# IJACSA

Technically Co Sponsored By:

IJACSA Publications Indexed @

# IJACSA Editorial

## From the Desk of Managing Editor…

It is a pleasure to present to our readers the sixth issue of International Journal of Advanced Computer Science and Applications (IJACSA).

The unique purpose of this journal is to publish peer-reviewed research reports and theoretical articles on the use of innovative technologies in Advanced Computer Science and Applications.

With monthly feature peer-reviewed articles and technical contributions, the *Journal*'s content is dynamic, innovative, thought-provoking and directly beneficial to readers in their work.

IJACSA is published by The Science and Information Organization (SAI) that brings together relevant people, projects, and events in the domain.

The reader community has been drawn from various settings like industrial and commercial organizations, governments and administrations, or educational Institutions.

In order to publish high quality papers, Manuscripts are evaluated for scientific accuracy, logic, clarity, and general interest. Each Paper in the Journal will not merely summarize the target articles, but will evaluate them critically, place them in scientific context, and suggest further lines of inquiry. As a consequence only 33% of the received articles have been finally accepted for publication.

IJACSA emphasizes quality and relevance in its publications. In addition, IJACSA recognizes the importance of international influences on Computer Science education and seeks international input in all aspects of the journal, including content, authorship of papers, readership, paper reviewers, and Editorial Board membership

The success of authors and the journal is interdependent. I am extremely thankful to the editorial board members and peer reviewers for their constant encouragement that helped in the progress of the journal and earning credibility amongst all the reader members.

We hope that the relationships we have cultivated will continue and expand.

**Thank You for Sharing Wisdom!**

**Managing Editor**
**IJACSA**
**December 2010**
**editorijacsa@thesai.org**

# Editorial Board

# IJACSA Reviewer Board

- **Prof. Rakesh L**

  Professor, Department of Computer Science, Vijetha Institute of Technology, India

- **Dr.Sagarmay Deb**

  University Lecturer, Central Queensland University, Australia

- **Mr. Lijian Sun**

  Research associate, GIS research centre at the Chinese Academy of Surveying and Mapping, China

- **Dr. Abdul Wahid**

  University level Teaching and Research, Gautam Buddha University, India

- **Mr. Chakresh kumar**

  Assistant professor, Manav Rachna International University, India

- **Mr.Zhao Zhang**

  Doctoral Candidate in the Department of Electronic Engineering, City University of Hong Kong, Kowloon, Hong Kong

- **Dr. Parminder Singh Reel**

  Assistant Professor in Electronics and Communication Engineering at Thapar University, Patiala

- **Md. Akbar Hossain**

  Doctoral Candidate, Marie Curie Fellow, Aalborg University, Denmark and AIT, Greeceas

- **Prof. Dhananjay R.Kalbande**

  Assistant Professor at Department of Computer Engineering, Sardar Patel Institute of **Technology**, Andheri (West),Mumbai, India.

- **Hanumanthappa.J**

  Research Scholar Department Of Computer Science  University of Mangalore, Mangalore, India

- **Arash Habibi Lashakri**

  University Technology Malaysia (UTM), Malaysia

- **Vuda Sreenivasarao**

  Professor  & Head in the Computer Science and Engineering at St.Mary's college of Engineering & Technology, Hyderabad, India.

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

  Professor in the Dept. of Information Technology (IT), SNIST, India.

- **Suhas J Manangi**

  Program Manager, Microsoft India R&D Pvt Ltd

- **M.V.Raghavendra**

  Head, Dept of ECE at Swathi Institute of Technology & Sciences,India.

- **Dr. V. U. K. Sastry**

  Dean (R & D), SreeNidhi Institute of Science and Technology (SNIST), Hyderabad, India.

- **Pradip Jawandhiya**

  Assistant Professor & Head of Department

- **T V Narayana Rao**

  Professor and Head, Department of C.S.E –Hyderabad Institute of Technology and Management, India

- **Shubha Shamanna**

  Government First Grade College, India

- **Totok R. Biyanto**

  Engineering Physics Department - Industrial Technology Faculty, ITS Surabaya

- **Dr. Smita Rajpal**

  ITM University Gurgaon,India

- **Dr. Nitin Surajkishor**

  Professor & Head, Computer Engineering Department, NMIMS, India

- **Dr. Rajiv Dharaskar**

  Professor & Head, GH Raisoni College of Engineering, India

- **Dr. Juan Josè Martínez Castillo**

  Yacambu University, Venezuela

# CONTENTS

# Pause Time Optimal Setting for AODV Protocol on RPGM Mobility Model in MANETs

Sayid Mohamed Abdule, Suhaidi Hassan, Osman Ghazali, Mohammed M. Kadhum

InterNetWorks Research Group
College of Arts and Sciences
University Utara Malaysia
06010 UUM Sintok, Malaysia
sayidabdule@internetworks.my|{suhaidi|osman|kadhum}@uum.edu.my

*Abstract*— **For the last few years, a number of routing protocols have been proposed and implemented for wireless mobile Ad hoc network. The motivation behind this paper is to discover and study the pause time effects on Ad hoc on Demand Distance Vector (AODV) protocol and find out the node pause time optimal setting for this protocol where Reference Point Group Mobility (RPGM) model uses as a reference model. In order to come across the best performance of a particular routing protocol, there a need to examine a number of parameters with different performance and analyze the optimal setting of that protocol and its network configuration environment. This experiment, the speed is fixed with 20 ms in all scenarios while the pause time is varying from scenario to another to observe the optimal setting of the pause time on protocol's performance in this configuration. The outcome of the experiment are analyzed with different parameters such as varying number of nodes, increasing connections, increasing pause time and discussed the effects of the pause time. The results have shown that the value of the pause time can be affecting the performance of the protocol. In the experiment, we found that the lower pause time give better performance of the protocol. However, this paper is a part of ongoing research on AODV protocol in link failure. Thus, it is important to figure out the factors which can be involved the performance of the protocol.**

*Keyword-MANETs, AODV, pause time, optimal settin, RPGM.*

## I. INTRODUCTION

A mobile ad hoc network is a number of devices that communicates each other without any central administration and each of them acts as a router as it receives packet and forwards if it is not the destination. Actually Mobile Ad hoc Networks (MANETs) is a wireless communication and it uses a wireless interface to send and receives packet data, but in different way from infrastructure wireless. However, Mobility models are important building blocks in wireless networks and ad hoc developers can choose from a range of models that have been developed in the wireless both infrastructure less and infrastructure [1]. Ad hoc networks can be considered as a flexible application and uses for some special occasions that are inconvenient or unable to infrastructure-network. These applications which can be used in this communication are as follows: disaster aid, police operations in particular areas, group conferences, rescue operations, military deployment in

aggressive environment etc. In practically, there are many applications that can be applied by Ad hoc networks. Therefore, the Mobile Ad hoc Networks (MANETs) are of much interest to network developers both public and the privet sectors because of its potential applicable to establish an easy communication network in any situation that involves both emergencies and normal applications [2]. There are a number of mobility models which have been widely used for evaluating the performance of relevant protocols or algorithms in traditional Ad hoc networks. Among of these models are Random Waypoint (RW) model, Reference Point Group Mobility (RPGM) model, Freeway Mobility (FW) model, Manhattan Grid (MG) model, Gauss-Markov (GM) and many others [3]. In this paper, is used the Reference Point Group Mobility (RPGM) to generate the mobility scenarios to evaluate the performance of the protocol in this particular configuration. Most of MANETs simulations are used Random Waypoint (RW) as a reference mobility model [4], [5]. Thus, there is a need to provide additional mobility models in order to examine many different MANET applications.

## II. MOBILITY MODELS

MANETs, mobile nodes move from point to another point. The main role of mobility models is to emulate this movement of real mobile nodes. Mobility models are based on setting out different parameters related to node movement in order to evaluate in different metrics performance in different routing protocols. These parameters are for examples the starting location, the nodes movement direction, velocity range, speed changes over time and so on. Mobility models can be categorized into two types and they are as follows [6]:

- Entity Models

- Group Models

When the mobile nodes are going to apply Entity models, the movements of the mobile nodes are completely independently from each other. In contrary, when the mobile nodes are using in group models, the mobile nodes are dependent on each other and follow predefined leader node.

One of the prominent group models is Reference Point Group Mobility (RPGM) model which represents the random motion of a group of mobile nodes and their random individual motion within the group. In this model, the group motion behavior is determined by a logical group center where all group members have to follow that group leader. The individual mobile nodes movements or the entity mobility models should be specified the way of the movement of the individual mobile nodes. These movements are independent and each every of mobile nodes can move to its way randomly with randomly velocity. The principle of that the logical group center for RPGM model is to guide group of nodes continuously calculating group motion vector $\overline{GM}$ in order to define; speeds and directions for mobile nodes. However, once the updated reference point RP (t+1) has been updated they are combined with random motion vector $\overline{RM}$ values to represent the random motion of each mobile node around its reference point [7].

### III. ROUTING PROTOCOLS

In Ad hoc networks, entire network nodes are needed to perform as routing functions. To manage with the dynamic nature of the topology of Ad hoc networks, several routing protocols have been proposed. Taking into consideration of procedures for route establishment and update, MANET routing protocols can be categorized into three types:

- Proactive

- Reactive and

- Hybrid protocols

Proactive (table-driven) protocols maintain the routing information consistently up-to-date from each node to every other node in the network. The main function of proactive routing protocol maintains its table in order to store routing information. Up on changing in the network topology caused by anything just need to be reflected to this table and propagate the updating information throughout the network. Reactive (on demand) protocols are based on source-initiated on-demand reactive routing. The reactive routing protocol initiates routes only when mobile node requires a route to a destination. This method works like that the node broadcasts a route discovery initiative on demand and finishes when the destination is fund. Hybrid protocols are combination of proactive and reactive protocols. An example for proactive protocol is DSDV. This protocol routing protocol, all the possible destinations and the number of hops in the network are stored in a table. Because of ad hoc networks instabilities the table updating may change extremely dynamically and updating advertisements might be caused more network congestion and increases network routing overhead as well. AODV is a reactive protocol that improves the DSDV in the sense of minimizing the number of extremely advertisements and an unnecessary hello messages by creating routes on a demand mechanism [7]. Three main of Route Request process for AODV protocol are:

- Route request (RREQ)

- Route reply (RREP)

- Route error (RERR)

Figure 1shows that the source broadcasts a route request to its entire neighbors and each of these neighbors forwards the request until it reaches at the destination. Figure 2 illustrates the reverse route propagating to source. When the destination receives a routing packet from the source, it replies the shortest path and the source will use this path to send data packets. Figure 3 presents a route error. After the path between source and destination is established, the data is transferring. But due to the Ad hoc network characteristic of rapid topology change, the link between the source and the destination breaks. When such accident happens in Ad hoc networks using AODV routing protocol the current node prepares an error message propagate to the source. In this case the current node is the node 2 in Figure 3 and the link between node 2 and node 3 is broken, thus the node 2 prepared an error message and sent to the source.[8]



Figure 1.        RREQ broadcast



Figure 2.        RREP

Figure 3.        Route Failure

## IV.    NETWORK CONFIGURATION

In this section will be analyzed the protocol performance optimal setting using Reference Point Group Mobility model (RPGM) and set up network configuration. These experiments have been conducted to finalize the maximum pause time optimal setting for protocal performance.  In this set of simulations, our intention is to investigate the protocol's performance under RPGM when the pause time increases. The results present the performance of the protocol is very much depending on pause time value. From that point of view it is very important to run several simulations with different parameters in order to find out the suitable value of pause time for the Protocol. Network set up, we increase the pause time from 5 up to 40 sec 40 and keep all other parameters unchanged in first scenario. Keep in mind that for each scenario we keep running five times and these five times the pause time increases (5, 10, 20, 30 & 40 sec) while all other parameters are fixed.  In the 2nd, third and fourth scenarios, we increase the number of nodes by double (from 10 to 20 nodes, 20 to 40 and 40 to 80) and number of connection is also increases as well (4, 8, 30 & 40 connections) respectively and of course for each scenario the pause time increases while the rest of network parameters unchanged. The following tables are indicated the scenarios that have been tested in this experiment.

The following tables are four main scenarios and each of them contains five sub-scenarios. Each of these sub-scenarios, the pause time is varying while the rest of the parameters are constant or unchanged. When the first scenario has tested the five sub-scenarios with varying pause time, we continue, the second main scenario with increasing the number of nodes and number of connection with 20 nodes and 8 connections respectively as the table1-2 has shown. This second scenario has also tested same as tested the scenario one with five sub-scenarios except that increasing the number of nodes and connections and of course the pause time is also varying. Subsequently the third and fourth scenarios keep carrying on with increasing the number of nodes and connections 40 nodes and 80 nodes with 30 and 40 connections respectively as the table1-3 and table1-4 are indicated.

TABLE I.        VARYING PAUSE TIME (RPGM-MODEL)

| Parameters | Values |
|---|---|
| Simulation | Ns-2 |
| Protocol | AODV |
| Movement Model | RPGM |
| Traffic source | Constant Bit rate (CBR) |
| Simulation area | 1000 m x 1000 m |
| Simulation Time | 200 sec |
| # of nodes | 10 |
| Pause time | 5, 10, 20, 30 & 40 |
| Nodes speed | 20m/s |
| # node of sourceS | 4 |

TABLE II.        VARYING PAUSE TIME (RPGM-MODEL)

| Parameters | Values |
|---|---|
| Simulation | Ns-2 |
| Protocol | AODV |
| Movement Model | RPGM |
| Traffic source | Constant Bit rate (CBR) |
| Simulation area | 1000 m x 1000 m |
| Simulation Time | 200 sec |
| # of nodes | 20 |
| Pause time | 5, 10, 20, 30 & 40 |
| Nodes speed | 20m/s |
| # node of sources | 8 |

TABLE III.        VARYING PAUSE TIME (RPGM-MODEL)

| Parameters | Values |
|---|---|
| Simulation | Ns-2 |
| Protocol | AODV |
| Movement Model | RPGM |
| Traffic source | Constant Bit rate (CBR) |
| Simulation area | 1000 m x 1000 m |
| Simulation Time | 200 sec |

| # of nodes | 40 |
|---|---|
| Pause time | 5, 10, 20, 30 & 40 |
| Nodes speed | 20m/s |
| # node of sources | 30 |

TABLE IV.        VARYING PAUSE time (RPGM-MODEL)

| Parameters | Values |
|---|---|
| Simulation | Ns-2 |
| Protocol | AODV |
| Movement Model | RPGM |
| Traffic source | Constant Bit rate (CBR) |
| Simulation area | 1000 m x 1000 m |
| Simulation Time | 200 sec |
| # of nodes | 80 |
| Pause time | 5, 10, 20, 30 & 40 |
| Nodes speed | 20m/s |
| # node of sources | 40 |

## V.    PERFORMANCE METRICS

As RFC 2501 described, a number of performance metrics that can be used for evaluating the performance of a routing protocol for MANETs are four metrics. In addition, AODV developers were used these four metrics [9]. Thus, in this paper, we follow the general ideas described in RFC 2501 and we used similar metrics such as Packet Delivery Ratio, Average End-to-End Delay, Normalized Routing Load and Routing Overhead (Normalized MAC Load). All these performance metrics are important, but the packet delivery ratio and average end-to-end delay are most important for best-effort traffic. It does not mean that the other metrics are meaningless, but it means that the two first metrics have high priority than others. However, the normalized routing load is important as it will be used to evaluate the efficiency of the routing protocol. In order to calculate the Packet Delivery Ratio, we collect all received data packets which are delivered at the destination node and sum together than divide the total Data packets sent by source associated with the agent type AGT. As proactive protocols are normaly expected to have a higher control overhead than reactive protocols, the TCP data traffic application can cause additional packet loss and network congestion at the intermediate nodes thus we used CBR (continuous bit –rate) as data traffic generation in this paper. Calculating Average End-to-end delay, At the destination side  we summarize the all receive packets(Recv-Pkts) then store the transmission time for each successful data and extracts the receiving time of that packet. These times are included all possible delays in the network such as broadcasting latency, the intermediate nodes retransmission delay, processing delay, queuing delay and propagation delay. Based on this information we calculate the end-to-end delay to summarize all those delays and divide by the number of Recv-Pkts. Normalized MAC Load(NML), can be define as the fraction of all control packets (routing control packets, included  Request-To-Send (RTS) which means a request signal to next neighbor asking for if he can be ready to receive his request, Clear-To-Send (CTS) which means an acceptance for that request, Address Resolution Protocol (ARP) requests and replies, and MAC ACKs (incase using TCP ) divide the total number of received data packets. The normalized routing load (NRL), this metric is different from Normalized MAC load (NML) and sometimes is difficult to distinguish each other. The NRL concerns only received data packets at the destination while NML take into the account all establishing process for control packets in routing level not the agent (AGT) level. The way of calculation of NRL is defined as the fraction of all routing control packets sent by all nodes divide the number of received data packets at the destination nodes.

## VI.    RESULTS AND ANALYZES

### A.   Routing Overhead Using Rpgm

This metric presents how the network enquiry packets control for intermediate nodes is huge. To evaluate the network efficient, the routing overhead or Normalized MAC Load plays a significant roll. Thus it is vital for this paper to figure out the vectors which can be involved the increasing of the network routing overhead. However, minimizing the network overhead is maximizing the network resources utilization with better performance.  Figure 4 presents the experiment result for routing overhead. The Figure 4 shows that the overhead is direct proportional to pause time. It indicates that the overhead increases while the pause time is increased. This experiment tested five different values of pause time, namely 5, 10, 20, 30 & 40 sec. The figure 4 illustrated four scenarios and each has a different value than others. These scenarios have been analyzed one by one as we get a better overview.

