Since, WEP does not provide an effective key management technique, changing the keys on all devices is a time consuming and difficult task. Thus, if any devices are lost or stolen, the higher the chances of the key being compromised. This exposes the whole system to security breaches [12]. More importantly, the encryption algorithm RC4 used in WEP is flawed and encryption keys can be recovered through cryptanalysis [8]. Besides the default transmission key, the IV is short and can be easily sniffed by passive attack using freely available software tools. One of the other problems is that WEP is disabled by default and its use is optional, therefore, many users never turn on encryption. It is better to use of some form of encryption than no encryption at all [8][12] .

In order to eliminate all well-known attacks and address the significant security flaws in WEP, the Wi-Fi alliances developed IEEE 802.11i security standard in 2004 which is called Wi-Fi protected access (WPA) and subsequently WPA2. WPA uses the same encryption algorithm (RC4) used in WEP but improved by the use of a 48 bit temporary key integrity protocol (TKIP) sequence counter (TSC) instead of WEP's 24 bit key. Moreover, the 64 bit message integrity check (MIC) algorithm named Michael is used to ensure integrity [1]. Furthermore, to improve user authentication and access control, WPA uses the extensible authentication protocol (EAP) and the IEEE 802.1x standard port based access control. This method uses the Radius (Remote Authentication Dial-in User Service) server to authenticate each user on the network [8]. In the

absence of a Radius server, it uses a pre-shared key (PSK), which is called WPA-PSK or WPA-Personal and is mostly used in small-offices and by home users [1].

Although, WPA is considered stronger than WEP, it does reuse the WEP algorithm. As a result it is vulnerable to offline dictionary and brute force attacks against the 4-way handshake protocol [10]. More importantly, it is much more vulnerable to DoS attacks which are carried out over the MAC layer by sending out de-authentication and disassociation messages to the client or AP resulting in the legitimate user being denied access to the service [8].

WPA2 employs the Counter Mode Cipher Block Chaining Message Authentication Code Protocol (CCMP) instead of TKIP and uses Advance Encryption Standard (AES) block cipher. AES replaces the WPA's RC4 stream cipher [1]. Although WPA2-AES is still regarded as extremely secure, it is vulnerable to DoS, offline dictionary and internal attacks and fails to provide the availability aspect of the CIA triad [2].

More importantly, the robust encryption standard only applies to data frames and not currently to the management frames. All 802.11 management and control frames are vulnerable to replay or forgery, including the messages that are used to probe, associate, authenticates, disassociate, and de-authenticate users from the WLAN. Besides, unlike the software upgrade required for migration from WEP to WPA, WPA2 requires the replacement of older hardware, extra processing power and has a much higher cost [12].

Both WPA and WPA2, are extremely vulnerable to dictionary and bruit-force attack, regardless of whether they are operating in Personal or Enterprise mode,. Most home networks use pre-shared key authentication, allowing quick and easy control over who can use the network. This PSK's are both simple and limited and they are the same as a group password. They can be shared with outsiders, or the device can be lost or stolen. As a result, it is hard to guess whether the WLAN being used by legitimate users or foes. Furthermore, the attacker can still listen to frames that are being sent and received without even trying to authenticate even in the case of 802.1x port based authentication. In addition 802.1x's lightweight EAP or LEAP, implement password authentication in a way that is vulnerable to a dictionary attack [5].

802.11 networks are inherently vulnerable to radio frequency interference problems. Most of the wireless LAN standards operate on the 2.4 GHz channel frequency band, while many other devices such as Bluetooth, cordless phones and microwave signals also operate on the same frequency band. This can lead to signal interference and cause a legitimate user to be disconnected [4].

Our inability to effectively contain radio signals makes the WLAN vulnerable to a different set of attacks from wired LANs. Although businesses can position their access points and use antennas to focus their signals in a specific direction, it is hard to completely prevent wireless transmission from reaching an undesirable location like nearby lobbies, semi-public areas and parking lots. This makes it easier for intruders to sniff sensitive data [5].

MAC address filtering can be configured in an access point in order to allow only an authorized client in the network. However, the various available open source hacking tools i.e. Kismet, SMAC etc. can be used to passively sniff a large amount of network traffic, including the MAC addresses of authorized computers. These can then be changed to act as legitimate clients on the network. Moreover, in a large network the continually updated list of MAC address at the access point sometimes creates a security hole, if the list is not correctly updated [13].

SSID is an identification that allows the clients to communicate with an appropriate access point. The available access points on the market come with a default SSID name and password. This creates potential security vulnerabilities, if these are not changed by the administrator or user. For example some of the common default passwords are: "tsunami" (Cisco), "101" (3Com), "Compaq" (Compaq) etc. Furthermore, most hotspots and guest networks operate in an open system mode allowing any stations to connect to that network without requiring any form of authentication [14].

## V. GENERAL ATTACKS/THREATS TO WLAN TECHNOLOGY

An attack is an action that is carried out by an intruder in order to compromise information in an organization. Unlike wired networks, a WLAN uses radio frequency or infrared transmission technology for communication; thus, making them susceptible to attack. These attacks are aimed at breaking the confidentiality and integrity of information and network availability. Attacks are classified into the following two categories:

- Passive attacks.

- Active attacks

Passive attacks are those types of attack in which the attacker tries to obtain the information that is being transmitted or received by the network. These types of attacks are usually very difficult to detect as there is no modification of the contents by the attacker [15]. There are two types of passive attack and these are traffic analysis and eaves dropping.

On the other hand, active attacks where the attacker not only gains access to the information on the network but also changes the information/contents or may even generate fraudulent information on the network. This type of malicious act, results in great loss for any organization [15]. Following are a list of active attacks in WLAN technology:

- Unauthorized Access

- Rogue Access Point

- Man in the Middle Attack (MITM)

- Denial-of-Service

- Reply Attack

- Session High jacking

According to the CIA triad, information security should meet three main principles, which are confidentiality, integrity and availability. All three concepts are needed to some extent

to achieve true security. Otherwise, the network will be vulnerable to attack. Furthermore, two other principals involved i.e. access control and authentication.

- Confidentiality is the prevention of intentional/unintentional disclosure of data.

- Integrity is control over the intentional/unintentional modification of data.

- Availability is the control over provision of system resources on demand to authorized users/systems/processes.

- Access control is the control of access to the resources by a legitimate user.

- Authentication is the process by which a system verifies the identity of a user who wants to access it [16].

Based on the CIA triad, access control and the authentication definitions described, various types of attack/threats in a WLAN are discussed below. These attack categories can also fall in the above active or passive types.

A. *Confidentiality Attacks*

In this type of attack, intruders attempt to intercept highly confidential or sensitive information that has been sent over the wireless association either encrypted or in clear text by the 802.11 or higher layer protocols. Examples of passive attacks are Eavesdropping, Man-in-the-Middle attack, Traffic Analysis etc. Active attack categories are WEP Key Cracking, Evil Twin AP and AP Phishing etc. [16].

*1) Traffic Analysis: Also known as footprinting, is the first step which is carried out by most hackers before attempting further attacks. This is a technique whereby the attacker determines the communication load, the number of packets being transmitted and received, the size of the packets and the source and destination of the packet being transmitted and received. Thus, the overall network activity has been acquired by the traffic analysis attack [17]. To accomplish this attack, the attacker uses a wireless card that can be set to promiscuous mode and special types of antenna to determine the signal range e.g. yagi antenna, along with the global positioning mode (GPS). Furthermore, there are a number of freely available software that can be used e.g. Netstumbler, Kismet etc.*

The intruders obtain three forms of information through traffic analysis. First, they identify if there is any network activity on the network. Secondly, he or she identifies the number of access points and their locations in the surrounding area. If the broadcast SSID has not been turned off in the AP, then it broadcasts the SSID within the wireless network in order to allow wireless nodes to get access to the network. Even if it is turned off, a passive sniffer like Kismet can obtain all the information about the network including the name, location and the channel being used by any AP. Finally, the third piece of information the attacker can learn through traffic analysis is the type of protocol that is being used in the transmission, along with the size, type and number of packets

being transmitted. For example, analysis of the three-way handshake information of TCP [17].

*2) Eavesdropping: An Eavesdrop attack, enables an attacker to gain access to the network traffic and read the message contents that are being transmitted across the network. The attacker passively monitors the wireless session and the payload. If the message is encrypted, the attacker can crack the encrypted message later. The attacker can gather information about the packets, specially their source, destination, size, number and time of transmission. More importantly, there are many directional antennas available in the market which can detect 802.11 transmissions under the right conditions, from miles away. This is an attack that cannot be easily prevented using adequate physical security measures. Besides, this attack can be done far away from the premises of any organizations [17][18].*

*3) Man-in-the-Middle Attack: A man-in-the-middle attack can be used to read the private data from a session or to modify them, thus, breaking the confidentiality and integrity of the data. This attack also breaks indirect data confidentiality. However, an organisation could employ security measures such as a VPN or IPsec, which only protect against direct data confidentiality attacks. This is a real time attack which occurs during the target machine's session. There are multiple ways to implement this attack. For example, in step one, the attacker breaks the target's client session and requires them to re-associate with the access point. In step two, the target client attempts to re-associate with the access point but can only re-associate with the attacker's machine, which is mimicking the access point. In the meantime, the attacker associates and authenticates with the access point on behalf of the target client. If an encrypted tunnel is in place, the attacker establishes two encrypted tunnels, one between it and the target client and another to the access point. In short, in this type of attack, the attacker appears to be an AP to the target client and a legitimate user of the AP [17].*

*4) Evil Twin AP: An Evil Twin attack poses as great a danger to wireless users on public and private WLANs alike. In this type of attack, an attacker sets up a phony access point in the network that pretends to be a legitimate AP by advertising that WLAN's name i.e. extended SSID. Karma is an attack tool that is used to perform this attack by monitoring station probes, watching commonly used SSIDs and using them as its own. Even APs that do not send SSIDs in the beacon can also be accessed using NetStumbler, Kismet or another WLAN analyzer tool while posing as a legitimate user [16].*

B. *Access Control Attacks*

This attack attempts to penetrate a network by bypassing the filters and firewall to gain unauthorized access. war driving, rogue access points, MAC address spoofing and unauthorized access are the most common types of attack in this category.

*1) War Driving: While war driving, the attacker drives around in a car with a specially configured laptop that has software such as Netstumbler or Kismet installed which identifies the network characteristics. More importantly, an external antenna and a GPS can be used to clearly identify the location of a wireless network [19]. The attacker discovers wireless LANs i.e. all the APs, the physical location of each AP, the SSID and the security mechanisms etc. by listening to the beacon or by sending a probe request. This attack provides the launch point for further attacks [19].*

*2) Rogue Access Point: In this type of attack, an intruder installs an unsecured AP usually in public areas like airports, shared office areas or outside of an organization's building in order to intercept traffic from valid wireless clients, to whom it appears as a legitimate authenticator. As a result, this attack creates a backdoor into a trusted network. The attacker could fool the legitimate client by changing its SSID to the same as that used by the target organization. Furthermore, the attacker uses an unused wireless channel to set up this fake access point. It is easy to trick unsuspecting users into connecting to the fake access point. Thus, the credential information of a user could easily be stolen [20][21].*

*3) MAC addresses spoofing: In this type of attack, the attacker gains access to privileged data and various resources such as printers, servers etc. by assuming the identity of a valid user in the network. To do so, the attacker reconfigures their MAC address and poses as an authorized AP or station. This could be easily done, because 802.11 networks do not authenticate the source MAC address frames. Therefore, the attacker can spoof MAC addresses and hijack a session. Furthermore, 802.11 does not require an AP to prove it is a genuine AP [14].*

*4) Unauthorized Access: Here the attacker is not aiming at a particular user, but at gaining access to the whole network. The attacker can gain access to the services or privileges that he/she is not authorized to access. Moreover, some WLAN architecture not only allows access to the wireless network but also grants the attacker access to the wired component of the network. This can be done by using war driving, rogue access points or MAC spoofing attack. This attack gives the attacker the ability to do a more malicious attack such as a MITM [17].*

### C. Integrity Attacks

An Integrity attack alters the data while in transmission. In this attack, the intruder tries altering, deleting or adding management frames or data i.e. forged control packets to the network, which can mislead the recipient or facilitate another type of attack [22]. DoS attacks are the most common example of this type of attack which is described in section D. Other types include session hijacking, replay attacks, 802.11 frame injection, 802.11 data replay, and 802.11 data deletion etc.

*1) Session Hijacking: In Session Hijacking, an attacker takes an authorized and authenticated session away from the legitimate user of the network. The legitimate user thinks that the session loss may be a normal malfunction of the WLAN.*

*Thus, he/she has no idea that the session has been taken over by the attacker. This attack occurs in real-time and the attacker uses the session for whatever purpose he/she wants and can maintain the session for an extended period of time [17].*

In order to successfully execute a Session Hijacking attack, the attacker performs two tasks. Firstly, the attacker masquerades as the valid target to the WLAN. This requires a successful eavesdropping on the target communication to gather the necessary information. Secondly, the attacker deluges the air with a sequence of spoofed disassociate packets to keep the legitimate target out of the session [17].

*2) Replay Attack: This type of attack is not a real time attack and uses the legitimate authentication sessions to access the WLAN. The attacker first captures the authentication of a session or sessions. Later on, the attacker replays authenticated sessions to gain access to the network without altering or interfering with the original session or sessions [17].*

*3) 802.11 Frame Injection Attack: In a frame injection attack intruders capture or send forged 802.11 frames. They also inject their own Ethernet frames into the middle of the transmission. For example, an attacker could inject a frame while a user is trying to logon into a banking website. The website looks legitimate but it is not, as the attacker has injected Ethernet frames. Thus, all the login information will be recorded by the intruders [16].*

*4) 802.11 Data /802.11X EAP / 802.11 RADIUS replay attack: This attack involves the capture of 802.11/ 802.11X EAP/ 802.11 RADIUS data frame or authentication information and save it for later use. This information can be used for 80.1X EAP or for 802.1 X RADIUS authentication. Once the attacker captures and saves the authentication information, they can monitor traffic for another authentication in order to inject saved frames instead of the legitimate authentication frames to gain access to the system [22].*

*5) 802.11 Data deletion: This type of attack involves the attacker deleting the data being transmitted. An attacker could jam the wireless signal from reaching its intended target and provide acknowledgements (ACKs) back to the sources. As a result, data would never reach the legitimate target and the senders have no idea as they appear to receive ACKs [22].*

### D. Availablity Attacks

This attack prevents or prohibits the legitimate clients by denying access to the requested information available on the network. DoS attack is the most common type of availability attack which focuses on attacking a specific part of the network so the network becomes unreachable. There are several types of DoS attack which are described below:

*1) Denial-of-Service Attack: In this type of attack, an attacker tries to prevent or prohibit the normal use of the network communication by flooding a legitimate client with*

*bogus packets, invalid messages, duplicate IP or MAC address.*

*2) Radio frequency (RF) Jamming: An 802.11 network operates in the unlicensed 2.4 GHz and 5 GHz frequency band. In this type of attack, the attacker jams the WLAN frequency with a strong radio signal which renders access points useless [17]. As a result, legitimate users cannot access the WLAN.*

*3) 802.11 Beacon Flood: An intruder overloads the network by flooding it with thousands of illegitimate beacons so that the wireless AP is busy serving all the flooding packets and cannot serve any legitimate packets. Thus, making it very difficult for legitimate clients to find the real AP [16].*

*4) 802.11 Associate/Authentication Flood: In this type DoS attack, an attacker sends thousands of authentication/association packets from MAC addresses in order to fill up the target AP's association table. This makes it harder for a legitimate user to gain access in the network [16].*

*5) 802.11 De-authentication & Disassociation: The attacker pretends to be a client or AP and sends unauthorized management frames by flooding thousands of de-authentication messages or disassociation messages to the legitimate target. This forces them to exit the authentication state or to exit the association state [21].*

*6) Queensland DoS / Virtual carrier-sense attack: In this type of attack, an intruder exploits the clear channel assessment (CCA) by periodically claiming a large duration field in a forged transmission frame to make a channel appeared busy. This prevents other clients from gaining access to the channel [16].*

*7) Fake SSID: The attacker floods the air with thousands of beacon frames with fake SSIDs and all the access points become busy processing the fake SSIDs [21].*

*8) EAPOL flood: In this type of attack, the attacker deluges the air with EAPOL beacon frames with 802.11x authentication requests to make the 802.1x RADIUS server busy. Thus, legitimate client authentication requests are denied [21].*

*9) AP theft - This an attack where the attacker physically removes the access point from the public space making the network unavailable for the user [16].*

*E. Authentication Attack*

In an authentication attack, an intruder steals legitimate user's identities and credentials in order to gain access to the public or private WLAN and services. Dictionary attacks and brute force attacks are the most common techniques in this category. Once they have got the required information, the attacker impersonates or masquerades as an authorized user. Thereby gaining all the authorized privileges in the WLAN [5].

*1) Dictionary & Brute force attack: A brute force attack involves trying all possible key's in order to decrypt the message. On the other hand dictionary attacks only try the possibilities which are most likely to succeed, usually derived from a dictionary file. If the appropriate time is given, a brute force attack can crack any key. Whereas, Dictionary attacks will be unsuccessful if the password is not in the dictionary [23].*

Most access points use a single key or password that is shared with all connecting devices on the wireless LANs. A brute force attack can be applied on sniffing packets captured by the attacker in order to obtain the key.

Authentication attacks that are directly or indirectly involved with brute force and dictionary attack techniques after capturing the required information are discussed below [16]:

*1) Shared Key Guessing: The attacker attempts 802.11 shared key authentication with the cracked WEP keys or with the provided vendor default key.*

*2) PSK Cracking: In this type of attack, the cracker first captures the WPA-PSK key handshake frame, using open source tools such as Aircrack-ng, Kismet etc. Later, they run a dictionary or a brute force attack to recover the WPA-PSK key.*

*3) Application Login Theft: The cracker captures user credentials e.g. e-mail address and passwords etc. from clear text application protocols.*

*4) VPN Login Cracking: The attacker runs brute force attacks on the VPN authentication protocol in order to gain the user credentials e.g. PPTP (point to point tunnelling protocol) password or IPsec Preshared Secret Key etc.*

*5) Domain Login Cracking: The cracker runs a brute force or dictionary attack on NetBIOS password hashes. Thus accessing the user credentials e.g. windows login and password.*

*6) 802.1X Identity Theft: The attacker captures 802.1X identity response packets. Later they run the brute-force attack to recover user identities.*

*7) 802.1X LEAP Cracking: The intruder captures 802.1X lightweight EAP beacon frames and then runs a dictionary attack in order to recover user credentials.*

*8) 802.1X Password: The attacker repeatedly attempts 802.1X authentication to guess the user's password by using a captured user's identity [16].*

Beyond the above attack categories there are many more attacks pertaining to 801.11 technologies and describing all those is out of the scope of this paper. For example, a WLAN is vulnerable to upper layer threats. Fishing messages, mass mailing worms and Trojan downloaders can be carried over either wired or wireless networks. Attackers can poison ARP and DNS caches on wireless devices. Furthermore, there are other kinds of attack that try to exploit the wireless encryption standard. Examples are the Chopchop attack, the Original Beck-Tews attack, Halvorsen-Haugen attacks, the hole 196 attack and the Ohigashi-Morii attack etc. [1].

## VI. SECURING A WIRELESS LAN

The above vulnerabilities and threats come to the conclusion that it is very important to make sure that the wireless network is secure whether for a home user or an

enterprise network. However, still there is no true security solution that has been implemented and is presently available. But, the following steps could serve as a guideline to prevent most known vulnerabilities and some common threats:

The security of a WLAN should be considered throughout the WLAN development lifecycle, from the initial design and deployment stage through implementation, maintenance and monitoring. The Administrator should ensure that the organization's WLAN client devices and AP's have followed standard security configurations and are always compliant with the organization's security policies. Furthermore, the organization should implement continuous attack and vulnerability monitoring and perform periodic technical security assessment to measure overall security of the WLAN [6].

The use of strong encryption standards protect WLANs from the worst threats. The best practice would be to enable Wi-Fi protected access WPA/WPA2 rather than WEP. Furthermore, it is recommended to uses the WPA2, AES-CCMP protocol rather than to use WPA, as WPA-TKIP uses the WEP encryption algorithm for backward compatibility [10]. However, when using WPA2-PSK, it is important to ensure that the users are using strong, longer and hard to guess passwords for authentication. Moreover, larger organisations should consider using certificate-based authentication mechanism or RADIUS, allowing the users to access their own managed credentials in order to protect their network from sharing [24].

All manufacturers' default SSIDs, usernames and passwords are very well known to hackers. Therefore, changing the default SSID is a crucial step for securing home and enterprise network. More importantly, in the case of choosing the name, a user should try and use a unique name that doesn't give much information away about the owner, such as house number, street name or business name. This could enable the hacker to identify the exact location of the network [8].

By disabling SSID, this effectively hides the access point. This means that the user has to manually configure the network name and password in order to access the WLAN. This provides a very light defence, as by using readily available sniffing software tools anyone can discover the hidden network name. However, further security for routers can be managed where WEP is the only option available [10].

For connections to an open network such as a Wi-Fi hotspot and those commonly provided by hotels, Starbucks, McDonalds and so on, a virtual private network (VPN) can be a good security solution to deliver consistent protection over any internet connection and provide end-to-end security on wireless devices. Furthermore, large organisations can benefit by using a VPN to secure data that is sent over to a home or business partner WLAN without having to rely on a business partner to secure their part. Employees can use a VPN-enabled device which uses a secure tunnelling protocol such as IPSec or SSL to connect to company networks. Besides, a VPN can be useful to secure traffic that is sent by devices such as Smartphones which frequently roam between wireless and wired network [5].

Captive portal is a kind of authentication method used for guest access to a network. It is widely used in public internet networks such as hotels, conference centres, cafes and so on. Using this method, users automatically get redirected to the login page. Once the user's credentials are verified, the user would then successfully be able to access the network. This challenge response authentication is encrypted using SSL to prevent a hacker from sniffing user's credentials. However, some portals offer only authentication without any encryption of password or user data. It is very important to make sure that the portal offers an adequate security service [25].

Virtual Local Area Networks (VLAN) are another technology that can be used in corporate wireless network to enforce a security policy. VLANs work by tagging LAN frames assigned to different workgroups. Those tags actually decide where incoming frames can and cannot go within the corporate network. For example, if a business provides guest and consultant access, all traffic coming from that wireless LAN will be tagged so that traffic is limited to the public internet thus, keeping them away from corporate data and services [5].

Network Access Control (NAC) is another authentication technology that can be used in conjunction with the 802.1x and VLANs to enforce an extra layer of security. Instead of filtering traffic based on IP addresses and port numbers, NAC controls user access to network resources based on the sender's authenticated user identity, the state of the user's device and the configured policy. With NAC, network devices like Ethernet switches, APs, routers and firewalls all can still control access but they are enforcing decisions made by the NAC. For example, NAC decisions can be enforced by permitting or denying the use of a particular SSID or using 802.1x to direct wireless clients to particular subnets or VLANs [5].

A wireless intrusion detection and prevention system can be an essential tool for identifying intrusions and notifying the system administrator of attacks. There is no option to stop passive sniffing on the network with the traditional firewall. As a result, WIDS/WIPS can be deployed to act as a watchdog in order to detect and prevent new threats and any malicious activity. A VPN used with WIDS/WIPS can provide a good security measure by actively monitoring the network to identify anomalies. This adds another layer of assurance for data confidentiality [5].

## VI.  CONCLUSION

Securing the wireless network is an ongoing process. Realistically, still there is no single true security measure in place. When a new technology is first introduced, hackers study the protocol, look for vulnerabilities and then cobble together some program and scripts to try to exploit those vulnerabilities. Overtime those tools become more focused, more automated and readily available and published on the open source network. Hence, they can be easily downloaded and run by anyone. So, we never eliminate all threats and vulnerabilities and even if we do, we will probably end up wasting money by defeating some low probability and low impact attack. On the other hand, if we start eliminating the biggest security loopholes, attackers may turn to easier targets.

Thus, true WLAN security is always going to be a game of balancing acceptable risk and the countermeasure to mitigate those risks. Understanding business risk, taking action to deter most important and most frequent attacks and following industry good practices gives us better security solutions.

### REFERENCES

[1] F. Sheldon, J. Weber, S. Yoo, W. Pan, "The Insecurity of Wireless Networks." IEEE Computer Society, vol. 10, no. 4, July/August, 2012, pp. 54-61.

[2] L. Wang, B. Srinivasan, N. Bhattacharjee, "Security Analysis and Improvements on WLANs", Journal of Networks, vol. 6, no. 3, March 2011, pp. 470-481

[3] W. Ashford, "More than a quarter of London's Wi-Fi networks are poorly secured", 2012 http://www.computerweekly.com/news/2240162747/More-than-a-quarter-of-Londons-Wi-Fi-networks-are-poorly-secured, [Accessed on: 14/11/12].

[4] P. Feng, "Wireless LAN Security Issues and Solutions", IEEE Symposium on Robotics and Applications, Kuala Lumpur, Malaysia, 3-5 June, 2012, pp. 921-924.

[5] L. Phifer, "Wireless Lunchtime Learning Security School", 2009, http://searchsecurity.techtarget.com/guides/Wireless-Security-School, [Accessed on: 14/11/12].

[6] M. Souppaya, K. Scarfone, "U.S Department of Commerce - Guidelines for Securing Wireless Local Area Networks (WLANs)", Gaithersburg, MD 20899-8930: National Institute of Standards and Technology, 2012, SP 800-153.

[7] S. Reddy, K. Rijutha, K. Ramani, S. Ali, C. Reddy, "Wireless Hacking – A WiFi Hack By Cracking WEP", IEEE 2nd International Conference on Education Technology and Computer, Shanghai, China, 22-24 June, 2010, p. 189-193.

[8] J. Li, M. Garuba, "Encryption as an Effective Tool in Reducing Wireless LAN Vulnerabilities", Fifth International Conference on Information Technology: New Generations, Las Vegas, Nevada, 7-9 April, 2008, pp. 557-562.

[9] D. Shiyang,"Compare of New Security Strategy With Several Others in WLAN", IEEE 2nd International Conference on Computer Engineering and Technology, Chengdu, China, 16-18 April, 2010. pp. 24-28.

[10] T. Stimpson, L. Liu, J., Zhang, R. Hill, W. Liu, Y. Zhan, "Assessment of Security and Vulnerability of Home Wireless Networks", IEEE 9th International Conference on Fuzzy Systems and Knowledge Discovery, Chongqing, China, 29-31 May, 2012, pp. 2133-2137.

[11] H. Bulbul, I. Batmaz, M. Ozel, "Wireless Network Security: Comparison of WEP (Wired Equivalent Privacy) Mechanism, WPA (Wi-Fi Protected Access) and RSN (Robust Security Network) Security

[12] P. Kahai, S. Kahai, "Deployment Issues And Security Concerns With Wireless Local Area Networks: The Deployment Experience At A University" Journal of Applied Business Research, 2004, vol. 20, no. 4, pp. 11-24.

Protocols.", Proceedings of the 1st international conference on Forensic applications and techniques in telecommunications, information, and multimedia, Adelaide, Australia, Jan 21-23, 2008. Belgium: Institute for computer Sciences, Social-informatics and telecommunications Engineering (ICST).