Look at the scenario one, when the pause time is 5 sec the overhead is 54, when we increase the pause time to 10 sec the overhead became 186, then 330, 408 and finally when the pause time increase up to 40 sec the overhead is worst. Second scenario, this scenario is totally different from previous scenario because the overhead started 442 while scenario one the overhead started 54, because of varying node density in the network.  But the principle is stil same because when the pause time is increased the overhead is also increased. The Scenarios three and four are some as these two scenarios, increase the pause time, the overhead goes up. So, seeing the scenarios in the Figure 4 the overhead is definatly proportional to pause time and obviously the overhead will be low when the pause time is small. Thus, this experiment shows that choosing a smaller pause time, gives the optimal setting for node pause time and indicates that the average performance of AODV protocol is when the node pause time is between 5 sec and 10 sec.

Figure 4.        Routing overhead using

### A.   *Packet Delivery Fraction Using RPGM*

Figure 5 shows the result of Packet Delivery Fraction (PDF). This part of analyzing is important as it describes the rate of packets drop as well as it affects the overall network throughput that the network can support. Figure 5 shows that the packets delivery ratios is mostly constant whether the the pause time increases or not. Only the scenario one shows little bit flactuation where the pause time for 10 and 20 are 95.06% and 93.35 % respectively. Scenarios two, three and four have shown the best result and most values are between 98 % and 99%.  Nonetheless, the overall results for the packets delivery ratios based on RPGM model indicate that the packet delivery is very smooth for all scenarios.



Figure 5.        Packet Delivery Fractions

### B.   *Normalizing Routing Load Using RPGM*

Figure 6 shows the simulation for Normalized Routing Load with the RPGM model.  The objective of this part of analyzing is to investigate the impact of varying pause time on the normalizing route load in the network. The Figure 6 shows the scenarios of the simulation for normalizing routing load. It shows that as the pause time increases, the normalized routing load increases.  For instance, when the pause time is  5 in scenario one the normalized routing load  shows 0.05 % . Then we increased the pause time up to 10, 20, 30 and 40 in same scenario, the results show 0.18%, 0.33 %, 0.39 % and 0.40 % respectively. This analyzing indicates that as the pause time goes up, the result gets worst which meas that the network becomes more and more congested in terms of traffic load.

The rest of scenarios are almost same as previous scenario (scenario one) and show when ever the pause time increases the result of NRL increases, but scenarios two and three have shown when the pause time is 30 and 10 the result do not follow the previous results which were increasing when the pause time increases instead they decrease in this time as can be seen in the Figure 6 the results show 0.16 % and 0.34 %. Since the most of results are shown the same direction when the pause time is increased, we can assume these two cases are accidently happen because of the ad hoc network behavior is very changeable. NRL represents the number of routing packets transmitted per data packet delivered at the destination, as it to evaluate the efficient of protocol's performance. However, in this experiment indicates that when the pause time is between 5 and 10 sec is the optimal setting for AODV protocol when using RPGM model.



Figure 6.        Normalized Routing Load

### C.   *Result for Average end to End Delay*

Figure 7 illustrates average end to end delay result for RPGM model. The Figure 7 shows that the delay of scenario one where the pause times are varying from 5, 10, 20, 30 and 40 sec. It shows when the the pause time is 30 the average end

to end delay is 11.81 %. So, the pause time is increased to10 and the result for this scenario is also increased up to 113.23 %. Then we increased the pause time up to 20, so the delay is also increased up to 204.43%. Again, the pause time is increased from 20 to 30 but as the figure 7 indicates in scenario one the result decreased 178.44%. Finally, we increased the pause time up to 40 and result decreased as well 135.22%. It seems that the average end to end delay is not that much affected by the varying of pause time. We expected that the lower pause time will have a better performance average end to end delay than the higher pause time according to other parameters results in Figures 4, 5 & 6. Normally, the network parameters are depending on one other, and usually there is no one parameter that cannot be affected by another vector in order to be independent completely. However, the explanation lies in that the delay can be independent to some extent.



Figure 7.    Average end to end Delay

### D.  Packet Loss

The Figure 8 presents data packets drop. In ttheoretical, the reactive routing protocols and proactive routing protocols differ drastically in the fact that they belong to two different routing families. Reactive routing protocol generates or flood packets on demand only, in order to reduce routing loads. Proactive protocol frequently updates the routing tables regardless of need. This proactive periodically message exchanging, causes tremendous routing overhead. These overhead can lead dropped data packets in the network. In contrary, the proactive routing protocols are expected less routing overhead because there is no need periodically information exchange to update the tables. This minimizing of information exchange for reactive routing protocol is an advantage for reactive routing protocol and will lead minimizing dropped data packets. However, the proactive routing protocols have also different behavior. AODV has more overhead than other reactive routing protocols caused by AODV route query and routing overhead is proportional to the number of route queries. For instance, the source node for

AODV will send a RREQ message if it does not know the route to the destination. The source waits for a while, if the source sends the second transmission of the RREQ message and does not receive RREP message within a time interval, it will drop the first packet in the queue and repeats the same procedure for the second data packet in the queue.

In addition, if one of the forwarding nodes cannot success to find a valid route to the destination will drop all data packets from its queue. Here we are going to examine the pause time optimal setting for AODV protocol in packet drop using RPGM as a model. The results have shown that the dropped data packets in this experiment are very reasonable and much better than expected results. Because, as above mentioned the AODV protocol has heavy routing overhead which can cause a significant drop data packets. The Figure 8 shows that dropped data packets for the scenario one are 2, 59, 90, 55 &59. First experiment in this scenario has only 2 packets lost while the second third, fourth and fifth experiments have 59, 90, 55 and 59 packets lost respectively. Drops for scenario two are 49, 68, 83, & 158 while third scenario results are 81, 83, 61 & 170 and 167. However, the three first scenarios optimal setting is when the pause times are 5 and 10.



Figure 8.    Packet drop

### VII.  CONCLUSION

This paper, we have conducted several experiments and analyzed the effects of pause time varying to evaluate the performance of AODV protocol based on RPGM as a reference model. The simulation results evaluated the performance of the routing protocol with regard five performance metrics such as Packet Delivery Ratio, Average end to end delay, Packet Drop, Routing overhead and Normalized routing load. The experiment results for all scenarios present that the low pause time, the high performance of the protocol. The experimentation also suggests that several parameters such as traffic patterns, node

density are also affect the routing performance and need to be investigated with various scenarios. Further study also needs to be done with additional analysis with different mobility models.

### REFERENCES

[1] Yang Huan , Jiang Hong , Liu Lei "Performance Analysis of Mobility Models in Sparse Ad-hoc Networks" School of Information Engineering, Southwest University of Science and Technology, Mian Yang 621010,P.R.China, Proceedings of the 27th Chinese Control Conference July 16-18, 2008, Kunming,Yunnan, China,2008.

[2] Tariq M M B,Ammar M and Zegura E. "Message Ferry Route Design for Sparse Ad hoc Networks with Mobile Nodes", In Proc. The 6th ACM International Symposium on Mobile Ad Hoc Networking and Computing (Mobi-Hoc'06), 2006.

[3] C.P.Agrawal et al "Evaluation of Varrying Mobility Models & Network Loads on DSDV Protocol of MANETs" Raipur, India, International Journal on Computer Science and Engineering Vol.1(2), 2009, 40-46.

[4] S. Sesay et al., "Simulation Comparison of Four Wireless Ad Hoc Routing Protocols",Information Technology Journal 3(3), 2004, pp:219-226.

[5] S. Shah et al, "Performance Evaluation of Ad Hoc Routing Protocols Using NS2 Simulation",Conf. of Mobile and Pervasive Computing, 2008.

[6] T. Camp et al., "A Survey of Mobility Models for Ad Hoc Network Research", Wireless Comm. & Mobile Computing: Special issue on Mobile Ad Hoc Networking: Research, Trends and Applications, vol. 2, pp. 483-502, 2002.

[7] Valentina Timcenko, Mirjana Stojanovic, Slavica Bostjancic Rakas,"MANET Routing Protocols vs. Mobility Models: Performance Analysis and Comparison" Proceedings of the 9th WSEAS International Conference on Applied Informatics and Communications (AIC '09).

[8] Jungkeun Yoon, Mingyan Liu, Brian Noble "Random Waypoint Considered Harmful" Electrical Engineering and Computer Science Department University of Michigan Ann Arbor, Michigan 48109-2122, IEEE INFOCOM 2003.

[9] Das, S.R., Perkins, C.E. Royer, E.M.2000, "Performance comparison of two on-demand routing protocols for ad hoc networks ", INFOCOM 2000. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings IEEE Volume 1, pg. 3 – 12, 26-30 March 2000.

## AUTHORS PROFILE



**Sayid Mohamed Abdule,** received his BSc degree in Computer Science from Agder University, Norway. And, his M.Sc degree in Satellite Communication focusing Quality of Service of VOIP over Satellite from the University Sains Malaysia (School of Computer Science), MALAYSIA. Sayid currently attached to the InterNetWorks Research Group at the UUM College of Arts and Sciences as a doctoral researcher. He is currently pursuing his PhD research in Ad-hoc Mobile networking. His current research interest is on Ad-hoc mobile network routing protocol.



**Associate Professor Dr. Suhaidi Hassan** is currently the Assistant Vice Chancellor of the College of Arts and Sciences, Universiti Utara Malaysia (UUM). He is an associate professor in Computer Systems and Communication Networks and the former Dean of the Faculty of Information Technology, Universiti Utara Malaysia. Dr. Suhaidi Hassan received his BSc degree in Computer Science from Binghamton University, New York (USA) and his MS degree in Information Science (concentration in Telecommunications and Networks) from the University of Pittsburgh, Pennsylvania (USA). He received his PhD degree in computing (focusing in Networks Performance Engineering) from the University of Leeds in the United Kingdom. In 2006, he established the ITU-UUM Asia Pacific Centre of Excellence (ASP CoE) for Rural ICT Development, a human resource development initiative of the Geneva-based International Telecommunication Union (ITU) which serves as the focal point for all rural ICT development initiatives across Asia Pacific region by providing executive training programs, knowledge repositories, R&D and consultancy activities. Dr. Suhaidi Hassan is a senior member of the Institute of Electrical and Electronic Engineers (IEEE) in which he actively involved in both the IEEE Communications and IEEE Computer societies. He has served as the Vice Chair (2003-2007) of the IEEE Malaysia Computer Society. He also serves as a technical committee for the Malaysian Research and Educational Network (MYREN) and as a Council Member of the Cisco Malaysia Network Academy.



**Osman Ghazali, Ph.D.** is a Senior Lecturer at Northern University of Malaysia (Universiti Utara Malaysia). He obtained his Bachelor and Master degrees in Information Technology in 1994 and 1996 from the Northern University of Malaysia. In 2004 he obtained his PhD specializing in Computer Network from Northern University of Malaysia. As an academician, his research interests include congestion control, quality of services, wired and wireless network, transport layered protocols and network layered protocols. His works have been published in international conferences, journals and won awards on research and innovation competition in national and international level.

**Mohammed M. Kadhum, Ph.D.** is an assistant professor in the Graduate Department of Computer Science, Universiti Utara Malaysia (UUM) and is currently attached to the InterNetWorks Research Group at the UUM College of Arts and Sciences as a research advisor. He had completed his PhD research in computer networking at Universiti Utara Malaysia (UUM). His research interest is on Internet Congestion and QoS. He has been awarded with several medals for his outstanding research projects. His professional activity includes being positioned as Technical Program Chair for NetApps2008 and NetApps2010, a technical committee member for various well known journal and international conferences, a speaker for conferences, and a member of several science and technology societies. To date, he has published a number of papers including on well-known and influential international journals.

# Automating Legal Research through Data Mining

M.F.M Firdhous,

Faculty of Information Technology,
University of Moratuwa,
Moratuwa,
Sri Lanka.
Mohamed.Firdhous@uom.lk

*Abstract* —**The term legal research generally refers to the process of identifying and retrieving appropriate information necessary to support legal decision-making from past case records. At present, the process is mostly manual, but some traditional technologies such as keyword searching are commonly used to speed the process up. But a keyword search is not a comprehensive search to cater to the requirements of legal research as the search result includes too many false hits in terms of irrelevant case records. Hence the present generic tools cannot be used to automate legal research.**
**This paper presents a framework which was developed by combining several 'Text Mining' techniques to automate the process overcoming the difficulties in the existing methods. Further, the research also identifies the possible enhancements that could be done to enhance the effectiveness of the framework.**

*Keywords* —*Text Mining; Legal Research; Term Weighting; Vector Space*

## I. INTRODUCTION

Legal research is the process of identifying and retrieving information necessary to support legal decision-making. In its broadest sense, legal research includes each step of a course of action that begins with an analysis of the facts of a problem and concludes with the application and communication of the results of the investigation [1].

The processes of legal research vary according to the country and the legal system involved. However, legal research generally involves tasks such as finding primary sources of law, or primary authority, in a given jurisdiction (cases, statutes, regulations, etc.), searching secondary authority for background information about a legal topic (law review, legal treatise, legal encyclopedias, etc.), and searching non-legal sources for investigative or supporting information. Legal research is performed by anyone with a need for legal information including lawyers, law librarians, law students, legal researchers etc., Sources of legal information range from decided cases, printed books, to free legal research websites and information portals [2].

Manually performing legal research is time consuming and difficult. Because of that, some traditional tools have been introduced. There are essentially only two types of tools which help users find legal materials in the Internet, they are commonly known as catalogs and search engines. Combinations of catalogs and search engines in the same site are now becoming more common, and such combinations are often referred to as 'portals'. Despite the existence of these research aids, finding legal information on the internet is surprisingly difficult, partly because neither catalogs nor search engines used alone can provide a satisfactory solution. A general keyword search contains too much false hits. It is difficult to make searches precise enough to find only the relevant information.

In this case, the search should not be just a keyword search; instead an intelligent search should be carried out according to the meaning of the search text. This research presents a methodological framework based on text mining to automate legal research by focusing on the retrieval of exact information specifically necessary for legal information processing. The proposed approach uses a term-based text mining system and a vector space model for the development of the framework.

## II. RELATED WORK

Data mining is the process of sorting through large amounts of data and picking out relevant information. It is usually used by business intelligence organizations, and financial analysts, but is increasingly being used in the sciences to extract information from the enormous data sets generated by modern experimental and observational methods. It has been described as "the nontrivial extraction of implicit, previously unknown, and potentially useful information from data" and "the science of extracting useful information from large data sets or databases" [3].

Text mining, sometimes alternately referred to as text data mining, refers generally to the process of deriving high quality information from text. High quality information is typically derived through the division of patterns and trends through means such as statistical pattern learning. Text mining usually involves the process of structuring the input text (usually parsing, along with the addition of some derived linguistic features and the removal of others, and subsequent insertion into a database), deriving patterns within the structured data, and finally evaluation and interpretation of the output. 'High quality' in text mining usually refers to some combination of relevance, novelty, and interestingness. Typical text mining tasks include text categorization, text clustering, and concept/entity extraction, production of granular taxonomies, sentiment analysis, document summarization, and entity relation modeling [4–5].

### A.  Dependency Analysis based Text Mining

One of the approaches in text mining is dependency based analysis. In [6–8] several methodologies in this area are presented. In [6], a text simplification approach is presented, whereas [7–8] present dependency analysis. Long and complicated sentences pose various problems to many state-of-the-art natural language technologies. For example, in parsing, as sentences become syntactically more complex, the number of parses increases, and there is a greater likelihood for an incorrect parse. In machine translation, complex sentences lead to increased ambiguity and potentially unsatisfactory translations. Complicated sentences can also lead to confusion in assembly/use/maintenance manuals for complex equipment [6].

Articulation-points are defined to be those points where sentences may be split for simplification. Segments of a sentence between two articulation points may be extracted as simplified sentences. The nature of the segments delineated by the articulation points depends on the type of the structural analysis performed. If the sentences are viewed as linear strings of words, articulation points can be defined to be, say, punctuation marks. If the words in the input are also tagged with part of speech information, sentences can be split based on the category information, for instance at relative pronouns, with part of speech information, subordinating and coordinating conjunctions may also be detected and used as articulation points. However, with just this information, the span of the subordinating/coordinating clause would be difficult to determine. On the other hand, if the sentence is annotated with phrasal bracketing, the beginnings and ends of phrases could also be articulation points.

The sentences in the training data are first processed to identify phrases that denote names of people, names of places or designations. These phrases are converted effectively to single lexical items. Each training sentence Si, along with its associated j (simplified) sentences Si1 to Sij, is then processed using the Lightweight Dependency Analyzer (LDA) [7].

The resulting dependency representations of Si and Si1 through Sij are 'chunked'. Chunking collapses certain substructures of the dependency representation (noun phrases and verb groups) and allows defining the syntax of a sentence at a coarser granularity. Chunking also makes the phrasal structure explicit, while maintaining dependency information. Thus, this approach has the benefit of both phrasal and dependency representations.

LDA is a heuristic based, linear time, deterministic algorithm which is not forced to produce dependency linkages spanning the entire sentence. LDA can produce a number of partial linkages since it is driven primarily by the need to satisfy local constraints without being forced to construct a single dependency linkage that spans the entire input. This, in fact, contributes to the robustness of LDA and promises to be a useful tool for parsing sentence fragments that are rampant in speech utterances exemplified by the switchboard corpus [7].

### B.  Text Mining at the Term Level

Most efforts in Knowledge Discovery in Databases (KDD) have focused on knowledge discovery in structured databases, despite the tremendous amount of online information that appears only in collections of unstructured text. At abstract level, KDD is concerned with the methods and techniques for making sense of data. The main problem addressed by the KDD process is mapping low-level data into other forms that might be more compact, more abstract, or more useful. At the core of the process is the application of specific data-mining methods for pattern discovery and extraction. Previous approaches to text mining have used either tags attached to documents [9] or words contained in the documents [10].