[13] M. Mathews, R. Hunt, "Evolution of wireless LAN security architecture to IEEE 802.11i (WPA2)", University of Canterbury, New Zealand

[14] SANS Institute Infosec Reading Room, "Wireless LAN: Security Issues and Solution" 2003, US: SANS Institute, 1.4b.

[15] B. Forouzan, Data Communications & Networking. 4th edition. New York: McGraw-Hill, 2008

[16] Search Security, (2011) Information security tutorials [Online], Available at: http://searchsecurity.techtarget.com/tutorial/Information-security-tutorials, [Accessed on: 14/11/12]

[17] D. Welch, S. Lathrop, "Wireless Security Threat Taxonomy", Proceeding of the Information Assurance Workshop, 2003. IEEE Systems, Man and Cybernetics Society, U.S Military Academy, West Point, NY, 18-20 June, 2003, pp. 76-83

[18] N. Sunday, "Wireless Local Area Network (WLAN): Security Risk Assessment and Countermeasures", Thesis (MSc), Blekinge Institute of Technology, 2008.

[19] C. Hurley, F. Thornton, M. Puchol, R. Roger, "WarDriving Drive, Detect, Defend A Guide to Wireless Security", 2004, US: Syngress

[20] Y. Zahur, T. Yang, "Wireless LAN Security and Laboratory Designs" Journal of Computing Sciences in Colleges, vol. 19, no. 3, January 2004, pp. 44-60.

[21] Z. Tao, A. Ruighaver, "Wireless Intrusion Detection: Not as easy as traditional network intrusion detection", Tencon 2005 IEEE Region 10 Conference, Melbourne, Australia, 21-24 Nov, 2005, pp. 1-5

[22] M. Roche, "Wireless Hacking Tools", 2007, Available at: http://www.cse.wustl.edu/~jain/cse571-07/ftp/wireless_hacking/index.html, [Accessed on: 14/11/12]

[23] YouTube, "Dictionary vs. Bruteforce Attacks – Explained", 2008, Available at: http://www.youtube.com/watch?v=2hveQ8QZ9MQ, [Accessed on: 14/11/12].

[24] J. Lyne, "Hot Tipes for Securing Your Wi-Fi Network", 2012, http://www.sophos.com/en-us/medialibrary/Gated%20Assets/white%20papers/sophostipsforsecuringwifinetwork.pdf, [Accessed on: 14/11/12].

[25] K. Hole, E. Dyrnes, P. Thorsheim, P., "Securing Wi-Fi Networks", IEEE Computer Society, 2005, vol. 38, no. 7, pp. 28-34.

# Novel Fractional Wavelet Transform with Closed-Form Expression

K. O. O. Anoh, R. A. A. Abd-Alhameed, and S. M. R. Jones
Mobile and Satellite Communication Research Centre, University of Bradford, United Kingdom

O. Ochonogor
Dept. of Electrical and Electronic Engineering, University of Westminster

Y. A. S. Dama
An-Najah National University, Nablus, Palestine

*Abstract*—A new wavelet transform (WT) is introduced based on the fractional properties of the traditional Fourier transform. The new wavelet follows from the fractional Fourier order which uniquely identifies the representation of an input function in a fractional domain. It exploits the combined advantages of WT and fractional Fourier transform (FrFT). The transform permits the identification of a transformed function based on the fractional rotation in time-frequency plane. The fractional rotation is then used to identify individual fractional daughter wavelets. This study is, for convenience, limited to one-dimension. Approach for discussing two or more dimensions is shown.

*Keyword*—*Fractional Fourier transform; wavelet; fractional wavelet transform*

## I. INTRODUCTION

An introduction has been given to a new family of wavelets that are formed from the fractional Fourier order of the Fourier transform [1-5]. The fractional order of the Fourier transform is discussed based on discrete Fourier transform (DFT) as the fractional Fourier transform (FrFT) [6, 7] which is believed to be related to the chirp-Fourier transform [7, 8]. The chirp signal is highly concentrated in the fractional domain and a time delay leads to a fractional shift making the FrFT an efficient tool for separating the chirp signal [2]. In fact, the property of the FrFT tool enables that such delays choreographing into/as noise can be effectively filtered off [8]. Meanwhile, there are other possible characterizations of the fractional domain based tools that can be derived from the fractional Fourier tool. This was introduced as fractional wave packet transform (FRWPT) in [3]. The FRWPT was aimed at combining the advantages of the WT and FrFT, but this transform is computationally expensive [2]. In the recent times, this relationship has received wider attention in the discussion of wavelet families based on the fractional order of the DFT such as in [1, 9] and in [2]. It is called the fractional wavelet transform (FrWT). It is hoped that the proposed FrWT circumvents the computational cost available in [3]. Although the new wavelet transform discussed in [3] stemmed from the parent impulse filter property of the wavelet function, we approach the problem in a new fashion. For instance, in this work we extend this novelty into discussing wavelet transform using the quadratic phase function, as an example, earlier mentioned in [1] and based on FrFT by exploiting the dilation and translation properties of the mother wavelets demonstrated in [10].

At the moment, other studies have followed different methods to showing the exact closed-form expression for wavelet transform, for instance, by using raised cosine function [11]. The exact closed form expression for discrete wavelet transform based on FrFT is derived in this study. The complexity of the transform within a novel family of wavelet proposed here is also described. The computational gain exhibited in this new design is well spelt out and stressed. This would however revive the interest in deploying wavelet in signal processing, for instance.

We have organized the remaining parts of this paper as: In Section II, we familiarize the reader with the basic wavelet theory, then the proposed wavelet in Section III and the closed-form expression for the proposed wavelet is described in Section IV. The conclusion is presented in Section V.

## II. TRADITIONAL WAVELET THEORY

Wavelets are orthonormal functions derived from the parent scaling functions. For instance, consider an input signal $f(t)$, that modulates the transforming function, or scaling function, $\varphi(t)$. There are narrowband functions $\psi(t)$ derivable from $\varphi(t)$, which are orthogonal wavelets useful in the design of multicarrier systems. By the Fourier relation and Parseval's theory, the signal for band-limited case can be periodic with $\beta$, $-2\pi \leq \beta \leq 2\pi$, so that if $\psi_{l,m}(t)$ belongs to a set of orthonormal functions, then;

$$\int \psi_{l,m}(t)\psi_{k,n}(t)dt = \delta(l-k)\delta(m-n) \qquad (1)$$

where $\delta(\bullet)$ is a Dirac delta. Equation (1) defines a simple orthogonality condition between two daughter wavelets. Since $\psi(t)$ is obtained from the decomposition of $\varphi(t)$, we can express the relationship of the input signal with $\varphi(t)$ in discrete form as [12];

$$S_{DWT} = \sum_{m=0}^{M-1} f(m)\varphi_m(t) \qquad (2)$$

where $M$ is the length of the characteristic filter. $f(m)$ is the discrete equivalent of $f(t)$. The mother wavelet has a clear relation to the filters;

$$\psi(t) = \sqrt{2}\sum_m g(m)\varphi(2t-m) \qquad (3)$$

where $g(m)$ is a high-pass filter (HPF) and can be directly derived or constructed from a low-pass filter (LPF) that comes from the parent scaling function as;

$$\varphi(t) = \sqrt{2} \sum_m h(m)\varphi(2t-m) \qquad (4)$$

where $h(m)$ is the LPF and $\varphi(t)$ is the scaling function. Thus, as an alternative to modulating the input symbols by the sinusoids, these half-band filters can be used. The high-pass filters construct the detail coefficient part of the signal while the low-pass filters construct the approximate coefficient part of the signal. The high-pass filter can be formed from the low-pass filter as:

$$g(n) = (-1)^n h(M+1-n) \qquad (5)$$

where $M$ is the length of the filter and $n$ is the prevailing filter coefficient index. Both the high-pass and low-pass filters constitute the filter bank [13] required in multiresolution analysis (MRA). In signal processing for example, the multiplexing/processing function can be equivalently used orthogonal basis function such as [14]:

$$\varphi_{m,n}(t) = \begin{cases} 1 & n = m \\ 0 & elsewhere \end{cases}$$

where $m$ and $n$ are scales and shifts respectively. $\varphi_{m,n}(t)$ represents the complex orthogonal DWT basis function similar to the traditional multicarrier system.

### III. PROPOSED FRACTIONAL WAVELET TRANSFORM

It is a common knowledge that when a mother wavelet is translated and dilated, the daughter wavelets are born. The shifting and translation parameters of the father wavelet (or scaling function) can be well represented and approximated respectively, each of which gives rise to a uniquely different family of wavelets (see [10]).

Earlier, [15] identified discrete wavelets for any family by approximating the shift and translation parameters to discrete coefficients. Alternatively, the DFT roots can as well be exploited to define a new family of wavelets, namely, fractional wavelet transform. The fractional orders of the DFT define a uniquely different FrFTs which must also identify daughter wavelets [1] consequent on the DFT roots [16] that characterize the FrFT order.

Now, let us recall the definition of wavelet as represented in [1, 3]; for instance, let the mother wavelet be defined as,

$$\psi(t) = e^{i\pi t^2} \qquad (6)$$

But daughter wavelets are translated and shifted parts of Equation 6. So, let the daughter wavelets be defined as:

$$\psi_{(k,\tau)} = \frac{1}{\sqrt{k}} \psi\left(\frac{t-\tau}{k}\right) \qquad (7)$$

where $\tau$ and $k$ are shifts and scale parameters respectively. These parameters are defined as $k = 2^d$ and $\tau = 2^d n$ for traditional discrete wavelets [15].

Equation 7 (in discrete sense) becomes $\psi_{d,n} = 2^{-d/2}\psi(2^{-d}t - n)$ where $d$ and $n$ are equivalent shift and scale parameters respectively. Let Equation 7 be defined in terms of Equation 6 as:

$$\psi_{(k,\tau)}(t,k) = \frac{1}{\sqrt{k}} \exp\left[i\pi\left(\frac{t-\tau}{k}\right)^2\right] \qquad (8)$$

If we consider a signal $f(t)$ to be transformed by the wavelet transform, then the resulting transformation can be expressed as:

$$x(\tau,k) = \frac{1}{\sqrt{k}} \int_{-\infty}^{+\infty} \exp\left[i\pi\left(\frac{t-\tau}{k}\right)^2\right] f(\tau)\,d\tau \qquad (9)$$

Notice that for unit impulse [17], in other words, for a maximum of unit amplitude input signals:

$$\int_{-\infty}^{+\infty} f(t)\delta(t-\tau)d\tau = f(\tau)$$

Equation 9 is necessarily a Fourier transform of $f(t)$ that is shifted and translated by $\tau$ and $k$ respectively. For a rotation angle $\alpha$ which defines the fractional Fourier order of *a-fractional* rotation, $0 < |\alpha| < \pi$ (i.e. $0 < |a| < 2$), then the fractional Fourier transform when $\alpha$ is not a multiple of $\pi$-rad can be expressed as [1, 3, 16]:

$$B_a(t,\tau) = \frac{\exp\left[-i\{\pi\hat{\alpha}/4 - \alpha/2\}\right]}{\sqrt{\sin\alpha}}$$
$$\times \exp\left[i\pi(t^2\cot\alpha - 2t\tau\sec\alpha + \tau^2\cot\alpha)\right] \qquad (10)$$

where $\hat{\alpha} = \text{sgn}(\sin\alpha)$ and $\alpha = \dfrac{a\pi}{2}$. Thus, there exist a unique fractional Fourier transform for every order $a$. From Equation 10, let $u = t\sec\alpha$ such that,

$$x(u,\tau) = f(t) = f\left(\frac{u}{\sec\alpha}\right) \qquad (11)$$
$$= A_\alpha D_\alpha \int_{-\infty}^{+\infty} \exp\left[i\pi\left(\frac{u-\tau}{\tan^{1/2}\alpha}\right)^2\right] f(\tau)d\tau$$

where,

$$A_\alpha = \frac{\exp\left[-i(\pi\hat{\alpha}/4 - \alpha/2)\right]}{\sqrt{\sin\alpha}} \text{ and } D_\alpha = \exp(-i\pi u^2\sin^2\alpha)$$

Figure 1 shows an asymptotic representation of the behaviour of DFT roots which implies that the possible Fourier roots that govern the fractional Fourier order can never be zero. Although at $a = 1$, the traditional DFT is obtained. Thus, the possible roots that define the kernel of the FrFT or the resulting fractional wavelet cannot be obtained from a zero root.
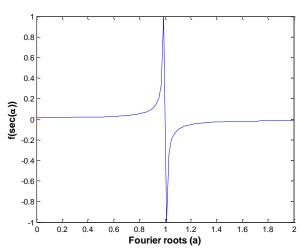
Fig. 1.    A graph of the possible fractional Fourier order

From Equation 11 $\tan^{1/2}\alpha$ is the scale parameter. Thus, Equation 11 is a fractional Fourier transform in the resemblance of a wavelet transform. If the wavelet is of the quadratic phase function where $\psi(t) = \exp(i\pi t^2)$, as of Equation 1 where the coordinate is scaled by $\tan^{1/2}\alpha$ and the amplitude is scaled by $A_\alpha$, then the discussion can proceed for signal of interest to exploit the FrFT property and wavelet property also. So, the convolution in the integral of Equation 11 is a wavelet transform. Notice that Equation 11 characterizes one–dimensional (1-D) wavelet transform only. However, for two-dimensional (2-D) and so on wavelet transforms, the following definitions must be followed [9]:

$$\psi(t,z) = e^{i\pi(t^2+z^2)} \tag{12}$$

Equation 12 is a typical 2-D wavelet transform. The scaled and dilated equivalent of Equation 12 can be expressed as:

$$\psi_{(k,\tau,\xi)}(t,z) = \exp\left[i\pi\left(\left(\frac{t-\tau}{k}\right)^2 + \left(\frac{z-\xi}{k}\right)^2\right)\right] \tag{13}$$

where $\tau$ is the shift parameter respective to *t*-coordinate, $\xi$ is the shift parameter respective to the *z*-coordinate. Recall the quadratic phase function based wavelet transform in Equation 9 and the fractional Fourier transform based wavelet constructed in Equation 11. We can define the phase function wavelet transformation based on Equation 11 by substituting the scaling factor *k* in Equation 6 into Equation 9 so that:

$$x(\tau,k) = \frac{1}{\tan^{1/2}\alpha}\int_{-\infty}^{+\infty}\exp\left[i\pi\left(\frac{t-\tau}{\tan^{1/2}\alpha}\right)^2\right]f(\tau)d\tau \tag{14}$$

where $k = \tan^{1/2}\alpha$. Thus in 1-D respect, if the function $f(\tau)$ is translated/dilated by $\tau$ then the *t* coordinate will be scaled by $\tan^{1/2}\alpha$. This gives a unique wavelet transform for every possible fractional Fourier order. This suggests that, instead of scaling $f(\tau)$ by some discrete factors (such as $k = 2^d$), the signal $f(\tau)$ can be scaled by the rotation factor of the

fractional Fourier order ($k = \tan^{1/2}\alpha$). It can be further stressed that for every frequency content, the shift and fractional rotation (or frequency fractional) content is well localized. The properties of the new wavelet function can be studied based on the fractional Fourier transform wavelet (fractional-wave) discussed in [3] and extended in [2, 5]. It exploits the time-frequency MRA advantage of the WT and the fractional frequency domain explanation of an input signal advantage of the FrFT. Thus, the new transform can provide information on a signal of interest in a fractional scale during the MRA in the time-frequency plane.