The exploitation of untagged, full text documents therefore requires some additional linguistic pre-processing, allowing the automated extraction from the documents of linguistic elements more complex than simple words. Normalized terms are used here, i.e. sequences of one or more lemmatized word forms (or lemmas) associated with their part-of-speech tags. "stock/N market/N" or "annual/Adj interest/N rate/N" are typical examples of such normalized terms [11]. In [11], an approach is presented to text mining, which is based on extracting meaningful terms from documents. The system described in this paper begins with collections of raw documents, without any labels or tags. Documents are first labeled with terms extracted directly from the documents. Next, the terms and additional higher-level entities (that are organized in a hierarchical taxonomy) are used to support a range of KDD operations on the documents. The frequency of co-occurrence of terms can provide the foundation for a wide range of KDD operations on collections of textual documents, such as finding sets of documents whose term distributions differ significantly from that of the full collection, other related collections, or collections from other points in time.

The next step is the Linguistic Preprocessing that includes Tokenization, Part-of- Speech tagging and Lemmatization. The objective of the Part-of-Speech tagging is to automatically associate morpho-syntactic categories such as noun, verb, adjective, etc., to the words in the document. In [12], a Transformation-Based Error-Driven Learning approach is presented for Part-of- Speech tagging. The other modules are Term Generation and Term Filtering.

In the Term Generation stage, sequences of tagged lemmas are selected as potential term candidates on the basis of relevant morpho-syntactic patterns (such as "Noun Noun", "Noun Preposition Noun", "Adjective Noun", etc.). The candidate combination stage is performed in several passes. In each pass, association coefficient between each pair of adjacent terms is calculated and a decision is made whether they should be combined. In the case of competing possibilities (such as (t1 t2) and (t2 t3) in (t1 t2 t3)), the pair having the better association coefficient is replaced first. The documents are then updated by converting all combined terms into atomic terms by concatenating the terms with an underscore. The whole procedure is then iterated until no new terms are generated [11].

In generating terms, it is important to use a filter that preserves higher precision and recall. The corpus is tagged, and a linguistic filter will only accept specific part-of-speech sequences. The choice of the linguistic filter affects the precision and recall of the results: having a 'closed' filter, which is strict regarding the part-of-speech sequences it accepts, (like only Noun +...) will improve the precision but will have a bad effect on recall [13]. On the other hand, an 'open' filter, which accepts more part-of-speech sequences, such as prepositions, adjectives and nouns, will have the opposite result. In [14], a linguistic filter is chosen, staying somewhere in the middle, accepting strings consisting of adjectives and nouns:

*(Noun | Adjective) + Noun*

However, the choice of using this specific filter depends on the application: the construction of domain-specific dictionaries requires high coverage, and would therefore allow low precision in order to achieve high recall, while when speed is required, high quality would be better appreciated, so that the manual filtering of the extracted list of candidate terms can be as fast as possible. So, in the first case we could choose an 'open' linguistic filter (e.g. one that accepts prepositions), while in the second, a 'closed' one (e.g. one that only accepts nouns). The type of context involved in the extraction of candidate terms is also an issue. At this stage of this work, the adjectives, nouns and verbs are considered [14].

### C. Term Weighting

The Term Generation stage produces a set of terms associated with each document without taking into account the relevance of these terms in the framework of the whole document collection. The goal of term weighting is to assign to each term found a specific score that measures the importance, with respect to a certain goal, of the information represented by the term.

The experimental evidence accumulated over the past several years indicates that text indexing systems based on the assignment of appropriately weighted single terms produce retrieval results that are superior to those obtainable with other more elaborate text representations. These results depend crucially on the choice of effective term weighting systems. The main function of a term weighting system is the enhancement of retrieval effectiveness. Effective retrieval depends on two main factors: on the one hand, items likely to be relevant to the user's needs must be retrieved; on the other hand, items likely to be extraneous must be rejected [15]. The research in this area have found various term weighting schemes. In [15], the method called "TF-IDF Weighting Scheme" is described and [14] describes a method called "C-Value Method".

### III. APPROACH

The input to the proposed system is a collection of law reports. Law reports consist of two sections; namely the head and the detail section. The head section summarizes the whole law report and the detail section contains the detailed information about the case. Only the head section is used for automated processing as it contains sufficient details for the purpose.

Figure 1 shows the overall architecture of the proposed system.



Figure 1: Overall Architecture of the Proposed System

The proposed system consists of two main components, namely;

a. The mining process
b. The research process

The mining process is the main process in the framework and to be completed prior to the research process. The mining process is carried out on the entire collection of the law reports of the repository. In this process, each document is analyzed and information that should be used for legal research is recorded in the processed law reports repository. Then the research process is carried out on the processed law reports. In this process, the text block is analyzed and the required information is extracted and compared with each law report to identify the matching reports.

The proposed approach uses text mining. More precisely, it uses terms level text mining, which is based on extracting meaningful terms from documents [11]. Each law report is represented by a set of terms characterizing the document.

### A. The Mining Process

Figure 2 shows the architecture of the mining process. The mining process goes through the stages shown in Figure 2 and is based on the approach presented in [11]. The law reports are stored in a repository and the mining process works on that. Basically, one report at a time is taken, processed and the required information are stored back along with the report.

Figure 2: Architecture of the Mining Process

The first step is the Linguistic Preprocessing that includes Tokenization and Part-of-Speech tagging. Here, the head part of the law report is tokenized into parts and morpho-syntactic categories such as noun, verb, adjective etc., are associated into the words of the text. Each word in the corpus is given a grammatical tag as noun, verb, adjective, adverb or preposition, which is a prerequisite step to the next stage. A list of predefined tags is used for tagging the words. Table A-1 in Appendix lists the set of speech tags used.

Along with part-of-speech tagging, the other preprocessing stage is 'Chunking'. A chunk is a syntactically correlated part of a language (i.e. noun phrase, verb phrase, etc.). Chunking is the process of identifying those parts of language. Part-of-speech tagging and chunking are bound together. As with part-of-speech tagging, chunking also uses a tag list and is given in Table A-2 in Appendix. A chunk tag is defined in the following format.

<Prefix> - <Chunk Type>

Figure 3 shows a sample block of text that has been speech tagged and chunked.



Figure 3: A sample Part-of-Speech Tagged Text Block

## B. Term Generation

In the Term Generation stage, sequences of tagged lemmas are selected as potential term candidates on the basis of relevant morpho-syntactic patterns (such as "Noun Noun", "Adjective Noun", etc.).

For the retention of higher accuracy and recall, only two word terms are generated as the following.

(Adj | Noun) + Noun

When the term includes more number of words, then it will have higher precision but lower recall. If the term includes only one word, then it will have higher recall but lower precision. To obtain a balanced precision and recall, by staying in the middle only two word terms are used in this approach [11],[13].

A stop word is a very common word that is useless in determining the meaning of a document. Stop words do not contribute a value for the meaning when selected as a term. So, when that kind of words are included in the generated term set, those terms are removed from the term set to preserve a clean valuable term set and avoid misleading meanings. In the context of legal research, the words contained in Table 1 are considered as stop words.

TABLE I: STOP WORDS IN THE CONTEXT OF LEGAL RESEARCH

| Stop Words |
| --- |
| Petitioner |
| Complainant |
| Plaintiff |
| Plaintiff-respondent |
| Court |

The Term Generation stage produces a set of terms associated with each law report without taking into account the relevance of these terms in the framework of the whole law report collection.

The goal of term weighting is to assign each term found a specific score that measures the importance, relative to a certain goal, of the information represented by the term. A term with a higher weight is a more important term then other terms. The TF-IDF weighting method presented in [11] is used in this framework due to its performance and easiness.

## C. The Research Process

Research process shown in Figure 4 comes into action when the mining process has been completed. This is the search process used in the proposed framework to find out the relevant law reports, according to the user's legal research criteria.

Figure 4: The Research Process

The research process is based on the input text block for the legal research given by the user. The process works on the mined repository of law reports, thus the mining process is compulsory before the research stage. The outcome is the set of matching law reports, sorted according to the relevance of the reports to the input research text.

The research process mainly works on the input text block for the legal research. A similar process as of mining the law reports is carried out on the input text block. All the steps in the mining process including linguistic pre-processing, term generation, term weighting are done on the input text. In short, the input text is treated like another document.

After generation of terms and weight assignment, the next stage is document comparison. Vector space model is used for comparing the input text (treating it as a document) against the law reports in the repository. The law reports are represented in terms of vectors (keywords) in a multi-dimensional space as shown in Figure 5, and the theory called the "Cosine Similarity" is used for comparison [16–23]. Using the cosine formula shown in Figure 6, the cosine angle between the query document and each law report in the repository is computed. Then the matching reports are sorted in descending order based on the cosine value that ranges between 0 and 1. This process brings the most relevant law report to the query, to the top.



Figure 5: Representation of Documents in Vector Space

$$cos(\theta) = \frac{v_1 . v_2}{\|v_1\| \|v_2\|}$$

$$cos(\theta) = \frac{(x_1 \cdot x_2) + (y_1 \cdot y_2)}{(x_1^2 + y_1^2)^{1/2} + (x_2^2 + y_2^2)^{1/2}}$$

Figure 6: Formulas for Computing Cosine Similarity

## IV. PROTOTYPE IMPLEMENTATION

Prototype application was built using Java Platform, Standard Edition. For part-of-speech tagging, text processing library was used in a high-level manner, and the other parts of the architecture are implemented using object orientation.

Figure 7 shows the user interface of the software. The user interface can be divided in to two parts. The left side of the interface is used for the storage of law reports and text mining, while the right side of the interface has facilities for the research process.



Figure 7: User Interface of the Prototype Application

The text processing library used for part-of-speech tagging uses a model based approach for the functionality. At the startup of the application, the required components are initiated and the models and other resources are loaded. Once loaded, they can be used throughout the application reducing the loading time during operations. This improves the efficiency of the operations. Figure 8 shows the resource loading at the startup.

Figure 8: Resource Loading at the Startup of Application

When the law reports have been loaded to the repository, mining should be carried out. The prototype has the facilities for viewing intermediate results such as part-of-speech tagging details, terms, weights etc., However, when a new law report is loaded to the repository the mining process should be executed in order to update the mining results according to the updated state of the repository.

Figure 9 shows the result of the mining process with the intermediate results window open.



Figure 9: Results of Mining Process

Once the mining process is completed, the data is ready for the research process. The user text based on which the research process is to be carried out needs to be supplied to the application through the user interface. During the research process, first the user text will be analyzed and tagged as in the mining process. Then the law reports in the repository are

compared against the user text in order to find out the relevant law reports. The interface application has the capability to show the intermediate results of the research process as well.

Figure 10 shows the results after the completion of the research process which is the last stage in the entire legal research process. Similar to the results of the mining process, the text tagging of the user text can also be viewed as intermediate results at the end of the research process. At the end of the research process, the results are presented on the same interface with all the law reports with relevant information. The law reports were organized in a descending order based on the similarity score with the most relevant law report on top. By selecting a law report in the interface, it is also possible to see the information (matching terms) based on which the user text and the law report were matched highlighted and the similarity score computed for the reports. The bottom most pane of the interface shows the verdict of the case. This makes it easier for a reader to immediately know the final judgment of the case without going through the reports separately.



Figure 10: Final Results of the Legal Research Process

## V. EVALUATION

The evaluation of the framework was carried out using the prototype application described under Section IV. For the purpose of evaluation, only the fundamental rights cases filed at the Supreme Court of Sri Lanka were used. The reason for using the fundamental rights cases was the relative easiness of finding the case records as the Supreme Court and the Court of Appeal are superior courts and the courts of record in Sri Lanka and the Supreme Court is the final appellate court and its rulings bind all the lower courts in Sri Lanka. Also only the Supreme Court has the jurisdiction to hear fundamental rights petitions in Sri Lanka. Hence, for the fundamental rights petitions, the Supreme Court is the court of first instance.

All the other types of cases like civil cases, criminal cases and intellectual property right violation cases need to be filed at lower courts such as the District Court, High Court and the

Commercial High Court respectively. These cases will be heard by the Court of Appeal and the Supreme Court as appeal applications against the verdict of lower courts. Hence, for proper of legal research to be carried out, multiple case records of the same case filed and argued at different courts need to be analyzed. Hence, it was decided to concentrate only on fundamental rights cases for simplicity.

Several fundamental rights case records were downloaded in their raw forms from the Lawnet web site [24] that hosts case records of all reported cases in Sri Lanka. All these case records or law reports were input to system and the legal research process was carried out using different user input text as search text.

The results of the evaluation showed that the accuracy of the reports retrieved were high. It could also be observed that the ordering of the law reports based on the similarity score was acceptable as most of the time the most relevant case record had the highest similarity score and came on top. This shows the high precision of the system finding the correct case record based on the user input. Figure 11 shows the results of the system for precision against the manual judgment for the most relevant case record.



Figure 11: Result of Evaluation for Precision

When the search text was modified without changing meaning, it also resulted in the same set of case records most of the time. This shows the high recall of case records for the same or similar user input. Figure 12 shows the results of the evaluation for recall.



Figure 12: Results of Evaluation for Recall

## VI. CONCLUSIONS

This paper presents the results of the research carried out to develop a framework to automate the often tedious time consuming process of legal research. The end result of the research is a framework which is based on a combination of several text mining techniques. Finally the framework developed was tested for accuracy using a prototype application and fundamental rights case records. Accuracy of the results in terms of precision and recall were shown to be very high.

As the future work, this can be extended to handle all types of law reports in addition to the fundamental rights cases. The accuracy can be further enhanced by using a comprehensively updated stop words list and a set of predefined terms to be introduced along with the dropping of candidate terms. The weighting scheme can also be upgraded to include context information.

## REFERENCES

[1] "Legal research in the United States: Wikipedia, the free encyclopedia" wikipedia.org;
http://en.wikipedia.org/wiki/Legal_Research_In_the_United_States.
[Accessed: June, 2009]

[2] "Legal research: Wikipedia, the free encyclopedia" wikipedia.org; http://en.wikipedia.org/wiki/Legal_research.
[Accessed: November, 2009]

[3] W.J. Frawley, G. Piatetsky-Shapiro, and C.J. Matheus, "Knowledge Discovery in Databases: an Overview" http://www.kdnuggets.com/gpspubs/aimag-kdd-overview-1992.pdf.
[Accessed: Nov, 2009]

[4] "Text Mining: Wikipedia, the free encyclopedia" wikipedia.org; http://en.wikipedia.org/wiki/Text_mining.
[Accessed: June, 2009]

[5] M. Sharp, "Text Mining" Rutgers University, School of Communication, Information and Library Studies; http://www.scils.rutgers.edu/~msharp/text_mining.htm.
[Accessed: December, 2009]

[6] R. Chandrasekar, and B. Srinivas, "Automatic induction of rules for text simplification", Knowledge-Based Systems, 1997.

[7] B. Srinivas, "A lightweight dependency analyzer for partial parsing", Natural Language Engineering, 1996.

[8] B. Srinivas, and K.J. Aravind, "Supertagging: An Approach to Almost Parsing", Computational Linguistics, 1999.

[9] R. Feldman, and I. Dagan, "KDT - Knowledge Discovery in Textual Databases", in Proceedings of 1st International Conference on Knowledge Discovery (KDD), 1995.

[10] B. Lent, R. Agrawal, and R. Srikant, "Discovering Trends in Text Databases", in Proceedings of 3rd International Conference on Knowledge Discovery (KDD), 1997.

[11] R. Feldman, M. Fresko, Y. Kinar, Y. Lindell, O. Liphstat, M. Rajman, Y. Schler, and O. Zamir, "Text Mining at the Term Level", in Proceedings of 2nd European Symposium on Principles of Data Mining and Knowledge Discovery, 1998.

[12] E. Brill, "Transformation-based error-driven learning and natural language processing: A case study in part-of-speech tagging", Computational Linguistics, 1995.

[13] I. Dagan, and K. Church, "Termight: identifying and translating technical terminology", in Proceedings of the 4th conference on Applied natural language processing, 1994.

[14] K. Frantzi, "Incorporating context information for the extraction of terms", in Proceedings of 8th Conference of the Association for Computational Linguistics, 1997.

[15] G. Salton, and C. Buckley, "Term-weighting Approaches in Automatic Text Retrieval", Information Processing and Management, 1988.

[16] E.W. Weisstein, "Law of Cosines: MathWorld--A Wolfram Web Resource" Wolfram Research, Inc., ; http://mathworld.wolfram.com/LawofCosines.html.
[Accessed: August, 2008]

[17] M. Ransom, "Vector Dot Products" algebralab.org; http://www.algebralab.org/lessons/lesson.aspx?file=Trigonometry_Trig VectorDotProd.xml.
[Accessed: September, 2009]

[18] E.W. Weisstein, "Dot Product: MathWorld--A Wolfram Web Resource" Wolfram Research, Inc., ; http://mathworld.wolfram.com/DotProduct.html.
[Accessed: August, 2009]

[19] "The Dot product" math.mit.edu; http://www-math.mit.edu/18.013A/HTML/chapter03/section03.html
[Accessed: September, 2009]

[20] E.W. Weisstein, "Vector: MathWorld--A Wolfram Web Resource" Wolfram Research Inc., ; http://mathworld.wolfram.com/Vector.html.
[Accessed: Aug, 2009]

[21] M.W. Berry, "Introduction to Vector-Space Models" Department of Computer Science, The University of Tennessee; http://www.cs.utk.edu/~berry/lsi++/node4.html.
[Accessed: September, 2009]

[22] E.W. Weisstein, "Vector Space: MathWorld--A Wolfram Web Resource" Wolfram Research, Inc.; http://mathworld.wolfram.com/VectorSpace.html.
[Accessed: August, 2009]

[23] "Vector space: Wikipedia, the free encyclopedia" wikipedia.org; http://en.wikipedia.org/wiki/Vector_space
[Accessed: August, 2009]

[24] Lawnet website: http://www.lawnet.lk
[Accessed: January 2010]

AUTHOR PROFILE

Mohamed Fazil Mohamed Firdhous is a senior lecturer attached to the Faculty of Information Technology of the University of Moratuwa, Sri Lanka. He received his BSc Eng., MSc and MBA degrees from the University of Moratuwa, Sri Lanka, Nanyang Technological University, Singapore and University of Colombo Sri Lanka respectively. In addition to his academic qualifications, he is a Chartered Engineer and a Corporate Member of the Institution of Engineers, Sri Lanka, the Institution of Engineering and Technology, United Kingdom and the International Association of Engineers. Mohamed Firdhous has several years of industry, academic and research experience in Sri Lanka, Singapore and the United States of America.