### IV.   CLOSED-FORM EXPRESSION OF THE PROPOSED WAVELET TRANSFORM

We can proceed to finding the exact discrete closed-form expression of Equation 14 while assuming that $k = \tan^{1/2}\alpha$ for brevity. Without loss of generality, recall the WT defined in Equation 9 can be expanded to accommodate the FrWT starting from:

$$x(\tau,k) = \frac{1}{k}\int_{-\infty}^{+\infty}\exp\left[i\pi\left(\frac{t}{k}-\frac{\tau}{k}\right)^2\right]f(\tau)d\tau$$

$$= \frac{1}{k}\int_{-\infty}^{+\infty}\exp\left[i\pi\left\{\left(\frac{t}{k}\right)^2 + \left(\frac{\tau}{k}\right)^2 - \frac{2\tau t}{k^2}\right\}\right]f(\tau)d\tau$$

$$= \frac{1}{k}\int_{-\infty}^{+\infty}\exp\left[i\pi\left(\frac{t}{k}\right)^2\right].\exp\left[i\pi\left\{\left(\frac{\tau}{k}\right)^2 - \frac{2\tau t}{k^2}\right\}\right]f(\tau)d\tau$$

$$= \frac{1}{k}\exp\left[\frac{i\pi t^2}{k^2}\right]\times\int_{-\infty}^{+\infty}\exp\left[i\pi\left\{\frac{\tau^2}{k^2} - \frac{2\tau t}{k^2}\right\}\right]f(\tau)d\tau$$

So,

$$x(\tau,k) = \frac{1}{k}\exp\left[\frac{i\pi t^2}{k^2}\right]\times\int_{-\infty}^{+\infty}\exp\left[\frac{i\pi}{k^2}\left\{\tau^2 - 2\tau t\right\}\right]f(\tau)d\tau$$

Let $y = \frac{i\pi}{k^2}\left(\tau^2 - 2\tau t\right)$ such that,

$$x(\tau,\tan^{1/2}\alpha) = \frac{e^{i\pi t^2/k^2}}{k}\int_{-\infty}^{+\infty}\exp[y]f(\tau)d\tau$$

$$= \frac{e^{i\pi t^2/k^2}}{k}\int_{-\infty}^{+\infty}\exp[y]d\tau\,f(\tau) \tag{15}$$

But from traditional derivative theory knowledge and taking the first derivative of *y*,

$$\frac{dy}{d\tau} = \frac{i\pi}{k^2}(2\tau - 2t)$$

$$= \frac{2i\pi}{k^2}(\tau - t) \tag{16a}$$

From Equation 16a, it can be rewritten that:

$$\Rightarrow d\tau = \frac{dy}{\left(\frac{2i\pi}{k^2}(\tau - t)\right)} \tag{16b}$$

Now, substituting for $d\tau$ : Put Equation 16b into Equation 15,

$$x(\tau,k) = \frac{e^{i\pi t^2/k^2}}{k} \times \int_{-\infty}^{+\infty} \exp[y] \frac{dy}{\left(\frac{2i\pi}{k^2}(\tau - t)\right)} f(\tau)$$

$$= G \times \int_{-\infty}^{+\infty} \exp[y] \frac{dy}{\left(\frac{2i\pi}{k^2}(\tau - t)\right)} f(\tau)$$

where $G = \dfrac{e^{i\pi t^2/k^2}}{k}$ , and remember that $k = \tan^{1/2}\alpha$ :

$$x(\tau,k) = G \times \int_{-\infty}^{+\infty} \exp[y] \times \frac{1}{\left(\left(\frac{2i\pi}{k^2}\right)(\tau - t)\right)} \times dy\, f(\tau)$$

$$= \frac{1}{\left(\left(\frac{2i\pi}{k^2}\right)(\tau - t)\right)} \times G \times \int_{-\infty}^{+\infty} \exp[y]\, dy\, f(\tau)$$

$$= \frac{1}{\left(\left(\frac{2i\pi}{k^2}\right)(\tau - t)\right)} \times G \times \left[e^y\right]_{\infty}^{+\infty} f(\tau)$$

The lower limit of the above relation is 0 (This is because $e^{-\infty} = 0$). However, recall that $y = \frac{i\pi}{k^2}\left(\tau^2 - 2\tau t\right)$, so

$$x(\tau,k) = \frac{1}{\left(\frac{2i\pi}{k^2}(\tau - t)\right)} \times G \times \exp\left[\frac{i\pi}{k^2}\left\{\infty^2 - 2t.\infty\right\}\right] f(\tau)$$

For discrete values of above expression, say $0 \le n < N$. Then,

$$x(n,k) = \sum_{n=0}^{N-1} \frac{1}{\left(\frac{2i\pi}{k^2}(n - t)\right)} \times G \times \exp\left[\frac{i\pi}{k^2}\left(n^2 - 2tn\right)\right] f(n) \tag{17}$$

By factoring the $n$ terms accordingly,

$$x(n,k) = \sum_{n=0}^{N-1} \frac{1}{\left(\frac{2i\pi(n-t)}{k^2}\right)} \times G \times \exp\left[\frac{ni\pi}{k^2}(n - 2t)\right] f(n) \tag{18}$$

Substituting for $G$ in Equation 18 with $k = \tan^{1/2}\alpha$,

$$x(n,k) = \sum_{n=0}^{N-1} \frac{(k)^2}{2i\pi(n-t)} \times \frac{e^{i\pi t^2/(k^2)}}{k} \times \exp\left[\frac{ni\pi}{k^2}\{n - 2t\}\right] f(n)$$

Eliminating the $k$ terms appropriately, Equation 18 can be written in a compact form such that:

$$x(n,k) = \sum_{n=0}^{N-1} \frac{(k)}{2i\pi(n-t)} \cdot e^{i\pi t^2/(k)^2} \times \exp\left[\frac{ni\pi}{k^2}\{n - 2t\}\right] f(n)$$

On the other hand,

$$x(n,k) = k \sum_{n=0}^{N-1} \frac{e^{i\pi t^2/(k)^2}}{2i\pi(n-t)} \times \exp\left[\frac{ni\pi}{k^2}\{n - 2t\}\right] f(n)$$

Factorizing the terms accordingly,

$$x(n,k) = \frac{k\, e^{i\pi t^2/\tan\alpha}}{2i\pi} \sum_{n=0}^{N-1} \frac{1}{n - t} \times \exp\left[\frac{ni\pi}{\tan\alpha}\{n - 2t\}\right] f(n) \tag{19}$$

From Equation 19, the term $(k\, e^{i\pi t^2/\tan\alpha})/2i\pi$ can be seen as a normalization parameter in the discrete sense. Similarly to the formulation of Equation, from Equation 11, the discrete relation can be expressed as:

$$x(u,\tau) = \sum_{n=0}^{N-1} A_\alpha D_\alpha \frac{G_{new}}{\frac{2i\pi(n-u)}{k^2}} \times \exp\left[\frac{ni\pi}{k^2}\{n - 2u\}\right] f(n) \tag{20}$$

Or,

$$x(u,\tau) = A_\alpha D_\alpha \sum_{n=0}^{N-1} \frac{k^2 G_{new}}{2i\pi(n-u)} \times \exp\left[\frac{ni\pi}{k^2}\{n - 2u\}\right] f(n)$$

where,

$$G_{new} = \frac{e^{i\pi t^2/k^2}}{k} = \frac{e^{i\pi u^2/(k)^2}}{k}$$

So that,

$$x(u,\tau) = A_\alpha D_\alpha \sum_{n=0}^{N-1} \frac{k^2}{2i\pi(n-u)} \times \frac{e^{i\pi u^2/(k)^2}}{k} \times \exp\left[\frac{ni\pi}{k^2}\{n - 2u\}\right] f(n)$$

By eliminating the $k$ terms appropriately, we obtain that:

$$x(u,\tau) = A_\alpha D_\alpha \sum_{n=0}^{N-1} \frac{k \, e^{i\pi u^2/(k)^2}}{2i\pi(n-u)} \times \exp\left[\frac{ni\pi}{k^2}\{n-2u\}\right] f(n) \qquad (21)$$

Now, recall that: $\hat{\alpha} = \text{sgn}(\sin\alpha) = \dfrac{\sin\alpha}{|\sin\alpha|}$, then expanding $A_\alpha$ and $D_\alpha$:

$$(AD)_\alpha = \frac{\exp[-i(\pi\hat{\alpha}/4 - \alpha/2)]}{\sqrt{\sin\alpha}} \times \exp(-i\pi u^2 \sin^2\alpha)$$

$$= \frac{\exp\left[-\left(i\pi\sin\alpha\left\{\dfrac{1}{4\times|\sin\alpha|} - u^2\sin\alpha\right\} - \dfrac{\alpha}{2}\right)\right]}{\sqrt{\sin\alpha}} \qquad (22)$$

Now, factorizing Equation 21:

$$x(u,\tau) = A_\alpha D_\alpha \frac{k \, e^{i\pi u^2/(k)^2}}{2i\pi} \sum_{n=0}^{N-1} \frac{1}{(n-u)} \times \exp\left[\frac{ni\pi}{k^2}\{n-2u\}\right] f(n) \qquad (23)$$

Substituting $(AD)_\alpha$ of Equation 22 for $A_\alpha D_\alpha$ in Equation 23:

$$x(u,\tau) = \frac{\exp\left[-\left(i\pi\sin\alpha\left\{\dfrac{1}{4\times|\sin\alpha|} - u^2\sin\alpha\right\} - \dfrac{\alpha}{2}\right)\right]}{\sqrt{\sin\alpha}} \times \frac{k \, e^{i\pi u^2/(k)^2}}{2i\pi}$$

$$\times \sum_{n=0}^{N-1} \frac{1}{(n-u)} \times \exp\left[\frac{ni\pi}{k^2}\{n-2u\}\right] f(n)$$

This can be written in a more compact form as:

$$x(u,\tau) = k \times \frac{\exp\left[-\left(\left\{\dfrac{i\pi\sin\alpha}{4\times|\sin\alpha|} - iu^2\sin^2\alpha\right\} - i\dfrac{\alpha}{2}\right) + \dfrac{i\pi u^2}{k^2}\right]}{2i\pi\sqrt{\sin\alpha}}$$

$$\times \sum_{n=0}^{N-1} \frac{1}{(n-u)} \times \exp\left[\frac{ni\pi}{k^2}\{n-2u\}\right] f(n)$$

But,

$$\frac{k}{\sqrt{\sin\alpha}} = \frac{\tan^{1/2}\alpha}{\sqrt{\sin\alpha}} = \frac{(\tan\alpha)^{1/2}}{\sqrt{\sin\alpha}} = \frac{\sqrt{\tan\alpha}}{\sqrt{\sin\alpha}} = \sqrt{\frac{1}{\cos\alpha}} = \frac{\sqrt{1}}{\sqrt{\cos\alpha}}$$

Clearly, Equation 23 can be exposed as:

$$x(u,\tau) = \frac{\exp\left[-\left(\left\{\dfrac{i\pi\sin\alpha}{4\times|\sin\alpha|} - iu^2\sin^2\alpha\right\} - i\dfrac{\alpha}{2}\right) + \dfrac{i\pi u^2}{k^2}\right]}{2i\pi\sqrt{\cos\alpha}}$$

$$\times \sum_{n=0}^{N-1} \frac{1}{(n-u)} \times \exp\left[\frac{ni\pi}{k^2}\{n-2u\}\right] f(n)$$

If $\quad \gamma = \dfrac{\exp\left[-\left(\left\{\dfrac{i\pi\sin\alpha}{4\times|\sin\alpha|} - iu^2\sin^2\alpha\right\} - i\dfrac{\alpha}{2}\right) + \dfrac{i\pi u^2}{k^2}\right]}{2i\pi\sqrt{\cos\alpha}}$ is

taken to be the normalization parameter for the alternative earlier fractional wavelet transform comparable to the proposed FrWT in Equation 19, then Equation 23 becomes:

$$x(u,\tau) = \gamma \times \sum_{n=0}^{N-1} \frac{1}{(n-u)} \times \exp\left[\frac{ni\pi}{\tan\alpha}\{n-2u\}\right] f(n) \qquad (24)$$

Now, comparing Equations 24 and 19, it can be observed that the complexity overhead associated with Equation 24 is more than that of Equation 19 (proposed). The computational cost of implementing FrWT based on Equation 11 is discouraging which is greatly reduced in the case of Equation 14. Meanwhile, that there are different wavelet for each fractional Fourier order, the idea of compact support stressed in [18] can be well accommodated. Also in the combination of the MRA property of the wavelet transform and the fractional Fourier property of the FrFT, the proposed well addresses a new wavelet that can achieve MRA in a fractional domain sense.