APPENDIX

TABLE A-1: LIST OF SPEECH TAGS USED

| Tag | Description |
|---|---|
| CC | Coordinating conjunction |
| CD | Cardinal number |
| DT | Determiner |
| EX | Existential *there* |
| FW | Foreign word |
| IN | Preposition or subordinating conjunction |
| JJ | Adjective |
| JJR | Adjective, comparative |
| JJS | Adjective, superlative |
| LS | List item marker |
| MD | Modal |
| NN | Noun, singular or mass |
| NNS | Noun, plural |
| NNP | Proper noun, singular |
| NNPS | Proper noun, plural |
| PDT | Predeterminer |
| POS | Possessive ending |
| PRP | Personal pronoun |
| PRP$ | Possessive pronoun |
| RB | Adverb |
| RBR | Adverb, comparative |
| RBS | Adverb, superlative |
| RP | Particle |
| SYM | Symbol |
| TO | *to* |
| UH | Interjection |
| VB | Verb, base form |
| VBD | Verb, past tense |
| VBG | Verb, gerund or present participle |
| VBN | Verb, past participle |
| VBP | Verb, non-3rd person singular present |
| VBZ | Verb, 3rd person singular present |
| WDT | Wh-determiner |
| WP | Wh-pronoun |
| WP$ | Possessive wh-pronoun |
| WRB | Wh-adverb |

TABLE A-2: LIST OF CHUNK TAGS USED

| Prefix | Chunk Type |
|---|---|
| B = beginning noun phrase | NP = noun phrase |
| I = in noun phrase | VP = verb phrase |
| O = other | PP = prepositional phrase |
| | O = other |

# Quantization Table Estimation in JPEG Images

Salma Hamdy, Haytham El-Messiry, Mohamed Roushdy, Essam Kahlifa
Faculty of Computer and Information Sciences
Ain Shams University
Cairo, Egypt
{s.hamdy, hmessiry, mroushdy, esskhalifa}@cis.asu.edu.eg

*Abstract*— **Most digital image forgery detection techniques require the doubtful image to be uncompressed and in high quality. However, most image acquisition and editing tools use the JPEG standard for image compression. The histogram of Discrete Cosine Transform coefficients contains information on the compression parameters for JPEGs and previously compressed bitmaps. In this paper we present a straightforward method to estimate the quantization table from the peaks of the histogram of DCT coefficients. The estimated table is then used with two distortion measures to deem images as untouched or forged. Testing the procedure on a large set of images gave a reasonable average estimation accuracy of 80% that increases up to 88% with increasing quality factors. Forgery detection tests on four different types of tampering resulted in an average false negative rate of 7.95% and 4.35% for the two measures respectively.**

*Keywords: Digital image forensics; forgery detection; compression history; Quantization tables.*

## I. INTRODUCTION

Due to the nature of digital media and the advanced digital image processing techniques provided by image editing software, adversaries may now easily alter and repackage digital content forming an ever rising threat in the public domain. Hence, ensuring that media content is credible and has not been "retouched" is becoming an issue of eminent importance for both governmental security and commercial applications. As a result, research is being conducted for developing authentication methods and tamper detection techniques. Mainly, active authentication include digital watermarking and digital signatures, while passive methods tend to exploit inconsistencies that in the natural statistics of digital images occur as a result of manipulation.

JPEG images are the most widely used image format, particularly in digital cameras, due to its efficiency of compression and may require special treatment in image forensics applications because of the effect of quantization and data loss. Usually JPEG compression introduces blocking artifacts and hence one of the standard approaches is to use inconsistencies in these blocking fingerprints as a reliable indicator of possible tampering [1]. These can also be used to determine what method of forgery was used. Many passive schemes have been developed based on these fingerprints to detect re-sampling [2] and copy-paste [3,4]. Other methods try to identify bitmap compression history using Maximum Likelihood Estimation (MLE) [5,6], or by modeling the distribution of quantized DCT coefficients, like the use of Benford's law [7], or modeling acquisition devices [8]. Image

acquisition devices (cameras, scanners, medical imaging devices) are configured differently in order to balance compression and quality. As described in [9,10], these differences can be used to identify the source camera model of an image. Moreover, Farid [11] describes JPEG *ghosts* as an approach to detect parts of an image that were compressed at lower qualities than the rest of the image and uses to detect composites.

In this paper we present a straightforward method for estimating the quantization table of single JPEG compressed images and bitmaps. We verify the observation that while ignoring error terms, the maximum peak of the approximated histogram of a DCT coefficient matches the quantization step for that coefficient. This can help in determining compression history, i.e. if the bitmap was previously compressed and the quantization table that was used, which is particularly useful in applications like image authentication, artifact removal, and recompression with less distortion.

After estimating the quantization table, both average distortion measure and blocking artifact measure are calculated based on the estimated table to verify the authenticity of the image.

All simulations were done on images from the UCID [12]. Performance for estimating $Q$ for single JEPG images was tested against two techniques that are relevant in how the quantization steps are acquired; MLE [5,6], and power spectrum [1]. For the other abovementioned techniques (e.g. Benford's), they are said to work on bitmaps. Investigating performance for previously compressed bitmaps can be found in [13]. The rest of the paper is organized as follows. In section 2 we begin with a brief review of the JPEG baseline procedure and then show how the quantization steps can be determined from the peaks of the approximated histogram of DCT coefficients. We also present the two distortion measure used in evaluation. Testing and performance evaluation are discussed and section 3, where we demonstrate the use of estimated quantization table with the distortion measures in classifying test images and exposing forged parts. Finally, section 4 is for conclusions.

## I. A STRAIGHTFORWARD APPROACH FOR QUANIZATION TABLE ESTIMATION IN JPEG IMAGES

The JPEG standard *baseline* compression for color

(a)                                        (b)

**Fig. 1.** Histogram of (a) $|X_q(3,3)|$ formed as periodic spaced peaks and (b) $|X^*(3,3)|$ formed as periodic spaced sets of peaks. The DCT coefficients were quantized with step size $Q(3,3)$=10 during compression.

photographs consists of four lossy transform steps and yields a compressed stream of data:

(1) RGB to $YC_bC_r$ color space conversion.

(2) $C_bC_r$ subsampling.

(3) Discrete Cosine Transform (DCT) of 8×8 pixel blocks.

(4) Quantization: $X(i, j) = round(D(i, j)/Q(i, j))$ , where at frequency *(i,j)*, D is the DCT coefficient, Q is the *(i,j)*<sup>th</sup> entry in the quantization table, and X is the resulting quantized coefficient.

Equivalently, a decompression process involves

(1) Dequantization: $X_q(i, j) = X(i, j)Q(i, j)$ .

(2) Inverse Discrete Cosine Transform (IDCT).

(3) $C_bC_r$ interpolation.

(4) $YC_bC_r$ to RGB color space conversion.

One of the most useful aspects in characterizing the behavior of JPEG compressed images is the histogram of DCT coefficients which typically has a Gaussian distribution for the DC component and a Laplacian distribution for the AC components [5,6]. The quantized coefficients are recovered, in step (1) of the dequantizer above, as multiples of *Q(i,j)*. Specifically, if $X_q(i,j)$ is a dequantized coefficient in the DCT domain, it can be expressed as *kQ(i,j)*, where Q*(i,j)* is the *(i,j)*<sup>th</sup> entry of the quantization table, and *k* is an integer. The estimation of *Q(i,j)* is direct from the histogram of $X_q(i,j)$ but $X_q(i,j)$ is an intermediate result and is discarded after decompression. Theoretically, $X_q(i,j)$ can be recalculated as $DCT(X_q(i,j))$ since IDCT is reversible. Nevertheless in reality, the DCT of an image block usually generates $X^*(i,j)$, which is not exactly $X_q(i,j)$, but an approximation of it. In our experiments, we show that *Q(i,j)* can also be directly determined from histogram of $X^*(i,j)$. **Fig. 1(a)** and **(b)** show a typical absolute discrete histogram of *X(3,3)* and *X\*(3,3)* respectively, for all blocks of an image. Nonzero entries occur mainly at multiples of *Q(3,3)*=10.

There are two main sources of error introduced during the IDCT calculation, mainly rounding and clipping, to keep the pixel levels integral and within the same range as a typical 8-

bit image (0-255). The decaying envelopes of the histograms in **Fig. 1** are roughly Gaussian although have shorter tails, at which the approximation error $|X^*(i, j) - X_q(i, j)|$ is limited. The reason according to [6] is that as the rounding error for each pixel does not exceed 0.5, the total rounding error is



(a)



(b)

**Fig. 2.** (a) test image compressed with QF = 80 and (b) its corresponding quantization table.

bounded by

$$\Gamma = \left|X^*(i,j) - X_q(i,j)\right| \le B(i,j)$$
$$= \sum_{u,v} 0.5 \, c(u) \, c(v) \left|\cos \frac{(2u+1)i\pi}{16} . \cos \frac{(2v+1)j\pi}{16}\right| \quad (1)$$

where $c(\omega) = \begin{cases} 1/\sqrt{2} & for \; \omega = 0 \\ 1 & otherwise \end{cases}$ .

So, $\Gamma$ can be modeled as a truncated Gaussian distribution in the range ±*B* and zero outside that range [6].

Now if we closely observe the histogram of $X^*(i,j)$ outside the main lobe (zero and its proximity), we notice that the maximum peak occurs at a value that is equal to the quantization step used to quantize $X_q(i,j)$. This means that rounding errors has less significance and could be ignored in the estimation. On the other hand, clipping or truncation errors are more significant and cannot be compensated for. Hence in our experiments, we leave out saturated blocks when creating the histogram. **Fig. 2** shows a test image compressed with quality factor 80, and the corresponding quantization table. **Fig. 3(a)** and **(b)** show *H* the absolute histograms of DCT coefficients of the image from **Fig. 2(a)** at frequencies (3,3) and (3,4), respectively. Notice that the maximum peak for (3,3) occurs at 6 which is equal to Q(3,3). Also for (3,4), the highest peak is at value 10, which corresponds to the (3,4)<sup>th</sup> entry of the quantization table. Because in JPEG compression, the brightness (the DC coefficient) or shading across the tile (the 3 lowest AC coefficients) must be reproduced fairly accurately, there is enough information in the histogram data to retain *Q(i,j)* for low frequencies. We have verified that the highest peak outside the main lobe corresponds to *q*, for all low frequency coefficients.

And since the texture (middle frequency AC coefficients) can be represented less accurately, the histogram of these frequencies may not be a suitable candidate for our observation. Indeed, we found that, for coefficients highlighted in gray in **Fig. 2(b)**, the maximum peak occurs at a value that does not match the specific quantization step. However, we investigated peaks at parts of *H* where error is minimal, i.e. outside ±*B* and concluded that the observation still applies with condition. The maximum peak above *B*, (that is when $| X^*(i,j)|$>B) occurred at a value matching the $Q(i,j)$, **Fig. 3(c)** and (**d**). For a particular frequency *(i,j)*, it is possible

that no peaks are detected outside the main lobe. This occurs with heavy compression when the quantization step used is

TABLE I.         DIFFERENCE BETWEEN ESTIMATED AND ORIGINAL *Q*.

| QF = 75 | | | | | | | | | QF = 80 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | | 8 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | X | | 0 | 0 | 0 | 0 | 0 | 0 | 0 | X |
| 0 | 0 | 0 | 0 | 0 | 0 | X | X | | 0 | 0 | 0 | 0 | 0 | 0 | X | X |
| 0 | 0 | 0 | 0 | X | X | X | X | | 0 | 0 | 0 | 0 | X | X | X | X |

large and hence $X^*(i,j)$ becomes small and sometimes quantized to zeros for all blocks. The histogram decays rapidly to zero showing no periodic structure. Hence we do not have enough information to determine $Q(i,j)$. **Table 1** shows the difference between estimated *Q* table using the above method, and the original table for two quality factors. The X's mark the "undetermined" coefficients.

The next step is to use the estimated table to verify the authenticity of the image by computing a distortion measure and then comparing it to a preset threshold. One measure is the average distortion measure. This is calculated as a function of the remainders of DCT coefficients with respect to the original Q matrix:

$$B_1 = \sum_{i=1}^{8} \sum_{j=1}^{8} \mathrm{mod}\big(D(i,j), Q(i,j)\big) \qquad (2)$$



(a)                           (b)

(c)                           (d)

**Fig. 3.** Absolute histogram of (a) $X^*(3,3)$ where $H_{max}$ occurs at Q(3,3)=6. (b) $X^*(3,4)$ where $H_{max}$ occurs at Q(3,4) = 10 (c) $X^*(5,4)$ where $H_{max}$ occurs at Q(5,4)=22. (d) $X^*(7,5)$ where $H_{max}$ occurs at Q(7,5) = 41.

where $D(i,j)$ and $Q(i,j)$ are the DCT coefficient and the corresponding quantization table entry at position *(i,j)*. An image block having a large average distortion value indicates that it is very different from what it should be and is likely to belong to a forged image. Averaged over the entire image, this measure can be used for making a decision about authenticity of the image.

Another measure is the blocking artifact measure, BAM [1], which is caused by the nature of the JPEG compression method. The blocking artifacts of an image block will change

a lot by tampering and therefore, inconsistencies in blocking artifacts serve as evidence that the image has been "touched". It is computed from the *Q* table as:

$$B_2(n) = \sum_{i=1}^{8} \sum_{j=1}^{8} \left| D(i,j) - Q(i,j) \, round\left( \frac{D(i,j)}{Q(i,j)} \right) \right| \qquad (3)$$

*B(n)* is the estimated blocking artifact for testing block *n*, *D(i,j)* and *Q(i,j)* are the same as in (2).

## II.    EXPERIMENTAL RESUTLS AND DISCUSSION

### A. Estimation Accuracy

We created a dataset of image to serve as our test data. The set consisted of 550 uncompressed images collected from different sources (more than five camera models), in addition to some from the public domain Uncompressed Color Image Database (UCID), which provides a benchmark for image processing analysis [12]. For color images, only the luminance plane is investigated at this stage. Each of these images was compressed with different standard quality factors, [50, 55, 60,

65, 70, 75, 80, 85, and 90]. This yielded 550×9 = 4,950 *untouched* images. For each quality factor group, an image's histogram of DCT coefficients at one certain frequency was generated and used to determine the corresponding quantization step at that frequency according to section 2. This was repeated for all the 64 histograms of DCT coefficients. The resulting quantization table was compared to the image's known table and the percentage of correctly estimated coefficients was recorded. Also, the estimated table was used in equations (2) and (3) to determine the image's average distortion and blocking artifact measures, respectively. These values were recorded and used later to set a threshold value for distinguishing forgeries from untouched images.

The above procedure was applied to all images in the dataset. **Table 2** shows the accuracy of the used method for each tested quality factor averaged over the whole set of images. It shows that quality factor of 75 gives a percentage of around 80%. This is reasonable as this average quality factor yields the best image quality-compression tradeoff and hence the histograms have enough data to accurately define the quantization steps. As the quality factor decreases, estimation accuracy drops steadily. This, as explained earlier, is due to heavy quantization and corresponding large steps used with lower qualities. Histograms convey no data to predict the compression values. For higher quality factors, it is predictable that performance tend to improve which is apparent in the rising values in **Table 2**. Nevertheless, notice the drop in estimation for very high quality factors (95 and 100). This is

If most low frequency steps are 1 then we consider *QF* = 100 and output the corresponding table of 64 ones.

To verify that a wide range of quantization tables, standard and non standard can be estimated, we created another image set of 100 JPEG images from different sources as our arbitrary test set. Each image's quantization table was estimated and the percentage of correctly estimated coefficients recorded. This gave an average percentage of correct estimation of 86.45%.

Maximum Likelihood methods for estimating *Q* tables [5-6], tend to search for all possible *Q(i,j)* for each DCT coefficient over the whole image which can be computationally exhaustive. Furthermore, they can only detect standard compression factors since they re-compress the image by a sequence of preset quality factors. This can also be a time consuming process. Other methods [1, 8] estimate the first few (often first 3×3) low frequency coefficients and then search through lookup tables for matching standard tables. Ye et. al [1], proposed a new quantization table estimation based on the power spectrum, PS, of the histogram of DCT coefficients. They constructed a low-pass filtered version of the second derivative of the PS and found that the number of local minima plus one equals the quantization step. Only the first 32 coefficients are used in the estimation because high frequency DCT coefficients would be all zero when quantized by large step. The authors of that work did not provide filter specifications, and we believe through experimenting, that there are no unanimous low-pass filter parameters for all quality factors or for all frequency bands. This means that

TABLE II.     PERCENTAGE OF CORRECTLY ESTIMATED COEFFICIENTS FOR SEVERLA QFS

| 50 | 55 | 60 | 65 | 70 | 75 | 80 | 85 | 90 | 95 | 100 |
|---|---|---|---|---|---|---|---|---|---|---|
| 66.9 | 69.2 | 72.0 | 74.2 | 76.9 | 79.4 | 82.3 | 85.5 | 88.2 | 66.33 | 52.71 |

TABLE III.     AVERAGE ESTIMATION ACCURACY AGAINST OTHER METHODS FOR DIFFERENT QUALITY FACTORS.

| QF Method | 50 | 60 | 70 | 80 | 90 |
|---|---|---|---|---|---|
| MLE | 59.12 | 63.75 | 86.25 | 86.34 | 70.50 |
| Power Spectrum | 65.37 | 68.84 | 75.75 | 90.12 | 84.75 |
| Maximum Peak | 96.04 | 97.69 | 97.33 | 91.89 | 73.33 |

   due to very small quantization steps. The peaks of the histogram are no longer distinguishable as "bumps" outside the zero vicinity, but rather show as quick swinging. Moreover, most of the lower steps for such high qualities have the values of 1 or 2, which are very close to zero (for *QF*=100, all entries of the *Q* table are 1's and no compression takes place). In our method, we remove zero and its neighborhood which are all the next lower points until we hit a mount again. These values removed before estimation causes our method to always fail to estimate a step size of 1. The only case we manage to record a 1 is when the histogram of 1 is larger than the histogram of 0 which sometimes occur within lower frequencies. As for higher frequencies, they often give erroneous results. One way to correct them is to threshold the number of entries in the resulting table having the value of 1.

either we use different settings for each group of *Q* steps, or use one filter to get a few low frequencies and then retrieve the rest of the table through matching in lookup tables. We found that a 1×3 Gaussian filter with a large cutoff frequency gave the best possible results when tested on a number of images. We used the filter to estimate the first nine AC coefficients and recorded the percentage of correct estimation. **Tables 3 and 4** show the estimation time and accuracy of the MLE method and power spectrum method against our method for different quality factors averaged over 500 test images of size 640×480 from the UCID. While MLE requires double the time, the average time in seconds for the latter two methods is very close while the average accuracy of the power spectrum method using the specified filter was around 77%. We believe filter choice is crucial but since we could not optimize a fixed set of parameters, we did not investigate the method any further.

*B.   Forfery Detection*

To create the image set used for forgery testing, we selected 500 images from the untouched image set. Each of these images was processed in a way and saved with different quality factors. More specifically, each image was subjected to four kinds of common forgeries; cropping, rotation, composition, and brightness changes. Cropping forgeries were

done by deleting some columns and rows from the original image to simulate cropping from the left, top, right, and bottom. For rotation forgeries, an image was rotated by 270°. Copy-paste forgeries were done by copying a block of pixels randomly from an arbitrary image and then placing it in the original image. Random values were added to every pixel of the image to simulate brightness change. The resulting fake images were then saved with the following quality factors [60, 70, 80, and 90]. Repeating this for all selected images produced total of (500×4) × 4 = 8,000 images. Next, the quantization table for each of these images was estimated as above and used to calculate the image's average distortion,

(2), and the blocking artifact, (3), measures, respectively.

The scattered dots in **Fig. 4** show the values of the average distortion for 500 untouched images (averaged for all quality factors for each image) while the cross marks show the average distortion values for 500 images from the forged dataset. As the figure shows, the values are distinguished to distortion measure and hence the values for forged images tend to cluster higher than those for untampered images.

with the average FNR over all images. As expected, as *QF* increases, a better estimate of the quantization matrix of the original untampered image is obtained, and as a result the error percentage decreases. Notice that cropping needs to destroy the normal JPEG grid alignment in order to achieve high distortion and hence mark the image as possible fake. This is because if the picture happens to be aligned perfectly to the original grid after cropping, then the cropping forgery would go undetected in this case. Similarly, detecting copy-paste forgery is possible since the pasted part fails to fit perfectly into the original JPEG compressed image. As a result, when the distortion metric is calculated, it exceeds the detection threshold. Charts show that the blocking artifact measure recorded a lower threshold and usually lower FNR than average distortion measure. Generally, the performance of the two measures is relatively close for brightened and rotated images. However, BAM is more sensitive to cropping and compositing since it works on the JPEG's "grid" and these two manipulations tend to destroy that natural grid. Brightness manipulated images are the most ones likely to go undetected

TABLE IV. AVERAGE ESTIMATION TIME (FIRST 3×3) AGAINST OTHER METHODS FOR DIFFERENT QUALITY FACTORS.

| QF Method | 50 | 60 | 70 | 80 | 90 |
|---|---|---|---|---|---|
| MLE | 22.29 | 22.35 | 22.31 | 22.26 | 22.21 |
| Power Spectrum | 11.37 | 11.26 | 10.82 | 10.82 | 11.27 |
| Maximum Peak | 11.27 | 11.29 | 11.30 | 11.30 | 11.30 |

as they leave the grid intact.



**Fig. 4.** Average distortion measure for untouched and tampered images.

Through practical experiments we tested the distortion measure for untouched images against several threshold values and calculated the corresponding false positive rate FPR (the number of untouched images deemed as forged.), i.e., the number of values above the threshold. Optimally, we aim for a threshold that gives nearly zero false positive. However, we had to take into account the false negatives (the number of tampered images deemed as untampered) that may occur when testing for forgeries. Hence, we require a threshold value keeping both FPR and the FNR low. But since we rather have an untampered show up as tampered, rather than the other way round, we chose a threshold that is biased towards false positive rate. We selected a vale that gave FPR of 12.6% and a lower FNR as possible for the different types of forgeries. The horizontal line marks the selected threshold $\tau = 30$. Similarly, the same set of images was used with the BAM and the threshold was selected to be $\tau = 20$, with a corresponding FPR of 6.8%.

**Fig. 5** shows the false negative rate (FNR) for the different forgeries at different quality factors. The solid line represents the FNR for the average distortion measure, while the dashed line is for the blocking artifact measure. Each line is labeled



**Fig. 5.** False negative rate for two distortion measures, calculated for different forgery types.

**Fig. 6(a)** and **(b)** show two untouched images that are used to make a composite image **(c)**. Part of the car from the second images was copied and pasted into the first image and the result was saved with different compression factors. The resulting distortion measures for the composite image are shown in **Fig. 6(d)** through **(g)**. The dark parts denote low distortion whereas brighter parts indicate high distortion values. Notice the highest values corresponding to the part

pasted from the second image and hence marking the forged area. Apparently as the quality factor increases, detection performance increases. Moreover, if the forged image was saved as a bitmap, detecting inconsistencies becomes easier as no quantization, hence loss of data, takes place. This can help in establishing bitmap compression history.

### III. CONCLUSIONS

We showed in this paper that while ignoring quantization rounding errors, we still can achieve reasonably high quantization table estimation accuracy through computing a histogram once for each DCT coefficient. The maximum peak method, although straightforward gives good estimation results while neglecting rounding error. Hence, this reduces the need to statistically model rounding errors and hence reduces computations and time. It was tested against MLE method that models round off errors as modified Gaussian, and proved to require half the time with no degraded accuracy, if not better for some quality factors.

We have found through extensive test that the method estimates all low frequencies in addition to a good percentage of the high frequencies. Hence, this reduces the need for lookup tables and matching overtime as a large percentage of the table can be reliably estimated directly from the histogram (even some high frequencies). And by "large percentage" we mean enough entries to compute the distortion measure correctly without further searching in lookup tables. Also this means that arbitrary step sizes can be estimated which are often used in different brands of digital cameras.

The method was tested against the power spectrum method and proved to require nearly the same estimation time with improved accuracy. However, eliminating the need for lookup tables will naturally affect execution time since we will have to process all 64 entries not just the first 9.

Nevertheless, for images heavily compressed, the histogram fails to estimate high frequencies. In this case, we can always estimate the first few low frequency coefficients and then search lookup tables for a matching $Q$ table. Of course this works only for standard compression table. Also images with large homogenous areas may fail to give estimation if we choose to exclude uniform blocks when approximating the histogram. In addition, performance tends to drop when an image is further compressed with a different quality factor. In such cases, double quantization leaves its traces in the histogram and methods for estimating primary and secondary quantization tables can be used.

Maximum peak also works well for retrieving bitmaps previous compression tables and using them for forgery detection.

### FUTURE WORK

Investigating the chroma planes and further testing on bitmaps and multiple compressions is due as future work. Also, after classifying an image, we require and approach that can be used to identify which type of manipulations the image underwent.

*Estimation of color space quantization tables:* We have so far addressed gray scale images and the luminance channel of color images. A further study of the two chroma channels and the histograms of their DCT coefficients, and hence suggestion of possible methods for estimating the chroma tables, are natural extension to this work.

*Double Quantization*: Double compressed images contain specific artifacts that can be employed to distinguish them from single compressed images. When creating composites, the pasted portion will likely exhibit traces of a single compression, while the rest of the image will exhibit signs of double compression. This observation could in principle be used to identify manipulated areas in digital images.

*JPEG2000:* provides better compression rates with respect to quality compared to the standard JPEG compression. It is based on Wavelet transforms, and constitutes an interesting research topic in digital image forensics.

### REFERENCES

[1] Ye S., Sun Q., Chang E.-C., "Detection Digital Image Forgeries by Measuring Inconsistencies in Blocking Artifacts", in *Proc. IEEE Int. Conf. Multimed. and Expo.*, July, 2007, pp. 12-15.

[2] Popescu A., Farid H., "Exposing Digital Forgeries by Detecting Traces of Resampling", *IEEE Trans. Signal Process*, 53(2): 758–767, 2005.

[3] Fridrich J., Soukal D., Lukas J., "Detection of Copy-Move Forgery in Digital Images", *Proc. Digit. Forensic Res. Workshop*, August 2003.

[4] Ng T.-T., Chang S.-F., Sun Q., "Blind Detection of Photomontage Using Higher Order Statistics," in *Proc. IEEE Int. Symp. Circuits and Syst*, vol. 5, May, 2004, pp. 688-691.

[5] Fan Z., de Queiroz R. L., "Maximum Likelihood Estimation of JPEG Quantization Table in The Identification of Bitmap Compression History", in *Proc. Int. Conf. Image Process. '00*, 10-13 Sept. 2000, 1: 948–951.

[6] Fan Z., de Queiroz R. L., "Identification of Bitmap Compression History: JPEG Detection and Quantizer Estimation", in *IEEE Trans. Image Process.*, 12(2): 230–235, February 2003.

[7] Fu D., Shi Y.Q., Su W., "A Generalized Benford's Law for JPEG Coefficients and its Applications in Image Forensics", *in Proc. SPIE Secur., Steganography, and Watermarking of Multimed. Contents IX*, vol. 6505, pp. 1L1-1L11, 2007.

[8] Swaminathan A., Wu M., Ray Liu K. J., "Digital Image Forensics via Intrinsic Fingerprints", *IEEE Trans. Inf. Forensics Secur.*, 3(1): 101-117, March 2008.

[9] Farid H., "Digital Image Ballistics from JPEG Quantization," *Department of Computer Science, Dartmouth College*, *Technical. Repo*rt TR2006-583, 2006.

[10] Farid H., "Digital Ballistics from JPEG Quantization: A Follow-up Study," *Department of Computer Science, Dartmouth College, Technical. Repo*rt TR2008-638, 2008.

[11] Farid H., "Exposing Digital Forgeries from JPEG Ghosts," in *IEEE Trans. Inf. Forensics Secur.,* 4(1): 154-160, 2009.

[12] Schaefer G., Stich M., "UCID – An Uncompressed Color Image Database", *School of Computing and Mathematics*, *Technical. Report,* Nottingham Trent University, U.K., 2003

[13] Hamdy S., El-Messiry H., Roushdy M. I., Kahlifa M. E., "Retrieval of Bitmap Compression History", *IJCSIS Vol. 8 No. 8,* November, 2010.

(a) Original with QF = 80.



(b) Original with QF = 70.



(c) Composite image.



91.6285



35.5106

(d) QF = 70



60.7237



27.4198

(e) QF = 80



57.1557



28.023

(f) QF = 90



81.0483



50.9392

(g) BMP

**Fig. 6.** Two test images (a) and (b) used to produce a composite image (c). For each QF (d) through (g), the left column figures represents the average distortion measure while the right column figures represents the blocking artifact measure for the image in (c).

# Modified ID-Based Public key Cryptosystem using Double Discrete Logarithm Problem

Chandrashekhar Meshram
Department of Applied Mathematics,
Shri Shankaracharya Engineering College
Junwani, Bhilai (C.G) India
Email: cs_meshram@rediffmail.com

*Abstract*— **In 1984, Shamir [1] introduced the concept of an identity-based cryptosystem. In this system, each user needs to visit a key authentication center (KAC) and identify him self before joining a communication network. Once a user is accepted, the KAC will provide him with a secret key. In this way, if a user wants to communicate with others, he only needs to know the "identity" of his communication partner and the public key of the KAC. There is no public file required in this system. However, Shamir did not succeed in constructing an identity based cryptosystem, but only in constructing an identity-based signature scheme. Meshram and Agrawal [4] have proposed an id - based cryptosystem based on double discrete logarithm problem which uses the public key cryptosystem based on double discrete logarithm problem. In this paper, we propose the modification in an id based cryptosystem based on the double discrete logarithm problem and we consider the security against a conspiracy of some entities in the proposed system and show the possibility of establishing a more secure system.**

*Keywords- Public key Cryptosystem, Identity based Cryptosystem, Discrete Logarithm Problem, Double Discrete Logarithm Problem.*

## I. INTRODUCTION

In a network environment, secret session key needs to be shared between two users to establish a secret communication. While the number of users in the network is increasing, key distribution will become a serious problem. In 1976, Diffie and Hellman [6] introduced the concept of the public key distribution system (PKDS). In the PKDS, each user needs to select a secret key and compute a corresponding public key stored in the public directory. The common secret session key, which will be shared between two users can then be determined by either user, based on his own secret key and the partner's public key. Although the PKDS provides an elegant way to solve the key distribution problem, the major concern is the authentication of the public keys used in the cryptographic algorithm.

Many attempts have been made to deal with the public key authentication issue. Kohnfelder [7] used the RSA digital signature scheme to provide public key certification. His system involves two kinds of public key cryptography: one is in modular p, where p is a large prime number; the other is in modular n, where n = p q, and p and q are large primes. Blom [11] proposed a symmetric key generation system (SKGS based on secret sharing schemes. The problems of SKGS

however, are the difficulty of choosing a suitable threshold value and the requirement of large memory space for storing the secret shadow of each user.

In 1984, Shamir [1] introduced the concept of an identity-In this system; each user needs to visit a based cryptosystem. Key authentication center (KAC) and identify him self before joining the network. Once a user is accepted, the KAC will provide him with a secret key. In this way, a user needs only to know the "identity" of his communication partner and the public key of the KAC, together with his secret key, to communicate with others. There is no public file required in this system. However, Shamir did not succeed in constructing an identity-based cryptosystem, but only in constructing an identity-based signature scheme. Since then, much research has been devoted, especially in Japan, to various kinds of ID-based cryptographic schemes. Okamoto et al. [10] proposed an identity-based key distribution system in 1988, and later, Ohta [12] extended their scheme for user identification. These schemes use the RSA public key cryptosystem [18] for operations in modular n, where n is a product of two large primes, and the security of these schemes is based on the computational difficulty of factoring this large composite number n. Tsujii and Itoh [2] have proposed an ID- based cryptosystem based on the discrete logarithm problem with single discrete exponent which uses the ElGamal public key cryptosystem. Meshram and Agrawal [5] have proposed an ID- based cryptosystem based on the integer factoring and double discrete logarithm problem which uses the public key cryptosystem based on integer factoring and double discrete logarithm problem. Meshram and Agrawal [4] have also proposed an ID- based cryptosystem based on double discrete logarithm problem which uses the public key cryptosystem based on double discrete logarithm problem. Now we Modified this cryptosystem for discrete logarithm problem with distinct double discrete exponent because we face the problem of solving double and triple distinct discrete logarithm problem at the same time in the multiplicative group of finite fields as compared to the other public key cryptosystem where we face the difficulty of solving the traditional discrete logarithm problem in the common group.

In this paper , we present modification in an ID based cryptosystem based on the double discrete logarithm problem with distinct discrete exponent (the basic idea of the proposed

system comes on the public key cryptosystem based on double discrete logarithm problem) here we describe further considerations such as the security of the system, the identification for senders. etc. our scheme does not require any interactive preliminary communications in each message transmission and any assumption except the intractability of the discrete logarithm problem.(this assumption seems to be quite reasonable)thus the proposed scheme is a concrete example of an ID –based cryptosystem which satisfies Shamir's original concept [1] in a strict sense.

## II.    MODIFIED ID-BASED PUBLIC KEY CRYPTOSYSTEM

### A. Implementation of the ID –Based Cryptosystem Preparation for the center and each entity

*Step 1.* Each entity generates a k-dimensional binary vector for his $ID$ . We denote entity A's $ID$ by $ID_A$ as follows $ID_A = (x_{A1}, x_{A2}, .........., x_{Ak}), x_{Aj} \in \{0,1\}$ , $(1 \le j \le k)$ (1)

Each entity registers his $ID$ with the center, and the center stores it in a public file.

*Step 2.:* The center generate two random prime number $p$ and $q$ and compute

$$N = pq \qquad (2)$$

Then the center chooses an arbitrary random number $e, 1 \le e \le \varphi(N)$ , such that $\gcd(e, \varphi(N)) = 1$ where $\varphi(N) = (p-1)(q-1)$ is the Euler function of $N$ then center publishes $(e, N)$ as the public key. Any entity can compute the entity $A's$ extended $ID, EID_A$ by the following:

$$EID_A \equiv (ID)^e (\bmod N)$$
$$= (y_{A1}, y_{A2}, .........., y_{At}), x_{Aj} \in \{0,1\}$$
$(1 \le j \le t)$ (3)

where $t = |N|$ is the numbers of bits of $N$ .

*Step 3. Center's secrete information:* - The center chooses an arbitrary large prime number $p$ and $q$ and compute $N = pq$ and also generated n-dimensional vector $a$ and m-dimensional vector $b$ over $Z^*_{\varphi(N)}$ which satisfies

$$a = (a_1, a_2, .........., a_n), b = (b_1, b_2, .........., b_m) \qquad (4)$$
$$2 \le a_i b_l \le \varphi(N) - 1 , (1 \le i \le n), (1 \le l \le m), (m \le n)$$
$$abI \neq abJ(\bmod(p-1)), I \neq J \qquad (5)$$

Where $I$ and $J$ are n-dimensional binary vector and stores it as the centers secret information. The condition of equation (5) is necessary to avoid the accidental coincidence of some

entities secrete key. A simple ways to generate the vectors $a$ and $b$ is to use Merkle and Hellmans scheme [19].