## V. CONCLUSION

A new kernel for discussing the wavelet transform has been presented. It was derived from the fractional Fourier properties of the fractional Fourier transform to exploit the wavelet transform properties. The new wavelet combines the MRA and the fractional Fourier properties to discuss input signal in the fractional domain sense. Analytical results obtained were explicitly described in discrete and closed form solution unlike any work before. Also, it was identified that the fractional wavelet transform presented shows uniquely different wavelet for every particular fractional Fourier order. The new fractional wavelet obtained was shown to be computational efficient that the earlier fractional wavelet using explicit discrete definition shown.

## VI. ACKNOWLEDGMENT

REFERENCES

[1] H. M. Ozaktas, B. Barshan, D. Mendlovic, and L. Onural, "Convolution, filtering, and multiplexing in fractional Fourier domains and their relation to chirp and wavelet transforms," *JOSA A,* vol. 11, pp. 547-559, 1994.

[2] J. Shi, N. Zhang, and X. Liu, "A novel fractional wavelet transform and its applications," *Science China Information Sciences,* vol. 55, pp. 1270-1279, 2012.

[3] Y. Huang and B. Suter, "The fractional wave packet transform," in *Recent Developments in Time-Frequency Analysis*, ed: Springer, 1998, pp. 67-70.

[4] D. Mendlovic, Z. Zalevsky, D. Mas, J. García, and C. Ferreira, "Fractional wavelet transform," *Applied optics,* vol. 36, pp. 4801-4806, 1997.

[5] G. Bhatnagar, Q. M. J. Wu, and B. Raman, "Discrete fractional wavelet transform and its application to multiple encription," *Information Scieneces,* vol. 223, pp. 297 - 316, 2013.

[6] N. Cotfas and D. Dragoman, "New definition of the discrete fractional Fourier transform," *arXiv preprint arXiv:1301.0704,* 2013.

[7] L. B. Almeida, "The fractional Fourier transform and time-frequency representations," *IEEE Transactions on Signal Processing,* vol. 42, pp. 3084-3091, 1994.

[8] X.-G. Xia, "Discrete chirp-Fourier transform and its application to chirp rate estimation," *IEEE Transactions on Signal Processing,* vol. 48, pp. 3122-3133, 2000.

[9] L. Onural, "Diffraction from a wavelet point of view," *Optics letters,* vol. 18, pp. 846-848, 1993.

[10] L. Debnath, "WAVELET TRANSFORM AND THEIR APPLICATIONS," *PINSA - A,* vol. 64, A, No. 6, pp. 685 - 713, 1998.

[11] G. Walter and J. Zhang, "Orthonormal wavelets with simple closed-form expressions," *IEEE Transactions on Signal Processing,* vol. 46, pp. 2248-2251, 1998.

[12] K. O. O. Anoh, R. A. Abd-alhameed, J. M. Noras, and S. M. R. Jones, "Wavelet Packet Transform Modulation for Multiple Input Multiple Output Applications," *IJCA,* vol. 63 - Number 7, pp. 46 - 51, 2013.

[13] B. Negash and H. Nikookar, "Wavelet-based multicarrier transmission over multipath wireless channels," *Electronics Letters,* vol. 36, pp. 1787-1788, 2000.

[14] K. O. Anoh, R. A. Abd-Alhameed, M. Chukwu, M. Buhari, and S. M. Jones, "Towards a Seamless Future Generation Network for High Speed Wireless Communications," *International Journal of Advanced Computer Science & Applications,* vol. 4, 2013.

[15] I. Daubechies, *Ten lectures on wavelets* vol. 61: SIAM, 1992.

[16] E. Sejdić, I. Djurović, and L. J. Stanković, "Fractional Fourier transform as a signal processing tool: An overview of recent developments," *Signal processing,* vol. 91, pp. 1351-1369, 2011.

[17] K. A. Stroud and D. J. Booth, *Advanced engineering mathematics*: Palgrave macmillan, 2003.

[18] I. Daubechies, "Orthonormal bases of compactly supported wavelets," *Communications on pure and applied mathematics,* vol. 41, pp. 909-996, 1988.

# Cost Effective System Modeling of Active Micro-Module Solar Tracker

Md. Faisal Shuvo
Dept. of ECE,
Khulna University of Eng. & Tech.,
Khulna-9203, Bangladesh

Md. Abu Saleh Ovi
Dept. of ECE,
Khulna University of Eng. & Tech.,
Khulna-9203, Bangladesh

Md. Mehedi Hasan
Dept. of ECE,
Khulna University of Eng. & Tech.,
Khulna-9203, Bangladesh

Arifur Rahman
Dept. of ECE,
Khulna University of Eng. & Tech.,
Khulna-9203, Bangladesh

Mirza Md. Shahriar Maswood
Dept. of ECE,
Khulna University of Eng. & Tech.,
Khulna-9203, Bangladesh

*Abstract*— The increasing interests in using renewable energies are coming from solar thermal energy and solar photovoltaic systems to the micro production of electricity. Usually we already have considered the solar tracking topology in large scale applications like power plants and satellite but most of small scale applications don't have any solar tracker system, mainly because of its high cost and complex circuit design. From that aspect, this paper confab microcontroller based one dimensional active micro-module solar tracking system, in which inexpensive LDR is used to generate reference voltage to operate microcontroller for functioning the tracking system. This system provides a fast response of tracking system to the parameters like change of light intensity as well as temperature variations. This micro-module model of tracking system can be used for small scale applications like portable electronic devices and running vehicles.

*Keyword— micro-module, active solar tracker, LDR sensor, microcontroller.*

## I.    INTRODUCTION

Renewable energies can technically contribute to practically all sectors of energy demand, that is, fuel for transportation, electricity, and low temperature heat for space heating and hot water and, to a limited degree, to high temperature process heat. As other sources of energy are finite and someday will be depleted but renewable energy will not run out ever. Most renewable energy investments are spent on materials and workmanship to build and maintain the facilities, rather than on costly energy imports. Solar energy is one of the best ways to reduce the problem of energy deficiency and environment pollution in modern society. In the actuality a lot of research works have been conducted to improve the use of the solar energy. Tracking technologies can also be applied to rural and remote areas for small scale applications, where energy is often crucial in human development because photovoltaic power generation is clean and pollution less with respect to other available resources [1] having wide range of applications [2]. Different kinds of single axis or two axis solar `s and programs were presented in the previous year's [3], [4]. Some optimization of solar trackers concentrated on shading, cost, and sizing for photovoltaic systems and mass-flow rate for thermal systems were discussed [5-8]. While many solar energy tracking projects are large-scaled [9] having complex design with maximum power point calculation [10], [11] as well as solar panels are body mounted and fuzzy based [12], the tracking system is controlled by the photovoltaic sensors and its accuracy is determined by the accuracy of mechanical model [13], [14] and combination of integrated circuit as well as large array sensor system [15]. Single axis or one dimensional solar tracker is a device having one degree of freedom that orients the solar panel to the direction of incident rays of sunlight for particularly in one dimensional rotation for obtaining maximum power for longer duration. It will be effective more than dual axis tracker if we consider some other criteria like circuital power consumption, sensor accuracy, mechanical simplicity, wind loading and tolerance to misalignment. There are two types of tracker according to their functional properties like active tracker and passive tracker. Active solar tracker uses external sources of energies to power blowers, pumps and other types of equipment to collect, store and convert solar energy. Once energy is received it stored for later appliances. Small systems are used to furnish electricity for heating and cooling system in homes and other sections, where large system can furnish power for entire communities. As the installation cost and complexity of dual axis tracker is very much larger than single axis or one dimensional tracker corresponds to its efficiency as well as the mechanical module. On the other hand single axis tracker offers lower cost and higher reliability. The possible horizontal rotation can be possible manually if needed if we consider one dimensional tracker vertically for small scale applications. This active one dimensional tracker will provide the best power output relevant to less installation cost than any other tracking system.

This paper proposes an active micro-module solar tracking system focusing on cost effective circuital setup with less complexity based on LDR sensor with its maximum accuracy as well as microcontroller as function of ADC converter and comparator to attain maximum energy storage with less power consumption and having fast tracking response for small scale

applications as well as portable devices witsih maximum performance corresponding to different parameters. Besides these we also consider about the power consumption of motor as well as motor driver for maximum output response as well.

The rest of the paper organized as follows. Section II briefly overviews of our proposed micro-module model. Section III describes the proposed algorithm of our method. Software simulation as well as experiments and result analysis are described in Section IV and V. The overall discussion is represented in Section VI. We conclude the paper and future research directions in Section VII.
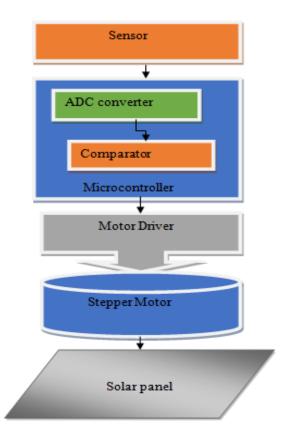
## II. OUR PROPOSED MODEL



Fig.1. Our proposed micro module tracking model

Fig.1 shows our improved solar tracker model in which we have tried to implement a complexion free tracking setup of vertical single axis or one dimensional module tracker using LDR (Light Dependent Resistors), microcontroller, a motor driver chip and a stepper motor which will control the tracking system of solar panel from east to west with necessary programs included in microcontroller as well.

### A. Sensing Protocol

Photo sensors are used in many projects involving sensing of light and shadow. In case of making an efficient project we must consider of choosing reasonable sensor from wide variety among different types of sensors.
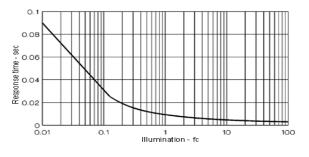


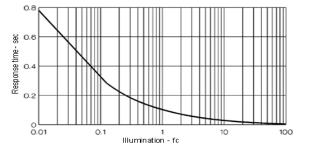Fig.2. Response time vs. illumination-decay time



Fig.3. Response time vs. illumination-rise time

We have used LDR as photo sensor that is cheap and rugged in nature which is suitable for external uses. They are basically resistors whose resistance depends on the intensity of light. By considering the properties of LDR which is suitable for applications where many different levels of light intensity are to be measured with its moderate response time. In here we have used the properties of changing resistance and voltage drop across the LDR. Fig.2 & Fig.3 shows speed of response at which a photocell responds to a change from light-to-dark or from dark-to-light. The rise time is defined as the time necessary for the light conductance of the photocell to reach $1-1/e$ or about 63% of its final value. The decay or fall time is defined as the time necessary for the light conductance of the photocell to decay to $1/e$ or about 73% of its illuminated state.

### B. Processing of Sensor Output with Microcontroller

As light signal or illumination is continuous analog signal so it is practically difficult to compare the output of two sensors we have used in east and west direction with real time to track sun's position without comparing it. For this comparison we need to convert analog voltage to discrete voltage using ADC conversion which is built in function of microcontroller we have used Atmega8L which reduces the external implementation ADC circuit. In it receives analog signal directly from the LDR sensor's voltage drop and then convert it to discrete value which can be determined numerically then compare the voltage difference of these sensors. This is done by small program loaded in it. This reduces addition of circuital arrangement of reference voltages. The program loaded in microcontroller have a small algorithm to check, compare and give necessary commands as output for taking the decision of rotation of tracker. To drive the motor we have used a motor driver L293D that will help to gear up the stepper motor of 12v.

## C. Reason of Using Stepper Motor for Maximum Efficiency

The power consumption of the motor depends on the moment of inertia which relates to the output torque of the motor and the moment of inertia relates to the factor of width and length of the solar panel as well as the tilt angle of the panel related to the step-angle of the motor and the basic equation of moment of inertia is given by the following equation.

$$I = \frac{\frac{1}{12}m(l^2 + w^2)}{cos\theta} \qquad (1)$$

Where the mass $'m'$, length $'l'$ and width $'w'$ of the solar panel is fixed and the only way to reduce the inertia by reducing the tilt angle of the solar panel which is similar to the step-angle of the motor.

$$I \propto \frac{1}{cos\theta} \qquad (2)$$

$$cos\theta \propto \frac{1}{\theta} \qquad (3)$$

If the step-angle θ decreases, $cos\theta$ increases and similarly inertia decreases. So, the stepper motor is used to reduce the inertia because its step-angle is $1.8°$ and the reason to reduce the inertia is to reduce the output torque of the motor because the torque is directly proportional to the inertia and it determines the power consumption of the motor.

$$\tau = I\alpha \qquad (4)$$

Where, $'\tau'$ is the output torque of the motor and $'\alpha'$ is the angular acceleration of the motor. The reduced inertia reduces the output torque and the reduced torque reduces the power consumption of the motor to achieve the maximum efficiency.

### III. TRACKING ALGORITHM

We use a simply modified algorithm based on LDR sensing and microcontroller processing.Fig.4 represents this algorithm provides the automatic tracking system from east to west in a day and for the other day it will adjust it position again automatically.