*Step 4:* The center also chooses $w$ which satisfies $\gcd(w, \varphi(N)) = 1$ and $w < \lfloor \varphi(N)/n \rfloor$ , where $\lfloor x \rfloor$ also denote the floor function which implies the largest integer smaller than compute $x$ .

The center chooses a super increasing sequences corresponding to $a$ and $b$ as $a'_i (1 \le i \le n)$ & $b'_l (1 \le l \le m)$ satisfies

$$\sum_{j=1}^{i-1, l-1} a'_j b'_j + v \prec \varphi(N) \text{ where } v = \lfloor \varphi(N)/w \rfloor \qquad (6)$$

$$\sum_{j=1}^{n} a'_j b'_j \prec \varphi(N), (m \le n) \qquad (7)$$

Then the centre computes

$$a_i b_l = a'_i b'_l w (\bmod \varphi(N))$$
$$c_i = a_i b_l (\bmod w)(1 \le i \le n)(1 \le l \le m)(m \le n) \qquad (8)$$

Where

$$a = (a_1, a_2, .........., a_n), b = (b_1, b_2, .........., b_m) \qquad (9)$$

Remark 1: it is clear that the vector $a$ and $b$ defined by (9) satisfies (4)-(5) the above scheme is one method of generating an $n$ and $m$ dimensional vectors $a$ and $b$ satisfies (4)-(5). In this paper, we adopt the above scheme. However, another method might be possible.

*Step 5*: The center also chooses an arbitrary integer $t$ such that $e = (e_1, e_2, .........., e_t)$ , satisfying $\gcd(e_i, \varphi(N)) = 1$, $(1 \le i \le t)$ and compute n-dimensional and m- dimensional vectors $D^j$ and $D^k$ respectively:

$$D^j = (d^j_1, d^j_2, ..... d^j_n)(1 \le j \le n)$$
$$d^j_l = e_l a_l (\bmod \varphi(N))(1 \le l \le n) \qquad (10)$$
$$D^k = (d^k_1, d^k_2, ..... d^k_m)(1 \le k \le m)$$
$$d^k_l = e_l b_l (\bmod \varphi(N))(1 \le l \le m)(m \le n) \qquad (11)$$

**Since** $D^j$ and $D^k$ are one to one system.

*Step 5 Center public information:* The center chooses two arbitrary generators $\alpha$ and $\beta$ of $Z^*_{\varphi(N)}$ and computes n-dimensional vector $h$ using generator $\alpha$ & m-dimensional

vector $g$ using generator $\beta$ corresponding to the vector $a$ and $b$ .

$$h = (h_1, h_2, \ldots\ldots, h_n), g = (g_1, g_2, \ldots\ldots, g_m) \quad (12)$$

$$h_i = \alpha^{a_i} (\mathrm{mod}\, N), (1 \le i \le n) ,$$

$$g_l = \beta^{b_l} (\mathrm{mod}\, N), (1 \le l \le m) \quad (13)$$

The center informs each entity $(N, \alpha, \beta, h, g)$ as public information.

*Step 6. Each entity secrete key:* Entity $A's$ secrete keys $s_a$ and $s_b$ are given by inner product of $a$ and $b$ (the centre's secret information) and $EID_A$ (entity $A's$ extended $ID$ , see eqn.3)

$$s_a \equiv d_l^{j} EID_A (\mathrm{mod}\, \phi(N))$$
$$= \sum_{1 \le j \le n} d_l^{j} y_{Aj} (\mathrm{mod}\, \phi(N))$$

(14)

$$s_b \equiv d_l^{k} EID_A (\mathrm{mod}\, \phi(N))$$
$$= \sum_{1 \le j \le n} d_l^{k} y_{Aj} (\mathrm{mod}\, \phi(N)) \quad (15)$$

### B. System Initialization Parameters
*Center Secrete information*

$a$ : n -dimensional vector and $b$ m-dimensional vector {see (8)-(9)}

*Center public information*

$h$ : n -dimensional vector & $g$ m-dimensional vector {see eqn.(12-13)} $p$ and $q$ :large prime numbers, $e$ : random integers , two generator $\alpha$ and $\beta$ of $Z^*_{\varphi(N)}$ .

Entity $A's$ secrete keys $s_a$ and $s_b$ = entity $A's$ public information = $ID_A$ ,k-dimensional vector.

### C. Protocol of the proposed cryptosystem

Without loss of generality suppose that entity B wishes to send message $M$ to entity A.

*Encryption*

Entity B generates $EID_A$ (Entity $A's$ extended ID, see eqn.3) from $ID_A$ . It then computes $\gamma_1$ and $\gamma_2$ from corresponding public information $h$ and $g$ and $EID_A$ .

$$\gamma_1 = \Big( \prod_{1 \le i \le n} h_i^{\,y_{Ai}} \Big)^{e_i} (\mathrm{mod}\, N)$$

$$= \Big( \prod_{1 \le i \le n} (\alpha^{a_i})^{\,y_{Ai}} \Big)^{e_i} (\mathrm{mod}\, N)$$

$$= \alpha^{\sum_{1 \le i \le n} e_i \alpha_i y_{Ai} (\mathrm{mod}\, \varphi(N))} (\mathrm{mod}\, N)$$

$$= \alpha^{\sum_{1 \le i \le n} d_i^{j} y_{Ai} (\mathrm{mod}\, \varphi(N))} (\mathrm{mod}\, N)$$

$$= \alpha^{s_a} (\mathrm{mod}\, N)$$

$$\gamma_2 = \Big( \prod_{1 \le l \le m} g_l^{\,y_{Al}} \Big)^{e_l} (\mathrm{mod}\, N)$$

$$= \Big( \prod_{1 \le l \le m} (\beta^{b_l})^{\,y_{Al}} \Big)^{e_l} (\mathrm{mod}\, N)$$

$$= \beta^{\sum_{1 \le l \le m} e_l \beta_l y_{Al} (\mathrm{mod}\, \varphi(N))} (\mathrm{mod}\, N)$$

$$= \beta^{\sum_{1 \le l \le m} d_l^{k} y_{Al} (\mathrm{mod}\, \varphi(N))} (\mathrm{mod}\, N)$$

$$= \beta^{s_b} (\mathrm{mod}\, N)$$

Entity B use $\gamma_1$ and $\gamma_2$ in Public key cryptosystem based on double discrete logarithm problem.

Let $M (1 \le M \le N)$ be entity B's message to be transmitted. Entity B select two random integer $u$ and $v$ such that $(2 \le uv \le \varphi(N) - 1)$ and computes

$$C_1 = \alpha^{u} (\mathrm{mod}\, N)$$

$$C_2 = \beta^{v} (\mathrm{mod}\, N)$$

$$E = M (\gamma_1)^{u} (\gamma_2)^{v} (\mathrm{mod}\, N)$$

$$= M (C_1^{S_a} C_2^{S_b}) (\mathrm{mod}\, N)$$

The cipher text is given by $C = (C_1, C_2, E)$ .

*Decryption*

To recover the plaintext $M$ from the cipher text

Entity A should do the following Compute

$$C_1^{\varphi(N) - s_a} (\mathrm{mod}\, N) = C_1^{-s_a} (\mathrm{mod}\, N)$$

And $C_2^{\varphi(N) - s_b} (\mathrm{mod}\, N) = C_2^{-s_b} (\mathrm{mod}\, N)$

Recover the plaintext $M = \big( C_1^{-s_a} C_2^{-s_b} E \big) (\mathrm{mod}\, N)$

### III. SECURITY ANALYSIS

The security of the proposed ID based cryptosystem is based on the intractability of the discrete logarithm problem. It is very difficult to give formal proofs for the security of a

cryptosystem, in the following; we analyze some possible attacks against the above schemes and show that the security of these attacks is based on the DLP assumption.

1.    An intruder should solve a discrete logarithm problem twice to obtain the private key given the public as following: In this encryption the public key is given by $\left(N, e, \alpha, \beta, \gamma_1, \gamma_2\right)$ and the corresponding secret key is given by $\left(s_a, s_b\right)$.

   To obtain the private key $\left(s_a\right)$ he should solve the DLP

$$s_a \equiv \log_\alpha \left(\alpha^{s_a}\right) (\bmod N)$$

   To obtain the private key $\left(s_b\right)$ he should solve the DLP

$$s_b \equiv \log_\beta \left(\beta^{s_b}\right) (\bmod N)$$

This information is equivalent to computing the discrete logarithm problem over multiplicative cyclic group $Z^*_{\varphi(N)}$ and corresponding secrete key $s_a$ and $s_b$ will never be revealed to the public.

2.   An attacker might try to impersonate user $A$ by developing some relation between $w$ and $w'$ since

$$\gamma_1 \equiv Y^{w s_a} (\bmod N) \text{ and } \gamma_1' \equiv Y^{w' s_a} (\bmod N) \text{ Similarly}$$

$$\gamma_2 \equiv Y^{w s_b} (\bmod N) \text{ and } \gamma_2' \equiv Y^{w' s_b} (\bmod N) \text{ by}$$

knowing $\gamma_1, \gamma_2, w, w'$ the intruder can derive $\gamma_1'$ and $\gamma_2'$ as $\gamma_1' = \gamma_1^{w^{-1} w'} (\bmod N)$ and $\gamma_2' = \gamma_2^{w^{-1} w'} (\bmod N)$ without knowing $s_a$ and $s_b$ however trying to obtain $w$ from $\alpha$ and $\beta$ is equivalent to compute the discrete logarithm problem.

## IV.   CONCLUSION

    In this paper present the modification in an ID-based cryptosystem based on double discrete logarithm problem with distinct discrete exponents in the multiplicative group of finite fields. The proposed scheme satisfies Shamir's original concepts in a strict sense, i.e. it does not require any interactive preliminary communications in each data transmission and has no assumption that tamper free modules are available. This kind of scheme definitely provides a new scheme with a longer and higher level of security than that based on a double discrete logarithm problem with distinct discrete exponents. The proposed scheme also requires minimal operations in encryption and decryption algorithms and thus makes it is very efficient. The present paper provides the special result from the security point of view, because we face the problem of solving double and triple distinct discrete logarithm problem at the same time in the multiplicative group of finite fields as compared to the other public key

cryptosystem, where we face the difficulty of solving the traditional discrete logarithm problem in the common groups.

## REFERENCES

[1]   A. Shamir "Identity-based cryptosystem and signature scheme," Advances in Cryptology: Proceedings of Crypto' (Lecture Notes in Computer Science 196). Berlin, West Germany: Springer-Verlag, vol. 84 pp. 47-53,1985.

[2]   S. Tsujii, and T. Itoh "An ID-Based Cryptosystem based on the Discrete Logarithm Problem"IEEE Jounral on selected areas in communications vol. 7 pp 467-473, 1989.

[3]   T. ElGmal "A Public Key Cryptosystem and a Signature Scheme Based on Discrete Logarithms", IEEE Trans. Inform. Theory, vol. 31, pp 469-472, 1995

[4]   C.S.Meshram and S.S.Agrawal "An ID-Based Public key Cryptosystem based on the Double Discrete Logarithm Problem" International Journal of Computer Science and Network Security, vol.10 (7) pp.8-13,2010.

[5]   C.S.Meshram and S.S.Agrawal "An ID-Based Public key Cryptosystem based on Integer Factoring and Double Discrete Logarithm Problem" Information Assurance and Security Letters, vol.1 pp.029-034,2010.

[6]   W. Diffie and M.E. Hellman, "New direction in Cryptography", IEEE Trans.Inform.Theory, vol. 22, pp 644-654,1976.

[7]   L. M. Kohnfelder, "A method for certification," Lab. Comput. Sci. Mass. Inst. Technol.. Cambridge, MA, May 1978.

[8]   S. Tsujii, T. Itoh, and K. Kurosawa, "ID-based cryptosystem using discrete logarithm problem," Electron. Lett., vol. 23. no. 24, pp 1318-1320,1987.

[9]   S. C. Pohlig and M. E. Hellman, "An improved algorithm for com puting logarithms over GF (p) and its cryptographic significance," IEEE Trans. Inform. Theory, vol. IT-24, pp. 106-110,1978.

[10]   E. Okarnoto and K. Tanaka, "Key distribution system based on identification information," IEEE J. SeIecr. Areas Commun., vol. 7, pp.481485, May 1989.

[11]   R. Blorn, "An optimal class of symmetric key generation systems." In Proc. Eurocryp '84, Pans, France, Apr. 9-11, pp. 335-338,1984.

[12]   K. Ohta, "Efficient identification and signature schemes." Electron. Lett., vol. 24, no. 2, pp. 115-116,1988.

[13]   Wei-Bin Lee and Kuan-Chieh Liao "Constructing identity-based cryptosystems for discrete logarithm based cryptosystems" Journal of Network and Computer Applications,vol. 27, pp. 191–199,2004.

[14]   Min-Shiang Hwang, Jung-Wen Lo and Shu-Chen Lin "An efficient user identification scheme based on ID-based cryptosystem" Computer Standards & Interfaces,vol. 26,pp. 565–569,2004.

[15]   Eun-Kyung Ryu and Kee-Young Yoo "On the security of efficient user identification scheme" Applied Mathematics and Computation 2005, vol.171, pp. 1201–1205.

[16]   Mihir Bellare , Chanathip Namprempre and Gregory Neven "Security Proofs for Identity-Based Identification and Signature Schemes" J. Cryptol.,vol. 22, pp. 1–61, 2009.

[17]   S. C. Pohlig and M. E. Hellman, "An improved algorithm for computing logarithms over GF (p) and its cryptographic significance," IEEE Trans. Inform. Theory, vol. IT-24, pp. 106-110,1978.

[18]   R. L. Rivest, A. Shamir And L. Adelman, "A method for obtaining digital signatures and public-key cryptosystem," Comrnun. ACM., vol. 21, no. 2, pp. 120-126,1978.

[19]   R. C. Merkle and M. E. Hellman, "Hiding information and signatures in trapdoor knapsacks" IEEE Trans. Inform. Theory, vol. IT- 24, pp. 525-530,1978.

[20]   C.S.Laih and J.Y.Lee "Modified ID-Based Public key Cryptosystem using Discrete Logarithm Problem" Electronic Letters, vol.24 (14) pp.858-859,1988.

AUTHORS PROFILE

**Chandrashekhar Meshram** received the M.Sc and M.Phil degrees, from Pandit Ravishankar Shukla University, Raipur (C.G.), India in 2007 and 2008, respectively. Presently he is teaching as an Assistant Professor in Department of Applied Mathematics, Shri Shankaracharya Engineering College, Junwani, Bhilai, (C.G.) India. He is doing his research in the field of Cryptography and its Application. He is a member of International Association of Engineers, Hong Kong, Computer Science Teachers Association (CSTA)USA, Association for Computing Machinery (ACM) USA ,International Association of Computer Science and Information Technology (IACSIT), Singapore, European Association for Theoretical Computer Science (EATCS) Greece, International Association of Railway Operations Research (IAROR) NetNetherland, International Association for Pattern Recognition (IAPR) New York and International Federation for Information Processing (IFIP) Austria, International Mathematical Union (IMU) and Life -time member of Internet Society (ISOC) USA ,Indian Mathematical Society ,Cryptology Research Society of India and Ramanujan Mathematical Society of India (RMS).

# Efficient Implementation of Sample Rate Converter

Charanjit singh
University College of Engineering
Punjabi university
Patiala, India
charanjit@pbi.ac.in

Manjeet Singh patterh
University College of Engineering
Punjabi university
Patiala, India

Sanjay Sharma
ECED
Thapar University
Patiala, India

*Abstract*—**Within wireless base station system design, manufacturers continue to seek ways to add value and performance while increasing differentiation. Transmit/receive functionality has become an area of focus as designers attempt to address the need to move data from very high frequency sample rates to chip processing rates. Digital Up Converter (DUC) and Digital Down Converter (DDC) are used as sample rate converters. These are the important block in every digital communication system; hence there is a need for effective implementation of sample rate converter so that cost can be reduced. With the recent advances in FPGA technology, the more complex devices providing high-speed as required in DSP applications are available. The filter implementation in FPGA, utilizing the dedicated hardware resources can effectively achieve application-specific integrated circuit (ASIC)-like performance while reducing development time cost and risks. So in this paper the technique for an efficient design of DDC for reducing sample rate is being suggested which meets the specifications of WiMAX system. Its effective implementation also ensures the pathway for the efficient applications in VLSI designs. Different design configurations for the sample rate converter are explored. The sample rate converter can be designed using half band filters, fixed FIR filters, poly-phase filters, CIC filters or even farrow filters.**

*Keywords: WiMAX; Half Band filter; FIR filter; CIC Filter; Farrow Filter; FPGA.*

## I. INTRODUCTION

Sample rate conversion (SRC) is the process of changing the sampling rate of a data stream from a specific sampling rate (e.g. the input/output hardware rate) to another sampling rate (e.g. the application rate). With the conversion of communication and software markets, SRC is becoming a necessary component in many of today's applications including Digital Mixing Consoles and Digital Audio Workstations, CD-R, MD Recorders, Multitrack Digital Audio and Video Tape Recorders, Digital Audio Broadcast Equipment, Digital Tape Recorders, Computer Communication and Multimedia Systems. In most of these applications, a very high quality sample rate converter is required. Most high quality SRCs currently available on the market employ a digital filter that provides the required quality by up-sampling the data to a very high sampling rate followed by down-sampling to the required output sampling rate. The digital filters have also emerged as a strong option for removing noise, shaping spectrum, and minimizing inter-symbol interference in communication architectures. These filters have become popular because their precise reproducibility allows design engineers to achieve performance levels that are difficult to obtain with analog filters. The Cascaded Integrator Comb (CIC) filter is a digital filter which is employed for multiplier-less realization. This type of filter has extensive applications in low-cost implementation of interpolators and decimators. However, there is a problem of pass-band droop, which can be eliminated using compensation techniques. The Farrow filters is another class of digital filters which are used extensively in arbitrary sample rate conversions and fractionally delaying the samples. They have poly-phase structure and are very efficient for digital filtering. In addition to this, Field-Programmable gate Array (FPGA) has become an extremely cost-effective means of off-loading computationally intensive digital signal processing algorithms to improve overall system performance.

In this paper for Sample Rate Converter, CIC filter with and without compensation technique are implemented on FPGA. Farrow filters are also implemented for fractional delay and arbitrary change in sample rate conversion. Both of these filter configurations provide a better performance than the common filter structures in terms of speed of operation, cost, and power consumption in real-time. These filters are implemented in Altera Stratix-II-EP2S15F484C3 FPGA and simulated with the help of Quartus II v9.1sp2.

### A. Cascaded Integrator Comb (CIC) Filters

The CIC filter is a multiplier free filter that can handle large rate changes. It was proposed by Eugene Hogenauer in 1981 [1]. It is formed by integrating basic 1-bit integrators and 1-bit differentiators. It uses limited storage as it can be constructed using just adders and delay elements. That's why it is also well suited for FPGA and ASIC implementation. The CIC filter can also be implemented very efficiently in hardware due to its symmetric structure.

The CIC filter is a combination of digital integrator and digital differentiator stages, which can perform the operation of digital low pass filtering and decimation at the same time. The transfer function of the CIC filter in z-domain is given in equation (1) [1].

$$H(z) = \left( \frac{1 - z^{-K}}{1 - z^{-1}} \right)^{L} \tag{1}$$

In equation (1), K is the oversampling ratio and L is the order of the filter. The numerator $(1 - z^{-K})$ represents the transfer function of a differentiator and the denominator $(1/(1 - z^{-1})^{L}$ indicates the transfer function of an integrator.

Figure 1: CIC Filter [2]

The CIC filter first performs the averaging operation then follows it with the decimation. A simple block diagram of a first order CIC filter is shown in Figure 1. Here the clock divider circuit divides the oversampling clock signal by the oversampling ratio, M after the integrator stage. The same output can be achieved by having the decimation stage between integrator stage and comb stage. By dividing the clock frequency by M the delay buffer depth requirement of the comb section is reduced. In Figure 1, the integrator operates at the sampling clock frequency, while the differentiator operates at down sampled clock frequency of $f_s/$ M. By operating the differentiator at lower frequencies, a reduction in the power consumption is achieved.

A very poor magnitude characteristic of the comb filter is improved by cascading several identical comb filters. The transfer function H (z) of the multistage comb filter composed of K identical single stage comb filters is given by

$$H(z) = \left(\frac{1}{N}\frac{1 - z^{-N}}{1 - z^{-1}}\right)^K \qquad (2)$$

Figure 2 shows how the multistage realization improves the selectivity and the stop-band attenuation of the overall filter: the selectivity and the stop-band attenuation are augmented with the increase of the number of comb filter sections. The filter has multiple nulls with multiplicity equal to the number of the sections (K). Consequently, the stop-band attenuation in the null intervals is very high. Figure 3 illustrates a monotonic decrease of the magnitude response in the pass-band, called the pass-band droop.



Figure 2: CIC filter gain responses: single-stage K= 1, two-stage K=2, three-stage K=3 and four-stage K=4. [3]



Figure 3: CIC Filter Gain Response: Pass-band Droop for Figure 2. [3]

### B. CIC filter for sample rate conversion

The CIC filters are utilized in multirate systems for constructing digital Upconverter and Downconverter. The ability of comb filter to perform filtering without multiplications is very attractive to be applied to high rate signals. Moreover, CIC filters are convenient for large conversion factors since the low-pass bandwidth is very small. In multistage decimators with a large conversion factor, the comb filter is the best solution for the first decimation stage, whereas in interpolators, the comb filter is convenient for the last interpolation stage.

### C. CIC filters for Decimators

The basic concept of a CIC decimator is explained in Figure 4. Figure 4(a) shows the factor-of-N decimator consisting of the K-stage CIC filter and the factor-of-N down-sampler. Applying the third identity, the factor-of-N down-sampler is moved and placed behind the integrator section and before the comb section; see Figure 4(b). Finally, the CIC decimator is implemented as a cascade of K integrators, factor-of-N down-sampler, and the cascade of K differentiator (comb) sections. The integrator portion operates at the input data rate, whereas the differentiator (comb) portion operates at the N times lower sampling rate.

(a) Cascade of up-sampler and CIC filter.
(b) Cascade of comb section, down-sampler, and integrator section.
(c) Implementation structure consisting of the cascade of K differentiators, down sampler, and the cascade of K integrators. [3]

Figure 4: Block diagram representation of CIC decimator:
(a) Cascade of CIC filter and down-sampler.
(b) Cascade of integrator section, down-sampler, and comb section.
(c) Implementation structure consisting of the cascade of K integrators, down sampler, and the cascade of K differentiators. [3]

### D. CIC filters in Interpolators

The basic concept of a CIC interpolator is shown in Figure 5. Figure 5 (a) depicts the factor-of-N interpolator consisting of the factor-of-N up-sampler and the K-stage CIC filter. Applying the sixth identity, the factor-of-N up-sampler is moved and placed behind the differentiator (comb) section and before the integrator section; see Figure 5 (b). Finally, the CIC interpolator is implemented as a cascade of K differentiator (comb) sections, factor-of-N down-sampler, and the cascade of K integrators. The comb portion operates at the input data rate, whereas the integrator portion operates at the N times higher sampling rate.



Figure 5: Block diagram representation of CIC interpolator:

The configuration composed as a cascade of interpolators and differentiators (differentiators and interpolators) separated by a down-sampler (up-sampler) is called recursive realization structure, or a CIC realization structure. The advantage of CIC decimators and interpolators is the ability of sampling rate conversion without multiplying operations. This is of particular interest when operating at high frequencies. Considering the implementation aspects of CIC filters, one should expect the register overflow, since the integrator has a unity feedback. Actually, the register overflow is a reality in all integrator stages. It is shown that the register overflow is of no consequence if the two's complement arithmetic is used and the range of the number system is equal to or exceeds the maximum magnitude expected at the output of the composite filter. In order to avoid register overflow in the integrator section, the word-length has to be equal to or greater than

$(W_o + K \log_2 N)$ bits, where $W_o$ is the word-length in bits of the input signal.

### E. CIC filters in decimation and interpolation

A CIC filter can be used as a first stage in decimation when the overall conversion ratio M is factorable as

$$M = N \times R \qquad (3)$$

The overall factor-of-M sampling rate conversion system can be implemented by cascading a factor-of-N CIC decimator and a factor-of-R FIR decimator as shown in Figure 6(a). The corresponding single-stage equivalent is given in Figure 6(b).

When constructing an interpolator with a conversion factor L factorable as

$$L = R \times N \qquad (4)$$



Figure 6: Two-stage decimator composed of a CIC filter and an FIR filter:
(a) Cascade implementation.
(b) Single-stage equivalent

It might be beneficial to implement the last stage as a CIC interpolator. The first stage is usually implemented as an FIR filter. Figure 7(a) depicts the two-stage interpolator consisting of the cascade of a factor of R FIR interpolator and a factor of N CIC interpolator. The corresponding single-stage equivalent is shown in Figure 7(b)

In the two-stage solutions of Figures 6 and 7, the role of CIC decimator (interpolator) is to convert the sampling rate by

the large conversion factor N, whereas the FIR filter T(z) provides the desired transition band of the overall decimator (interpolator) and compensates the pass-band characteristic of the CIC filter.



Figure 7: Two-stage interpolator composed of an FIR filter in the first stage, and the CIC filter in the second stage:
(a) Cascade implementation.
(b) Single-stage equivalent [3]

Filter T $(Z^N)$ ensures the desired transition band, compensates the pass-band droop of the comb filter of the first stage. The CIC filter H(Z) has its two nulls just in the undesired pass-bands of the periodic filter T $(Z^N)$ that ensure the requested stop-band attenuation of the target two-stage decimator[3]. Finally, we compute the frequency response of the overall two-stage decimation filter,

$$H_1(z) = H(z).T(z^N) \qquad (5)$$



Figure 8: Gain responses of the CIC filter $H(z)$ (solid line), and that of periodic FIR filter T ($z^5$) (dashed line).



Figure 9: Gain responses of the two-stage decimator implemented as a factor of 5 and a factor of 2.

*F.   Filters with non integer Decimation factor: Farrow Filters*

When the decimation factor 1/R or the interpolation factor R is an integral value, then the conversion of sampling rate can be performed conveniently with the aid of fixed digital filters [8]. In case of a scenario where the factors are irrational, it will be impossible to use fixed digital filters directly. Moreover, if R is considered as the ratio of two relatively large prime integers, then, in the case of the conventional poly-phase implementation, it is quiet essential that the orders of the required filter become very large [8]. It means that a large number of coefficients need to be stored in coefficient memory. In sample rate conversion by non-integer factor, it is required to determine the values between existing samples. In this case, it is very convenient to use interpolation filters. Among them, polynomial-based filters are generally assumed to provide an efficient implementation form directly in digital domain. Such filters witness an effective implementation through Farrow structure or its higher version [8]-[13]. The main advantage of the Farrow structure is based on the presence of fixed finite-impulse response (FIR) filters as one of its ingredients. Thus eventually there is only one changeable parameter being the so-called fractional interval μ. Besides this, the control of μ is easier during the operation than in the corresponding coefficient memory implementations [6], and the concept of arithmetic preciseness; not the memory size limits the resolution of μ. These characteristics of the Farrow structure make it a very attractive structure to be implemented using a VLSI circuit or a signal processor [6].

Consider the diagram shown in Figure 10. The dashed line separates the filter into a section running at the input signal's sampling-rate and a section running at the output sampling rate [5]. Note that the output is re-labeled to be y[m] rather than y[n]. This is due to different input and output rates. Notably, the fractional delay now denoted $β_m$ will now change at every instant an output sample occurs.

Figure 10: Multi-rate Farrow Filter

Now consider a case where the sampling rate is increased by a factor of 2. Since for every input there are two outputs, the value held in the delay register will be used twice. The first time an input is used, $\beta_m$ will take on the value 0.5 and the output will be computed as

$$y[m] = 0.5(x[n-1] - x[n]) + x[n] = 0.5x[n-1] + 0.5$$

Before the input sample changes, one more output sample will be computed. $\beta_m$ will take the value 0 and the output will simply be

$$y[m+1] = x[n] \qquad (7)$$

Subsequently, the input sample will change; $\beta_m$ will be once again set to 0.5 and so forth.

In summary, when increasing the sampling rate by a factor of two, $\beta_m$ will cycle between the values {0.5, 0} twice as fast as the input, producing an output each time it changes.

In the general case, it is simply a matter of determining which values $\beta$ must take. The formula is simply

$$\beta_m = \left(\frac{mf_s}{f_s'}\right) \bmod 1 \qquad (8)$$

Where $f_s$ is the input sampling rate and $f_s^1$ is the output sampling rate.

In order to perform non integer SRC one can use Farrow structure or its modifications directly. However, in many cases, it becomes more efficient to use cascaded structures engineered by the modification of the Farrow structure and fixed FIR, or multistage FIR filter [12], [13]. The main advantage of using the cascaded structures instead of the direct modification of the Farrow structure lies in the fact that in case of joint optimization of the two building blocks the computational complexity to generate practically the same filtering performance is dramatically reduced. This is because of the following reasons. First, the implementation of a fixed linear phase FIR interpolator is not very costly, compared to the Farrow structure. Second, most importantly, the requirements for implementing the modification of the Farrow structure become significantly milder. This is mainly because of that the FIR filter takes care of pass-band and stop-band shaping, where the Farrow-based structure should only take care of attenuating images of FIR filter.

## II. DESIGN EXAMPLES

This section illustrates the properties of the proposed filters by means of design examples.

### A. Comparison of CIC filters with and without compensation:

The specifications for DDC for WiMAX are as follows:

- Input Sampling Frequency :91.392 MHz
- Output Sampling Frequency:11.424 MHz
- Pass-band Edge            :4.75 MHz
- Pass-band Ripple          :0.14 dB
- Stop-band Attenuation     :92 dB

*1)  Simulation results of CIC filter is as follows:*



Figure 11: Magnitude Response of DDC for WiMAX using CIC Filters

*2)  Simulation results of CIC filter with compensation is as follows:*



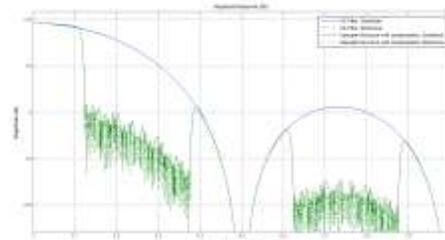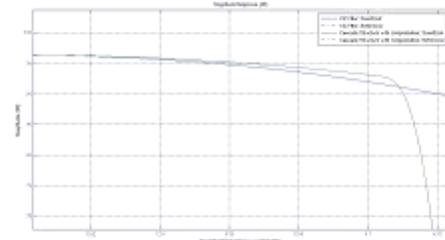Figure 12(a): Magnitude Response of DDC for WiMAX using CIC Filters using compensation



Figure 12(b): Magnitude Response of DDC for WiMAX using CIC Filters with Compensation

Both of the above filters are implemented in Altera's Stratix II FPGA family with device number EP2S15F484C3.

TABLE 1: Comparison of implementation cost and speed analysis of CIC filter and its cascaded structure with compensation:
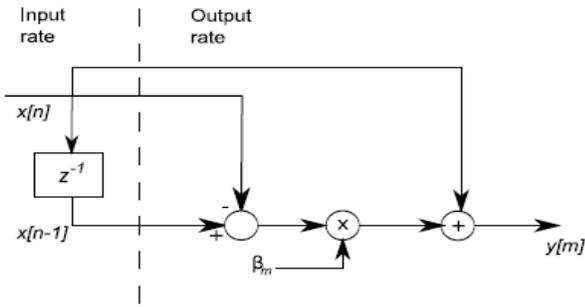
| PROPERTY | CIC Filter | Cascaded CIC Filter with Compensation |
|---|---|---|
| Logic Utilization | 29% | 63% |
| $t_{su}$ | 7.425 ns | 6.654 ns |
| $t_{cq}$ | 7.379 ns | 6.619 ns |
| Clock Frequency | 203.54 MHz | 21.91 MHz |

From the above simulation and implementation results, it can be concluded that the CIC filter are efficient for low-cost implementations. Due to absence of multipliers, they also have faster response. But the pass-band droop present in CIC filters restricts the scope of applications. With compensation technique, the response of CIC filter in pass-band is improved, but at the cost of extra hardware.

### B. Comparison of FIR and Farrow Sample Rate Conversions:

Consider a design example to change the sampling rate by a factor of 1.536 (192/125). The design to change the sample rate by an arbitrary factor is considered. The results obtained for the simulations are as follows:



Figure 13(a): Response of FIR SRC (blue) and Farrow SRC (green) for the factor of 1.536 changes in sample rates.



Figure 13(b): Response of FIR SRC (blue) and Farrow SRC (green)

From figure 13, it is clear that the response of farrow structure is better than the normal filter. Also, the efficiency of farrow filters is more as compared to other one. Both of these filters are also implemented efficiently in Altera's Stratix II FPGA family with device number EP2S15F484C3.

The MODELSIM simulation results of above designed filters are shown in figure 14-17.

### III. CONCLUSION

Sample rate conversion (SRC) is the process of changing the sampling rate of a data stream from a specific sampling rate (e.g. the input/output hardware rate) to another sampling rate (e.g. the application rate). Though the implementation of Signal Processing systems on ASICs provide better optimized devices, but the cost of such devices are rising. Also, the specification alteration requires the complete re-design of the system. With the recent advances in FPGA technology, the more complex devices providing high-speed as required in DSP applications are available. Also, the FPGA has advantage of reconfiguration which provides an upper hand over ASIC devices. The filter implementation in FPGA, utilizing the dedicated hardware resources can effectively achieve application-specific integrated circuit (ASIC)-like performance while reducing development time cost and risks.

In this paper, CIC filter and cascaded CIC filer for compensation is implemented in Altera's Stratix II FPGA for the specifications of Digital-Down Convertor of WiMAX. Though the CIC filter have upper hand on the basis of cost of implementation and speed, but the response of compensated CIC filter proves to be more reliable.

Apart from CIC filters, Farrow filters are also implemented on FPGA for SRC. Result shows that the Farrow filters are more efficient and have better response then FIR filters for arbitrary SRC Design.

### REFERENCES

[1] E. B. Hogenauer, "An economical class of digital filters for decimation and interpolation", IEEE Transactions on Acoustic, Speech and Signal Processing, ASSP29(2):155-162, 1981.

[2] U. Meyer-Bease, Digital Signal Processing with Field Programmable Gate Arrays, Springer, Third Edition, 2007.

[3] Ljiljana Milić, Multirate Filtering for Digital Signal Processing: MATLAB Applications, Hershey, PA: Information Science Reference, Jan 2009.

[4] R. E. Crochiere and L. R. Rabiner, Multirate Digital Signal Processing. Englewood Cliffs, NJ: Prentice-Hall, 1983.

[5] Ricardo Losada: Digital Filters with MATLAB, May 2008. [Online]. Available: www.mathworks.com/matlabcentral/fileexchange/19880.

[6] J. Vesma, Optimization and Applications of Polynomial Based Interpolation Filters, Doctoral Thesis, Tampere University of Technology, Publications 254, 1999.

[7] I. Løkken, The ups and downs of arbitrary sample rate conversion. [Online]. Available: www.iet.ntnu.no/courses/fe8114/slides/upsanddownsofasrc.pdf.

[8] D. Babic, A. Shahed Hagh Ghadam, M. Renfors, Polynomial-based filters with odd number of polynomialsegments for interpolation," to appear in IEEE Signal Processing Letters.