According to Fig.4 firstly we determine light intensity incident from the sunlight of East and West sensors. From the difference of two sensors we will decide the stepper motor will rotate or not. If the difference exceeds 0.3 volt which is defined in program the motor starts to move on. Then according to given instruction it compare and check statements for the East and West light intensity and step on to desired position to hold the panel in order to attain maximum solar energy.
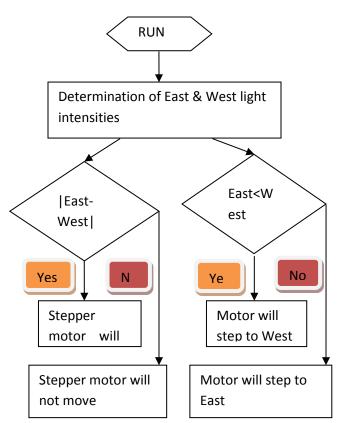


Fig.4. Flowchart of micro module tracking algorithm

### IV. SOFTWARE SIMULATION

Fig.5 shows our ISIS professional simulated circuit. We have developed a program for AVR on the basis of an algorithm. First of all we implemented two sensor circuits to take input from the outer environment like radiation intensity of sunlight. Two LDR are used to take input of light intensity. Those are connected in ADC0 and ADC1 pin. By using ADC pin of Atmega8L we took that input then compare both input data.

Then it will convert into binary data with respect its logical unit. Then we send binary input to output pin for corresponding logical program. The PD0-PD3 is set to output port. These output ports are connected to stepper motor through driver L293D.

This driver drive and rotate the stepper motor in both direction according to the program generated in the microcontroller. The program is simulated in Micro C AVR studio & converted into a hex file to load in microcontroller as well.
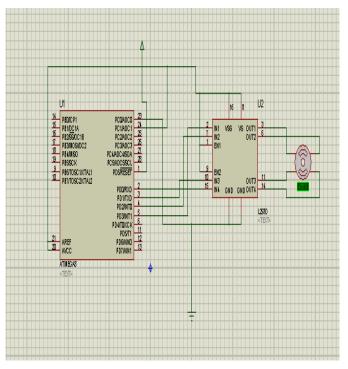
Fig.5. ISIS professional simulated circuit

## V. EXPERIMENTAL EVALUATION

TABLE I. PERFORMANCE ANALYSIS OF TRACKING RESPONSE OVER A DAY.

| Initial time | Final time | Resistance of LDR 1 (Ω) | Resistance of LDR 2 (Ω) | Volt. Drop of Sensor 1 (V) | Volt. Drop of Sensor 2 (V) |
|---|---|---|---|---|---|
| 8:00 Am | 8:10 am | 512 | 500 | 3.30 | 3.34 |
| 9:05 Am | 9:13 am | 178.5 | 166.5 | 4.24 | 4.29 |
| 10:30 am | 10:40 am | 112 | 100 | 4.49 | 4.55 |
| 11:45 am | 11:53 am | 83.4 | 71.5 | 4.61 | 4.66 |
| 12:53 pm | 12:58 pm | 64.6 | 52.6 | 4.69 | 4.75 |
| 1:30 Pm | 1:35 pm | 62 | 50 | 4.70 | 4.76 |
| 2:53 Pm | 2:58 pm | 74.5 | 62.6 | 4.65 | 4.71 |
| 3:55 Pm | 4:00 pm | 112 | 100.5 | 4.50 | 4.55 |
| 4:55 Pm | 5:00 pm | 262 | 250 | 3.96 | 4.00 |

As the light intensity increases, the resistance of two LDR decreases and for this reason the voltage drop across two LDR also decreases and the light resistance as well as voltage drop of LDR is determined from the following equation which is closest to our experimented value.

$$\text{Light resistance, } R_L = \frac{500}{lux} K\Omega \qquad (5)$$
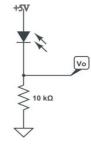
Voltage drop across LDR,

$$V_{LDR} = \frac{5 * R_L}{(R_L + 10)} volt \qquad (6)$$

As the voltage drop decreases with the increasing of light intensity, the current flow increases which causes the power consumption.

$$P_c = I^2 R \ watt \qquad (7)$$

For this reason, we use a 10K resistor in the sensor circuit to reduce this consumption and the sensor circuit is shown in the following figure.



Fig.6. Sensor circuit

10K is preferred to 100Ω or 1KΩ because of its lower power consumption but if we use 25KΩ power consumption can be reduced whether the sensitivity also decreases.

The output voltage of sensor circuit,

$$V_o = 5 - V_{LDR} \ volt \qquad (8)$$

As the voltage drop across LDR decreases, the output voltage drop of sensor circuit increases and this output voltage is used as the input of micro-controller for analog to digital conversion which reduces the power consumption. The data Table.1 shows the relationship between resistance of LDR-1 as well as LDR-2 and light intensity across day time from 8:00 am to 5:00 pm and Fig.7 shows the resistance characteristic curve of LDR-1 and LDR-2.
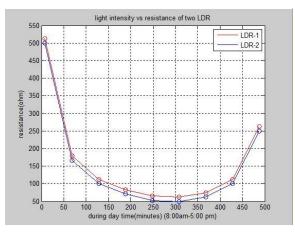


Fig.7. Matlab simulated resistance variation curve of LDR-1 &LDR-2

In Fig.8, we have also observed that the voltage drop of Sensor-2 (West) is always greater than Sensor-1 (East) during day time and for this reason the solar panel rotates from East to West during day time and at night it stays at the West direction and in the morning it moves to the East direction because at that time we place an another additional LDR-3 as switch at East which sense the raising sunlight and automatically hold the panel to the East direction by an additional program loaded in microcontroller.
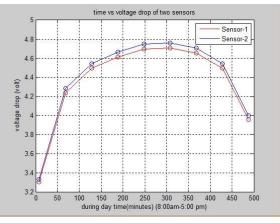


Fig.8.    Voltage drop of two sensors across time

## VI.    DISCUSSION

Now-a-days research on renewable energy for compensating our loss of resources of energy has been becoming more and more popular. Among different types of energy options the best and easiest option is solar energy is virtually inexhaustible than other natural resources. With economic consideration here we have implemented a new tracking scheme for covering wide range of uses of small applicable sections. We have observed different types of tracking system and their performance relevant to their installation cost and started our work on the basis of that with aiming at reduction of cost, complexion corresponding with maximum performance. We have modified our sensing part, replace large complicated control circuits in a single chip and evaluated it. We have also observed device response through over a long day and compared our experimental data as well as performance.

## VII.    CONCLUSION

This paper is proposed focusing on an improved version of tracking system in order to provide faster response and suitable size for obtaining maximum power for small scale applications. The circuit is implemented in such a way to reduce power consumption of tracker. Here we have included much additional circuit like comparator and voltage regulator in a single chip in order to minimize complexion as well as increase efficiency with possible lower cost. System performance has been evaluated through the simulation. Above all to provide constant performance, replacement of chip in case of damage with less installation cost for wide variety of appliances can be provided by this micro-module solar tracking system.

In future we are interested to work on designing of solar panel for small objects in order to minimize the power consumption of motor as well as driver for higher degree of accuracy.

REFERENCES

[1]    B. K. Bose, P. M. Szezesny, and R. L.Steigerwald, "Microcontroller control of residential photovoltaic power conditioning system," *IEEE Transaction on InustrialAppication.*, vol. 21, no. 5, pp. 1182–1191,1985.

[2]    ZouJian, Ji Xing, Du Haitao, "Research of a New Automatic Solar Tracking System [J]," *Photoelectron Technology*, pp.159-163, 2010.

[3]    J. M. Enrique, J. M. Andújar and M.A. Bohórquez, "A reliable, fast and low cost maximum power point tracker for photovoltaic applications", *Solar Energy*, vol. 84, pp. 79-89, 2010.

[4]    H. Arbab, B. Jazi, and M. Rezagholizadeh, "A computer tracking system of solar dish with two-axis degree freedom based on picture processing of bar shadow", *Renewable Energy*, vol. 34, pp. 1114-1118, 2009.

[5]    D. Weinstock and J. Appelbaum, "Optimization of Solar Photovoltaic Fields,"*ASME Journal of Solar Energy Engineering*, vol. 131, pp. 031003, Iun. 2009.

[6]    D. Weinstock and J. Appelbaum, "Optimization of Economic Solar Field design of Stationary Thermal Collectors", *ASME Journal of Solar Energy Engineering*, vol. 129, pp. 363-370, 2007.

[7]    S. Conti, G. Tina and C. Ragusa, "Optimal Sizing Procedure forStand-Alone Photovoltaic Systems by Fuzzy Logic," *ASME Journal of Solar Engineering,* vol. 124, pp. 77-82, 2002.

[8]    V. Badescu, "Optimal control of flow in solar collectors for maximum energy extraction," *International Journal of Heat and Mass Transfer*, vol. 50, pp. 4311-4322, 2007.

[9]    M. Davis, J. Lawler, J. Coyle, A. Reich, T. Williams Green Mountain Engineering, LLC,"MachineVision as  a Method for Characterizing Solar Tracker Performance," *In the proc. of 33rd IEEE Photovoltaic Specialists Conference*, pp. 1-6, May 2008.

[10]    S. Mekhilef, N.A. Rahim, and H.W. Ping, "Performance of maximum power point tracker in tropical climate," *In the proc. of 37th IEEE Power Electronics Specialists Conference*, pp. 1-4, 2006.

[11]    Il-Song Kim, Myung-Bok Kim, and Myung-JoongYoun, "New Maximum Power Point Tracker Using Sliding-Mode Observer for Estimation of Solar Array Current in the Grid-Connected Photovoltaic System," *IEEE Transactions on Industrial Electronics,* vol. 53, pp. 1027-1035, Jun. 2006.

[12]    Mohsen Taherbaneh,MohammadB.Menhaj, "A Fuzzy-Based Maximum Power Point Tracker for Body Mounted Solar Panels in LEO Satellites," *In the proc. of IEEE/IAS Industrial & Commercial Power Systems Technical Conference*, pp. 1-6, May 2007.

[13]    YizhuGuo,Jianzhong Cha, "A System Modeling Method for Optimization of a Single Axis Solar Tracker," *In the proc. of 2010 International Conference on Computer Application and System Modeling (ICCASM)*, vol. 11, pp. 30-34, Oct. 2010.

[14]    Alex Joseph,Kamala J, "Economic and Backlash Tolerable Solar Tracking System," *In the proc. of 2013 International Multi-Conference on Automation, Computing, Communication, Control and Compressed Sensing (iMac4s)*, pp. 748-753, Mar. 2013.

[15]    Daniel A. Pritchard, "Sun Tracking by Peak Power Positioning for Photovoltaic Concentrator Arrays," *IEEE Control Systems Magazine*, vol. 3, pp. 2-8, Aug. 1983.

AUTHORS PROFILE

**Md. Faisal Shuvo** is currently pursuing B.Sc degree program in electronics and Communication Engineering in Khulna University of Engineering and Technology, Khulna-9203, Bangladesh. His research interest includes antenna design, machine learning, pattern recognition, neural networks, image processing, and telecommunication and feature selection.

**Md. Abu Saleh Ovi** is currently pursuing B.Sc degree program in electronics and Communication Engineering in Khulna University of Engineering and Technology, Khulna-9203, Bangladesh. His research interest includes antenna design, machine learning, pattern recognition, neural networks, image processing, and telecommunication and feature selection.

**Md. Mehedi Hasan** received the B.Sc degree in Electronics and Communication Engineering from Khulna University of Engineering and Technology, Khulna-9203, Bangladesh, 2013. His research interest includes antenna design, biomedical image processing, machine learning, pattern recognition, neural networks, image processing, and telecommunication and feature selection.

**Arifur Rahaman** received the B.Sc degree in Electronics and Communication Engineering from Khulna University of Engineering and Technology, Khulna-9203, Bangladesh, 2013. His research interest includes antenna design, machine learning, pattern recognition, neural networks, image processing, and telecommunication and feature selection.

**Mirza Md. Shahriar Maswood** received his B. Sc degree in Electrical and Electronic Engineering from Khulna University of Engineering and Technology. He received M. Sc degree in ECE from Khulna University of Engineering .He is accomplishing my Ph.D. degree beginning at fall 2013 semester from UMKC (University of Missouri, Kansas City) which is located at the state of Missouri. . His research interest includes Network security: Enhancing security in wireless sensor network, ad-hoc network, Image Processing, Machine Learning etc.

# A Review of Scripting Techniques
# Used in Automated Software Testing

Milad Hanna

Department of Computer Science,
Faculty of Computers and Information,
Helwan University,
Cairo, Egypt

Nahla El-Haggar

Lecturer of Information Technology,
Faculty of Computers and Information,
Helwan University,
Cairo, Egypt

Mostafa Sami

Professor of Computer Science, HCI
lab, Faculty of Computers and
Information,
Helwan University,
Cairo, Egypt

*Abstract—* **Software testing is the process of evaluating the developed system to assess the quality of the final product. Unfortunately, software-testing process is expensive and consumes a lot of time through software development life cycle. As software systems grow, manual software testing becomes more and more difficult. Therefore, there was always a need to decrease the testing time. Recently, automation is as a major factor in reducing the testing effort by many researchers. Therefore, automating software-testing process is vital to its success. This study aims to compare the main features of different scripting techniques used in process of automating the execution phase in software testing process. In addition, an overview of different scripting techniques will be presented to show the state of art of this study.**

*Keyword— Software Testing; Automated Software Testing; Test Data; Test Case; Test Script; Manual Testing; Software Under Test; Graphical User Interface.*

## I. INTRODUCTION

Software testing has evolved since 1970's as an integral part of software development process. Through it, the final quality of the software can be improved by discovering errors and faults through interacting, checking behavior and evaluating the System Under Test (SUT) to check whether it operates as expected or not on a limited number of test cases with the aim of discovering errors that are found in the software and fixing them. According to Ilene Burnstein, software testing describes as a group of procedures carried out to evaluate some aspect of a piece of software [1]. Ehmer Khan [2] shortly defines it as a set of activities conducted with the intent of finding errors in software. In addition, according to Ammann and Offutt [3] software testing means evaluating software by observing its execution.

Since software-testing process is a very expensive process, complete testing is practically impossible and it is not acceptable to reduce testing effort by accepting quality reductions. Testing effort is often a major cost factor during software development. Many software organizations are spending up to 40% of their resources on testing [4]. Therefore, an existing open problem is how to reduce testing effort without affecting the quality level of the final software.

Automation is one major solution for reducing high testing effort. Automating certain manual tasks from software testing process can save a lot of testing time. It can help in performing repetitive tasks more quickly than manual testing.

## II. MANUAL TESTING VS. AUTOMATED TESTING

Software testing can be divided into two main categories, manual testing, and automated software testing. Both categories have their individual strengths and weaknesses.

With a manual testing, the more traditional approach, tester initiates each test, interacts with system, reports and evaluate the test results. To satisfy the test results manually, testers should prepare and execute test cases on SUT. These test cases will best test the system using defined processes trying to find bugs. So, they can be fixed before releasing the product to the public [5].

Automation is one of the more popular and available strategies to reduce testing effort. It develops test scripts that will be used later to execute test cases instead of human [6]. The idea behind automation is to let computer simulate what the tester is doing in reality when running test cases manually on SUT. AST is more suitable for repetitive tasks during different testing levels such as regression testing, where test cases are executed several times whenever the source code of SUT is modified or updated [7].

Katja Karhu [5] summarizes the difference between the two categories by suggesting that automated software testing should be used to prevent new errors in the already tested working modules, while manual testing is better used for finding new and unexpected errors. The two approaches are complementary to each other, automated testing can perform a large number of test cases in little time, whereas manual testing uses the knowledge of the tester to target testing to the parts of the system that are assumed to be more error-prone.

## III. SCRIPTING TECHNIQUES

Test scripts are the basic element of automation. Test script is a series of commands or events stored in a script language file to execute a test case and report the results. It may contain logical decisions that affect the execution of the script, creating multiple possible pathways, constant values, variables whose values change during playback. The advantage of test scripts development process is that scripts can repeat the same instruction many times in loops, each time with different data. There are many types of scripting techniques that can be used in automation. Fewster and Graham [8] listed five different types of scripting techniques that will be discussed in this section.

## A. Linear Scripting Technique

John Kent [9] explains the idea behind linear technique, which is simply to set the test tool to the record mode while performing actions on the SUT. The generated recorded script consists of a series of testing instructions using the programming language supported by the tool. Gerald Everett suggested that the linear scripts are being created by recording the actions that a user performs manually on interface of the system and then saving test actions as a test script. These test scripts can then be replayed back to execute the test again. So, linear scripting technique is called Record/Playback [10]. 2Figure 1 illustrates record/playback steps.
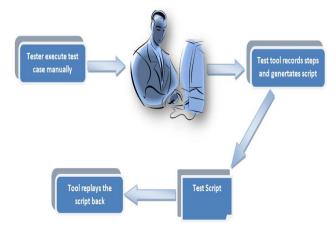


Figure 1: Record/Playback Steps

Microsoft® Visual Studio® Team Edition is an example for tool applying linear scripting technique. It enables testers to perform record and playback to be used to create and execute the tests [11].

Every time a test case is being automated using linear technique, new test script is generated. Thus, the more test cases are automated, the more lines of code are generated. This means that the number of Lines of Code (LOC) is proportional to the number of automated test cases [9]. Thus: Lines of code α Number of automated test cases

Figure 2 and Figure 3 shows a practical example for applying linear scripting technique on a simple login page, which followed by the recorded test script that presents the sequence of actions performed manually on that page.
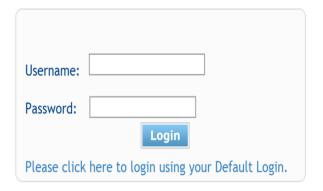


Figure 2: Login Page

```
                        //Fileds
/// <summary>
/// Go to web page 'http://TestSite/login.aspx' using new
browser instance
/// </summary>
public string UIWelcometoTestSiteWinWindowUrl
="http://testsite/login.aspx";

/// <summary>
/// Type 'test-user' in 'txtUserName' text box
/// </summary>
public string UITxtUserNameEditText = "test-user";

/// <summary>
/// Type '{Tab}' in 'txtUserName' text box
/// </summary>
public string UITxtUserNameEditSendKeys = "{Tab}";

/// <summary>
/// Type '********' in 'txtPassword' text box
/// </summary>
public string UITxtPasswordEditPassword =
"to+VpC5U2lKdiNhE9v4dzPA0ZmKuc60K";

/// <summary>
/// Type '{Enter}' in 'txtPassword' text box
/// </summary>
public string UITxtPasswordEditSendKeys = "{Enter}";

                        //Actions
// Go to web page the webpage using new browser instance
this.UIWelcometoTestSiteWinWindow.LaunchUrl(new
System.Uri(this.LoginParams.UIWelcometoTestSiteWinWindowU
rl));

// Type 'test-user' in 'txtUserName' text box
uITxtUserNameEdit.Text =
this.LoginParams.UITxtUserNameEditText;

// Type '{Tab}' in 'txtUserName' text box
Keyboard.SendKeys(uITxtUserNameEdit,
this.LoginParams.UITxtUserNameEditSendKeys,
ModifierKeys.None);

// Type '********' in 'txtPassword' text box
uITxtPasswordEdit.Password =
this.LoginParams.UITxtPasswordEditPassword;

// Type '{Enter}' in 'txtPassword' text box
Keyboard.SendKeys(uITxtPasswordEdit,
this.LoginParams.UITxtPasswordEditSendKeys,
ModifierKeys.None);
```

Figure 3: Linear Script for Login Page

John Kent [9] mentioned main advantages for linear scripting technique as listed below:

- It enables tester to start automating quickly as no planning is required, tester can just simply record any manual test case.

- The tester does not need to have any programming skills.

- It is good for demonstrating the SUT.

John Kent [9] mentioned the shortcomings for linear scripting technique as listed below:

- The generated scripts are very difficult to be maintained because they are made up of long lists of actions of objects interacting with interface, it contains its own hard-coded data, and this is not the best way for saving them.

- The recorded script can only work under exactly the same conditions as when it was recorded at the first time. If simple error happened or unexpected normal events (e.g. file not found) during a test run, it will not be handled correctly by the test script.

- Linear test scripts are not reliable enough, even if the application has not changed. They often fail on replay because other things occurred that did not happen when the test was recorded.

### B. Structured and Shared Scripting Techniques

Both structured and shared scripting techniques are being formed by using structured programming instructions that are used to control the flow of execution of the script.

Structured scripting technique uses structured programming instructions, which either be control structures or calling structures [8]. Control structures is used to control the different paths in the test script (e.g. If condition). Calling structures is used to divide large scripts into smaller and more manageable scripts. For example, one script can call another script to perform specific functionality and then return to the first script where the subscript was called. The most important advantage of structured technique is that the test script can validate for specific conditions to determine if the executed test passed or failed according to these conditions. However, the script has now become a more complex program and the test data still tightly coupled within the test script itself. Besides, implementing structured scripts require not only testing skills but also programming skills [8].

Figure 4 shows applying structured scripting technique on a simple login web page.

Shared scripting technique enables common actions to be stored in only one place. This implies that a scripting language that allows one script to be called by another one is required. The idea behind shared scripts is to generate separate script that performs one specific common task that other scripts may need to perform later.

Thus, different test scripts can call this common task whenever they needed and testers will not have to spend time for implementing common actions many times across all scripts [12]. It works well for small-scale systems to be tested using relatively few test scripts. Figure 5 illustrates using shared scripting technique [12].

```
[TestMethod]
public void Login_TestMethod()
{
    WatiN.Core.Settings.WaitForCompleteTimeOut = 120;
    IE ie = new IE("http://testsite/login.aspx", true);
    ie.TextField(Find.ById("txtUserName")).Value = "test-
user";
    ie.TextField(Find.ById("txtPassword")).Value =
"12345678";
    ie.Button(Find.ById("btnLogin")).Click();
    ie.WaitForComplete();
    // If "Welcome" message is displayed, then the test
is passed
    if (ie.Text.Contains("Welcome"))
    {
        Console.WriteLine("Testing Passed");
    }
    else
    {
        //If not, then the test is failed
        Console.WriteLine("Testing Failed");
    }
}
```
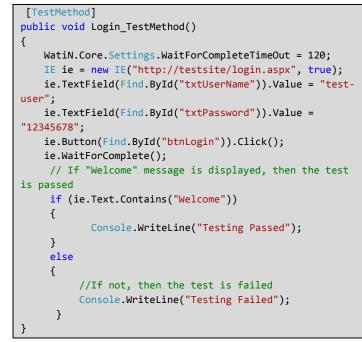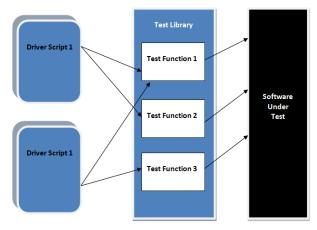
Figure 4: Structured Script for Login Page



Figure 5: Driver Scripts and a Test Library

For example, instead of having the same login action repeated in a number of scripts, tester could simply implement it once as shared script and each test script just have to call this common function as illustrated in Figure 6:

```
public IE Login()
{
    WatiN.Core.Settings.WaitForCompleteTimeOut = 120;
    IE ie = new
IE(Telco_Automation.Properties.Settings.Default.SiteURL,
true);
    ie.TextField(Find.ById("txtUserName")).Value =
"test-user";
    ie.TextField(Find.ById("txtPassword")).Value =
"amv1234!@#$";
    ie.Button(Find.ById("btnLogin")).Click();
    Assert.IsTrue(ie.Text.Contains("Welcome"));
    return ie;
}
```
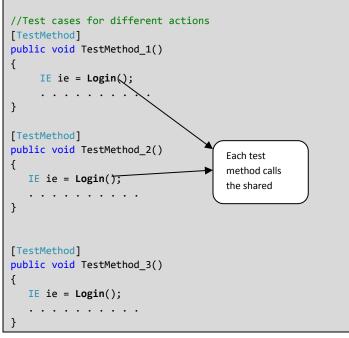
```
//Test cases for different actions
[TestMethod]
public void TestMethod_1()
{
    IE ie = Login();
    . . . . . . . . . .
}

[TestMethod]
public void TestMethod_2()
{
    IE ie = Login();
    . . . . . . . . . .
}

[TestMethod]
public void TestMethod_3()
{
    IE ie = Login();
    . . . . . . . . . .
}
```

Each test method calls the shared

Figure 6: Shared for Login Page

### C.  Data-Driven Scripting Technique

New additional scripting techniques are required to form test scripts in such a way that the maintenance costs of the test scripts can be reduced than in the previous scripting techniques. Data-Driven scripting technique proposes better organization of test scripts and hence lower maintenance costs of the test scripts. Bhaggan [13] demonstrates that test data is stored in a separate data file instead of being tightly coupled to the test script itself. While performing tests, test data is read from the external data file instead of being taken directly from the script itself. It allows both input data and expected results to be stored together separately from the script itself. For example, instead of having username and password data input values within the login script, these values can be stored in an external excel file and implement test script to read test data to use it while executing the test script.

In this technique, it is important that the external data file must be synchronized with the control script. This means that if any changes applied to the format of the data file, then the control script must be updated also to correspond to it.

To automate new test case, new control script has to be implemented with new data records inserted into external data file. Figure 7 illustrates data-driven scripting technique [12].
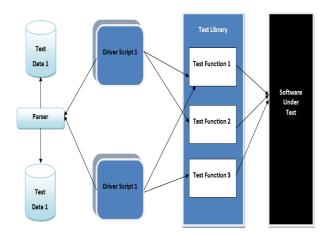


Figure 7: Data-Driven Scripting Technique

In data-driven scripting technique, the maintenance costs are lower than the costs of rerecording the tests from the beginning. Therefore, tests will not have to be rerecorded, but only maintained [13].

Linda G. Hayes [14] presents the main advantages of this approach as below:

- Similar tests can be added very quickly with different input data as the same script can be used to run different tests with different data. xIt may useful when testing large number of data values using the same control script.

- Data files are stored in easily and maintainable text records, so it can be updated.

- The format of data files can be modified to suit the testers with some modifications in the control script. For example, the data file can contain special column for comments that the control script will ignore while execution. This make the data file more readable, understandable and therefore maintainable.

The disadvantages of data-driven technique by Linda G. Hayes [14] are listed below:

- It requires high level of programming technical skills in the scripting language supported by the tool. Such tests need to be well managed, as it requires maintaining data files used by various test scripts. This may increase the cost for the project.

- One script is needed for every logically different test case. This can easily increase the amount of needed scripts dramatically. Laukkanen considered that this is the major problem in this technique [12].

The following test script with the external data file show applying data-driven scripting technique on a simple web page in Figure 8.



Figure 8: Sample Web Page

```
[TestMethod]
public void GeneratedTestMethod()
{
IE ie = new IE();
Application xlApp = new Application();
Workbook xlWorkbook =
xlApp.Workbooks.Open(@"D:\Data.xlsx");
Worksheet xlWorksheet = xlWorkbook.Sheets[1];
string url = ((Range)xlWorksheet.Cells[3,
2]).Text.ToString();
ie.GoTo(url);

ie.TextField(Find.ById(new Regex("txtTaskId"))).Value =
((Range)xlWorksheet.Cells[4, 2]).Text.ToString();

ie.TextField(Find.ById(new Regex("txtTaskORI"))).Value =
((Range)xlWorksheet.Cells[5, 2]).Text.ToString();

ie.SelectList(Find.ById(new
Regex("ddlTaskName"))).Options[int.Parse(((Range)xlWorksh
eet.Cells[6, 2]).Text.ToString())].Select();

ie.TextField(Find.ById(new
Regex("txtTaskComment"))).Value =
((Range)xlWorksheet.Cells[7, 2]).Text.ToString();

ie.Button(Find.ById(new Regex("btnClose"))).Click();
}
```

Figure 9: Generated Test Script for the Web Page



Figure 10: Output Data File Snapshot for the Web Page

### D. *Keyword-Driven Scripting Technique*

Keyword-Driven scripting technique is a very similar to manual test cases. The business functions of the SUT are stored in a tabular format as well as in step-by-step instructions for each test case. Keyword-driven approach separates not only test data for the same test as in data-driven scripts but also special keywords for performing business function in the external file. The tester can create a large number of test scripts simply using predefined keywords. All what the tester needs is just to know what keywords are currently available to be applied on SUT and what is the data that each keyword is expecting. Additional keywords can be added to the list of available programmed set of keywords to enlarge the scope of automation. It is more sophisticated than data-driven technique [12]. Fewster and Graham [8] state that the keyword-driven scripting technique is a logical extension of the data-driven scripting technique. A limitation of the data-driven technique is that the detailed steps of what the tests are doing are implemented within the control script itself. Therefore, keyword-driven technique takes out some of the intelligence from the script, put it into the external file with the test data, and leaves the task for reading both steps and data for the control script. Thus, instead of having data file in data-driven, complete test file is needed in keyword driven scripting technique. It doesn't contain test data only but also a complete description of the test case to be automated using a set of keywords to be read and interpreted later on while test case execution. The test file states what the test case will do, not how to do it.

Laukkanen [12] supposed that in order to execute the tabular automated test cases, there have to be a middle layer that converts the special keywords to the source code that interacts with SUT (the source code that implements the keywords are called "handlers"). The translation of keywords is implemented outside of the control script itself. Now, the control script only reads each keyword in order from the test file and calls corresponding supporting script. In addition, a driver script, which parses the test data and calls the appropriate keyword handlers, is needed. Figure 11 demonstrates these layers.
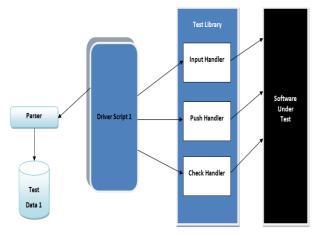


Figure 11: Handlers for Keywords

He also divides the keywords into two different levels of test keywords: high level and low-level keywords. Low-level keywords are more suitable for detailed testing on the interface level (e.g. Input, Click, and Select…etc.). High-level keywords are more suitable for testing higher-level functionality like SUT business logic (e.g. Create Account, Login…etc.). Multiple low level keywords can be combined together to form high-level keywords [12].

Zylberman and Shotten [7] show that keyword-driven technique is the next generation approach of automation that separates the task of automated test case implementation from the automation infrastructure. They state that keyword driven testing can be divided into two main layers:

1.  Infrastructure Layer: It is a combination of the three types of keywords. It receives the different keywords as inputs to perform operations on the SUT.

2.  Logical Layer: This layer helps manual testers to build new test scripts using the pre-defined keywords (that is already implemented in Infrastructure Layer).

They also divide the keywords into three different kinds (item/base level keywords, utility functions, and sequence/user keywords) which described below [7]:

1.  Item Operation: an action that performs a specific operation on a given GUI element. Parameters should be specified to perform an operation on a GUI item such as name of GUI item, operation to be performed and the values needed.

2.  Utility Functions: a script that executes a certain functional operation that is hard or ineffective to

implement as a sequence. For example: Run Application, Close Application, Wait X seconds, Retrieve Data from DB,...etc

3.  Sequence: a set of keywords that produces a business process such as "create customer" keyword. Sequence keyword is made by combining various items and functions.

They also suggest reducing the number of keywords by creating multi-function keywords. For example, "Update_Subscriber_Status" keyword is a better approach than creating two special keywords for "Activate_Subscriber" and "Deactivate_Subscriber" [7]. Although it can be argued that may be it is more useful not to combine keywords together because this allows using them again in creating another test script. For example, tester can use "Deactivate_Subscriber" keyword in another sequence of keywords (e.g. Delete_Subscriber).

Rantanen [15] suggests a new method for dividing system to multiple user stories. Each user story consists of one or multiple test cases. Each test case is to be mapped to the actual code interacting with the SUT while execution. Every test case contains one or more sentence format keywords. Every sentence format keyword consists of one or more user keywords which written in understandable text (they can be understood without technical skills). A user keyword consists of one or more base keywords. Finally, the base keywords contain the source code interacting with system to be tested. Figure 12 illustrates dividing SUT to multiple user stories.
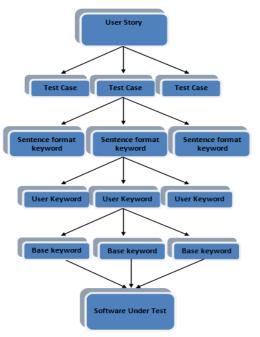


Figure 12: Mapping from User Story to SUT

Rashmi and Bajpai [16] proposed a new contribution for keyword-driven framework based on the concept of recording as shown in Figure 13. To start automating process, enter the URL of the system to be tested. Like linear scripting technique, the keyword driven testing framework records the steps while user navigates the web application manually. The

user name and password are entered in the appropriate text boxes, and then the user clicks the Log-In button. The tool records all the operations performed manually in the web browser until the test is stopped.

When finishing the recording, a corresponding test script file is generated that contains all user actions. The user actions consists of items clicked, items selected and value typed…etc. These steps are generated in tabular format, representing each operation performed in the form of keyword, value and operation.

When the test is finished and replayed back, the tool runs keywords that were saved in the output test script file. The SUT opens in a new web browser and all recorded steps are performed again automatically, as it was originally performed manually in the test.
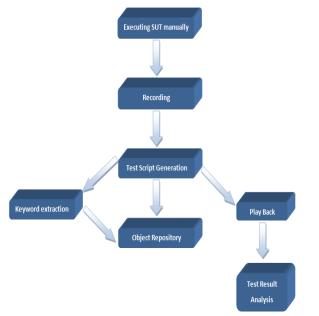


Figure 13: Keyword Driven Framework Based on Recording

After the new test run is completed, the test results are displayed to indicate the status of the test whether the test is passed or failed. The test results window displays two key elements of the test run for analysis purpose. The first one presents the steps that were performed while test execution while the second element presents the test result details.

Object Repository is a centralized place for storing the properties of available objects in SUT. Websites are developed using many different objects (e.g. textbox control, input tag). Each object is identified based on the object type. It has properties (e.g. name, title, caption, color, and size) and specific set of methods, which help in object identification.

Wissink and Amaro [6] state that the principal feature of the keyword-driven scripting technique is the separation of engineering tasks into a set of roles. These roles include test designer, automation engineer, and test executor. To automate

test cases in the keyword-driven scripting technique, the next steps are to be followed:

1. A set of actions need to be defined by the *test designer* and then documented in an external file with other keywords, input data, and expected results.
2. The *automation engineer* implements the different keywords defined above by the test designer in the programming language of the tool.
3. The *test executor* just runs the tests directly from the spreadsheet.

The advantages of keyword-driven scripting technique mentioned by Linda G. Hayes [14] are:

- Using keyword-driven scripting technique, the tester only needs to know keywords and learn how to use them.
- The number of generated scripts required for keyword-driven is dependent of the size of the SUT rather than the number of tests. This means that many more tests can be created without increasing the number of scripts.
- Like data-driven scripting technique, the way in which tests are created can be modified to suit the testers rather than the test tool, using the format and tools that the testers are most comfortable with.

The disadvantages of keyword-driven scripting technique mentioned by Linda G. Hayes [14] are:

- The costs for development of customized application specific functions (framework) are very high in terms of both time and human resources for technical skills. Such specific framework development can be considered as standalone software development that needs to be tested before using in testing other software.
- If the SUT requires more than just a few customized keywords, then testers should learn a high number of keywords.

## IV. DISCUSSION

According to the above review of the paper about the different scripting techniques demonstrated by Figure 14, which illustrates moving from linear to keyword scripting technique in addition to a comparison of the main features for each of them as in Table 1. We recommend applying the data-driven scripting technique for automating the execution phase through software testing process as it is considered as the most cost effective scripting technique.

It is necessary to spend time building the test to avoid high maintenance costs on the long run. If the tester spends more time to develop test scripts, maintenance cost will be lower. However, if tester uses the fastest way to create test scripts (record/playback), then the maintenance cost will be very high. The following Table 1 and Figure 14 present a comparison between different scripting techniques. Numbers used in the table range from 1 (Lowest) to 5 (Highest).

TABLE 1: COMPARISON BETWEEN DIFFERENT SCRIPTING TECHNIQUES

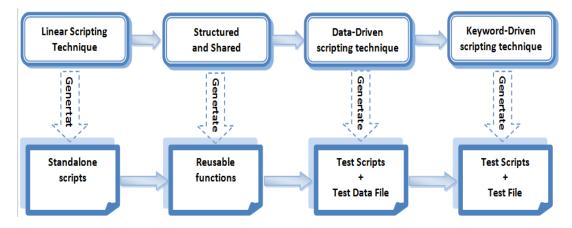| Property | Linear | Structured | Shared | Data-Driven | Keyword-Driven |
|---|---|---|---|---|---|
| **Ability to use reusable functions** | No | No | Yes | Yes | Yes |
| **Data separation from test script** | No | No | No | Yes | Yes |
| **Logic steps separation from test script** | No | No | No | No | Yes |
| **Access to code required** | No | Yes | Yes | Yes | Yes |
| **Use structured programming instructions** | No | Yes | Yes | Yes | Yes |
| **Ability to compare test results with expected** | No | Yes | Yes | Yes | Yes |
| **Ability to using script in regression testing** | No | Yes | Yes | Yes | Yes |
| **Special framework required** | No | No | No | No | Yes |
| **Programming skills level** | 1 (Low) | 2 | 3 | 4 | 5 (High) |
| **Effort needed to create test script** | 1 (Low) | 2 | 3 | 4 | 5 (High) |
| **Maintenance costs needed to update test script** | 5 (High) | 4 | 3 | 2 | 1 (Low) |
| **Reusability of test script** | 1 (Low) | 2 | 3 | 4 | 5 (High) |



Figure 14: Evolution of Test Automation

## V. CONCLUSION

Across many organizations, it is well known that testers lack the time needed to fully test the SUT within the time allocated to testing phase. This often happens because of unexpected environmental problems or problems in the implementation phase of development process. This normally shifts the software final delivery date. As a result to this delay, only two options is found, either to work longer hours or to add other resources to the test team to finalize testing in the required limited time. Automation can be one solution to this problem to accelerate testing and meet project deadline. Automation of testing phase offers a potential source of savings across all the life cycle. Automation using scripting techniques can save the costs for the overall software testing automation process, improve the speed of testing, shorten the product's launch cycle and it can achieve an amount of work that manual tests are impossible to finish.

REFERENCES

[1] I. Burnetein, "Practical Software Testing: process oriented approach," Springer Professional Computing, 2003.

[2] M. E. Khan, "Different Forms of Software Testing Techniques for Finding Errors," International Journal of Software Engineering (IJSE), vol. 7, no. 3, 2010.

[3] P. Ammann and J. Offutt, Introduction to Software Testing, New York: Cambridge University Press, 2008.

[4] F. Elberzhager, A. Rosbach, J. Münch and R. Eschbach, "Reducing test effort: A systematic mapping study on existing approaches," Information and Software Technology 54, p. 1092–1106, 2012.

[5] K. Karhu, T. Repo and K. Smolander, "Empirical Observations on Software Testing Automation," International Conference on Software Testing Verification and Validation, 2009.

[6] T. Wissink and C. Amaro, "Successful Test Automation for Software Maintenance," in 22nd IEEE International Conference on Software Maintenance (ICSM'06), 2006.

[7] A. Zylberman and A. Shotten, "Test Language: Introduction to Keyword Driven Testing," http://SoftwareTestingHelp.com, pp. 1-7, 2010.

[8]    M. Fewster, Software Test Automation: Effective Use of Test Execution Tools, Addison-Wesley Professional, 1999.

[9]    J. Kent, "Test Automation From RecordPlayback to Frameworks," http://www.simplytesting.com/, 2007.

[10]   M. Fewster, "Common Mistakes in Test Automation," Grove Consultants, 2001.

[11]   "How to: Generate a Coded UI Test by Recording the Application under Test," August 2013. [Online]. Available: http://msdn.microsoft.com/en-us/library/dd286608%28v=vs.100%29.aspx.

[12]   P. Laukkanen, "Data-Driven and Keyword-Driven Test Automation Frameworks," Helsinki University of Technology, Software Business and Engineering Institute, 2007.

[13]   K. Bhaggan, "Test Automation in Practice," Delft University of Technology, the Netherlands, 2009.

[14]   L. Hayes, The Automated Testing Handbook, Automated Testing Institute, 2004.

[15]   J. Rantanen, "Acceptance Test-Driven Development with Keyword-Driven Test Automation Framework in an Agile Software Project," Helsinki University of Technology, Software Business and Engineering Institute, 2007.

[16]   N. Bajpai, "A Keyword Driven Framework for Testing Web Applications," International Journal of Advanced Computer Science and Applications (IJACSA), vol. 3, 2012.