[9] C. Farrow, A continuously variable digital delay element, Proc. IEEE International Symposium on Circuits and Systems (ISCAS88), pages 2642-2645, 1998.

[10] T. Hentschel and G. Fettweis, Continuous-time digital filter for sample rate conversion in reconfigurable radio terminals, Proc. of European Wireless Conference, pages 55-59, 2000a.

[11] D. Babic and M. Renfors, Power efficient structures for conversion between arbitrary sampling rates, IEEE Signal Processing Letters, 12(1):1-4, 2005.

[12] D. Babic, T. Saramäki, and M. Renfors, "Sampling rate conversion between arbitrary sampling rates using polynomial-based interpolation filter," The Second International Workshop on Spectral Methods and Multirate Signal Processing SMMSP 2002, Toulouse, France, September 2002, pp. 57-64.

[13] D. Babic, J. Vesma, and M. Renfors, Decimation by irrational factor using CIC filter and linear interpolation, Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing, Vol. 6, May 2001, pp. 3677–3680.

[14] J. Vesma and T. Saramäki, Polynomial-Based Interpolation Filters—Part I: Filter Synthesis, Circuits Systems Signal Processing, Vol. 26, no. 2, 2007, pp. 115–146.

[15] W. Abu-Al-Saud and G. Stuber, Modified CIC filter for sample rate conversion in software radio systems, IEEE Signal Processing Letters, 10(5):152-154, 2003.

[16] T. Hentschel and G. Fettweis, Sample rate conversion for software radio, IEEE Communications Magazine, 38(8):142-150, 2000b.

[17] A. K. Z. Jiang and A. W. Jr., Application of filter sharpening to cascaded integrator-comb decimation filters, IEEE Transactions on Signal Processing, 45(2):457-467, 1997.

[18] G. Dolecek and S. Mitra, A new two-stage sharpened comb decimator, IEEE Transactions on Circuits and Systems-I, 52(7):1414-1420, 2005.

[19] H. Oh, S. Kim, G. Choi and Y. Lee, On the use of interpolated second order polynomials for efficient filter design in programmable down-conversion, IEEE Journal on Selected Areas in Communications, 17(4):551-560, 1999.

[20] G. Eason, B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," Phil. Trans. Roy. Soc. London, vol. A247, pp. 529–551, April 1955. *(references)*

[21] J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.

[22] I. S. Jacobs and C. P. Bean, "Fine particles, thin films and exchange anisotropy," in Magnetism, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.

[23] K. Elissa, "Title of paper if known," unpublished.

[24] R. Nicole, "Title of paper with only first word capitalized," J. Name Stand. Abbrev., in press.

[25] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, "Electron spectroscopy studies on magneto-optical media and plastic substrate interface," IEEE Transl. J. Magn. Japan, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetics Japan, p. 301, 1982].

[26] M. Young, The Technical Writer's Handbook. Mill Valley, CA: University Science, 1989.

# Cloud Computing Through Mobile-Learning

N.Mallikharjuna Rao
Associate Professor
Annamacharya P.G College of
Computer Studies,
Rajampet, AP, India
e-mail:drmallik2009@gmail.com

C.Sasidhar
Assistant Professor
Annamacharya P.G College of
Computer Studies,
Rajampet, AP, India
e-mail: sasicmca39@gmail.com

V. Satyendra Kumar
Assistant Professor
Annamacharya Institute of
Technology and Sciences,
Rajampet, AP, India
e-mail: sati2all@gmail.com

*Abstract*- **Cloud computing is the new technology that has various advantages and it is an adoptable technology in this present scenario. The main advantage of the cloud computing is that this technology reduces the cost effectiveness for the implementation of the Hardware, software and License for all. This is the better peak time to analyze the cloud and its implementation and better use it for the development of the quality and low cost education for all over the world. In this paper, we discuss how to influence on cloud computing and influence on this technology to take education to a wider mass of students over the country. We believe cloud computing will surely improve the current system of education and improve quality at an affordable cost.**

*Keywords— Cloud Computing, Education, SAAS, Quality Teaching, Cost effective Cloud, Mobile phone, Mobile Cloud*

## I. INTRODUCTION

Cloud Computing has been one of the most booming technology among the professional of Information Technology and also the Business due to its Elasticity in the space occupation and also the better support for the software and the Infrastructure it attracts more technology specialist towards it. Cloud plays the vital role in the Smart Economy, and the possible regulatory changes required in implementing better Applications by using the potential of Cloud Computing [1][2][3].

The main advantage of the cloud is that it gives the low cost implementation for infrastructure and some higher business units like Google, IBM, and Microsoft offer the cloud for Free of cost for the Education system, so it can be used in right way which will provide high quality education. In this paper, we discussed back ground of this paper in section 2, section 3 presented present scenario for existing systems, in section 4 discussed the proposed cloud system for education and section 5 illustrated merits and section 6 concluded with advantages.

## II. BACKGROUND

The term *cloud computing* is being bandied about a lot these days, mainly in the context of the future of the web. But cloud computing potential doesn't begin and end with the personal computer's transformation into a thin client - the mobile platform is going to be heavily impacted by this technology as well. At least that's the analysis being put forth by ABI Research. Their recent report, Mobile Cloud Computing, theorizes that the cloud will soon become a disruptive force in the mobile world, eventually becoming the dominant way in which mobile applications operate.

With a Western-centric view of the world, it can sometimes be hard to remember that not everyone owns a smart phone. There are still a large number of markets worldwide where the dominant phone is a feature phone. While it's true that smart phones will grow in percentage and feature phones will become more sophisticated in time, these lower-end phones are not going away anytime soon. And it's their very existence which will help drive the mobile cloud computing trend.

Not only is there a broader audience using feature phones in the world, there are also more web developers capable of building mobile web applications than there are developers for any other type of mobile device. Those factors, combined with the fact that feature phones themselves are becoming more capable with smarter built-in web browsers and more alternative browsers available for download, will have an impact on mobile cloud computing growth. As per the above statements, we are proposing to use any software applications on mobiles. They can use in even villages of rural areas in India, because, in India most of the part is covered in rural areas.

## III. PRESENT SCENARIO

In this scenario cloud computing is being looked upon by experts in various domains because of its advantages. Cloud has been used in the business oriented unit and in the current education system in India the teaching via web is not so widely available and adapted. Even if it is available, it is provided at a very high cost. This is mainly because of the high cost of data storage and the software they make use of. Cloud has generated many resources which can be used by various educational institutions and streams where their existing/proposed web based learning systems can be implemented at low cost.

### A. Benefits of Cloud Computing

The advantages that come with cloud computing can help resolving some of the common challenges one might have while supporting an educational institution. [4][5].

### 1) Cost

One can choose a subscription or in, some cases, pay-as-you-go plan –whichever works best with that organization business model.

### 2) Flexibility

Infrastructure can be scaled to maximize investments. Cloud

computing allows dynamic scalability as demands fluctuate.

*3) Accessibility*

This help makes data and services publicly available without make vulnerable sensitive information.



Figure 1: Cloud computing with various components

Some would resort to a cloud computing vendor because of the lack of resources while others have the resources to build their cloud computing applications, platforms and hardware. But either way, components have to be implemented with the expectation of optimal performance when we are using through mobile terminals [7].

*4) The Client – The End User*

Everything ends with the client (mobile). The hardware components, the application and everything else developed for cloud computing will be used in the client. Without the client, nothing will be possible. The client could come in two forms: the hardware component or the combination of software and hardware components. Although it's a common conception that cloud computing solely relies on the cloud (internet), there are certain systems that requires pre-installed applications to ensure smooth transition. In this work, all the pre-installed applications can view by mobile devices though clouds. The hardware on the other hand will be the platform where everything has to be launched. Optimization is based on two fronts: the local hardware capacity and the software security. Through optimized hardware with security, the application will launch seamlessly with mobile devices [7].

Cloud computing always has a purpose. One of the main reasons cloud computing become popular is due to the adoption of businesses as the easier way to implement business processes. Cloud computing is all about processes and the services launched through mobile cloud computing always has to deal with processes with an expected output.

*B. Services in Cloud Computing*

Infrastructure as a Service**.** One can get on-demand computing and storage to host, scale, and manage applications and services. Using Microsoft data centers means one can scale with ease and speed to meet the infrastructure needs of that entire organization or individual departments within it, globally or locally [6].

Platform as a Service. The windows azure cloud platform as a service consists of an operating system, a fully relational database, message-based service bus, and a claims access controller providing security-enhanced connectivity and federated access for on premise applications. As a family of on-demand services, the Windows Azure platforms offers organization a familiar development experience, on-demand scalability, and reduce time to market the applications.

Software as a Service. Microsoft hosts online services that provide faculty, staff, and students with a consistent experience across multiple devices.

Microsoft Live at edu provides students, staff, faculty, and alumni long-term, primary e-mail addresses and other applications that they can use to collaborate and communicate online— all at no cost to your education institution.

Exchange Hosted Services offers online tools to help organizations protect themselves from spam and malware, satisfy retention requirements for e-discovery and compliance, encrypt data to preserve confidentiality, and maintain access to e-mail during and after emergency situations.

Microsoft Dynamics CRM Online provides management solutions deployed through Microsoft Office Outlook or an Internet browser to help customers efficiently automate workflows and centralize information. Office Web Apps provide on-demand access to the Web-based version of the Microsoft Office suite of applications, including Office Word, Office Excel, and Office PowerPoint.'

*C. Cloud computing usage*

The cloud plays the main role in the business role and also it is the only elastic data centre which wrapped around various new technologies into it. The technology is most probably used in the business oriented scenario than the service motivated organization as per the survey did by us. According to the Survey made during the month of October 2010 based on the questionnaire prepared by us we found that a major part of the survey group knew about cloud computing (Figure. 1), 69% knew that cloud is used in business, 12% knew it is used in education, 88% agree to implement the cloud for education sector, 94% believes that the cloud technology can reduce the cost of high quality education system and most of them are unaware that the cloud is also offered at low cost.

The following chart describes the statistical report generated based on various surveys done on the cloud computing methodologies and also based on the usage of the cloud in various sectors like data sharing, web mail services, store personal photos online applications such as Google Documents or Adobe Photoshop Express, store personal videos online, pay to store computer files online and also for back up hard drive to an online site.

*D. Requirements for Cloud*

In the previous generation of the information technology the data sharing which led the path for the knowledge sharing was not used by the users globally, in this generation the various streams have the knowledge of e-Learning and the Mobile based learning. In this present context the usage of the central data centre is a easy process for the education system however the cost of implementation and the maintenance of the data storage space and also the load capability also software licensing depends on the real time usage of these systems. Business streams can make revenue out of those expenses whereas for educational institutions which really want to motivate the learners and want to offer a quality education at affordable cost can achieve this by spending a large amount. This can be overcome by the present cloud computing technology that is "Pay as Use" (PAU).

### IV. PROPOSED SYSTEM FOR EDUCATION

In this proposed system the main advantage is that the cloud computing has been used here to overcome all the drawbacks of the present system. Cloud computing is designed as that it can support any group of the users in all the criteria like the software, hardware, infrastructure, storage space everything is provided by the cloud and also the user has to pay as much as they need and don't want to pay for unused and not required as we do for the present data centers.

This cloud system has the cloud model and the client model for its implementation. In the cloud model it is designed as if its suits all the requirements and also the basic consideration. Here the various devices are used and also it has the device control, memory control, load control and several organizing units with the resources like networks, high bandwidth support also the main Security constrains with the filters and the firewall for the better maintenance of the data with proper backup methodology and now the system is ready to provide the data from the cloud.

The main process of this cloud model (Figure. 2) is that it provides the data manipulation operation with the load control and also high authentication and the authorization that too based on the external needs the size of the cloud usage varies so flexible and elastic in means of the data storage and the load accessibility but no compromise on security or the data backup.



Figure 2: The Cloud Model Infrastructure and its Integrated resources

In this client model (Figure. 3) where the data is to be distributed, so that knowledge resources will be used by the all sorts of user in the education streams. Here the applications are developed which will mainly concentrate on the invalid usage of cloud and the data has to be managed and send to the user based on the various data centers available so this client model will check for the registration and valid clients to login into the system and use the application and also the security is maintained in this model so that the data are safeguarded. This model will be used to access the data and share the knowledge.



Figure 3: The Client Model Infrastructure and it's integrate resources and also the various processing layers in the model

#### A. Mobile- Learning

This is a system which is implemented for education using cloud computing. The main objective of Mobile-Learning is that the learners can get the knowledge from the centralized shared resources at any time and any where they like to read that too at free of cost.

Mobile-learning is a system where one can learn through any source on topics of his choice without the need of storing everything in his device. As-you-pay and that much can you can use the services from the cloud data centers for learning selected topics over mobile phone even you in a small village or remote area. For example, if student want learn a JAVA technologies from his agricultural land and works.

#### B. Functionality of Mobile-Learning

The person who wants to make use of Mobile-learning (Figure. 4) has to register and get the credentials to use it via web. It can also be downloading as a mobile application which will be installed in the mobile and through the GPRS/WIFI connectivity they can access the content over the cloud and the user can select among various available topics the one he needs. The topic might be Text based

documents , audio and video files which will be buffered from the cloud to that mobile user and downloaded in the mobile if the memory is available in the mobile(if the user wishes to do so). The user can read the documents, look at the video tutorials, listen to lectures or seminars and finally they can take up self assessments. They will be given a results analysis so that they can evaluate their strengths & weaknesses on their own. This system helps to "Learn while you roam" and also education for all at any time any where globally. Experts can also share their valid tutorials in to the cloud for development of the education community.



Figure 4: The Process flow of Mobile-Learning Cloud Computing

## C. Mobile-Learning Cloud Model

In this cloud model (Figure 5) the User will access the cloud space using his/her credentials so that the required data will be shared from the cloud based on the client request only for the authenticated user.



Figure 5: The Process flow of Mobile-Learning Cloud Model

In this paper, the process flow of mobile learning cloud as we shown in figure 5 which is having 6 steps. Data storage is used for storing the huge data where users are retrieving or handle the data from the data centers. Memory management is organizing and manages the data which is coming from clouds to mobile subscribers and process layer is interacting with security firewalls and memory management.

## D. Mobile-Learning Client Model

In this client model (Figure 6) the user has to download this application and install in their Personal Digital Assistance (PDA) devices or in their mobile phones. The user has to connect to GPRS / Bluetooth / Wi-Fi and connect to the cloud network and get the required topics and based on the selected topic the materials will be downloaded to the mobile for the reading process.



Figure 6: The Process flow of Mobile-Learning Client Model

As we shown in figure 6, the mobile users retrieve the data either in the form of text/video/voice from the cloud center. The subscribers are select which they want to download or retrieve from the data centers with the help of self assistance and keep downloaded data at mobile database.

## E. Advantages with Cloud in Mobile-Learning

In this Mobile-Learning the cloud plays a vital role because the data sharing is the very important role of this learning system, so cloud takes the responsibility of data sharing security and also the load management during the peak hours of access without affecting the network band access. The cloud helps to increase the storage space if the data content are posted more by the users and also during peak hours the total number of user who uses the system will be increased so the load has to be tolerated automatically. Some of the companies offer the cloud at free of cost or at economy prices so this cloud computing will helps in offering the very quality high class education at affordable price.

Today there are more direct applications for teaching and learning as opposed to simple platform-independent tools and scalable data storage

The web search found that organizations of all sizes were using mobile devices for learning because technological advances meant that there was no longer the need for large infrastructure and support costs, and even small enterprises

could deliver mobile learning simply by structuring learning around web-based content that could be accessed from web-enabled mobile devices.

The economics of cloud computing provide a compelling argument for mobile learning. Cloud-based applications can provide students and teachers with free or low-cost alternatives to expensive, proprietary productivity tools. For many institutions, cloud computing offers a cost-effective solution to the problem of how to provide services, data storage, and computing power to a growing number of Internet users without investing capital in physical machines that need to be maintained and upgraded on-site.

Teachers don't have to worry about using outdated or different versions of software. As an Algebra department, we often use Classroom Performance System (CPS) and when we send our test files to each other, this often generates errors because we don't have the same version installed. The same goes with Microsoft Word documents. The cloud will take care of issues like these.

Students and teachers will have 24/7 access to not only their files, but their applications as well (provided they have Internet access). This means that if we need to create an assessment using the CPS software, we don't have to worry about having it installed on every single computer we need to access it from.

The Cloud Computing Opportunity by the Numbers, a multitude of interesting and convincing figures that back the claim that there is a huge opportunity arising in the cloud. With all due credit for the stats going to the article from Reuven's website: ElasticVapor, we wanted to share some of these interesting points:

- There are 50 Million servers worldwide today. By 2013 60% of server workload will be virtualized

- In 2008 the amount of digital information increased by 73%

- There were 360,985,492 internet users in 2000. In 2009 that number increased to 1,802,330,457. That's roughly 27% of the entire world population.

- 50% of the servers sold worldwide are destined for use in a data centre (the average data centre uses 20 megawatts, 10 times more than data centers in 2000 used.)

- Merrill Lynch predicts that the cloud computing market will reach $ 160 billion by 2011.

- IBM claims Cloud cuts IT labor costs by up to 50% and improves capital utilization by 75%.

V.     MERITS OF CLOUD COMPUTING MOBILE-LEARNING

We described so many advantages offered by cloud computing in mobile education. Following are the some of important merits with mobile computing.

 *A. Lower costs.*

You don't need a high-powered and high-priced computer to run cloud computing web-based applications, since applications run in the cloud, not on the desktop PC. In a Simple way we can run such high configured applications on your mobile with cheap cost. When you're using web-based applications on mobiles need not required any memory space and as no software programs have to be loaded and no document files need to be saved.

*B.   Improved performance*

With fewer overfed programs hogging your mobile memory, we will see better performance from your mobile device. Put it simply, mobiles in a cloud computing system boot and run faster because they have fewer programs and processes loaded into mobile memory.

*C.   Reduced software costs.*

Instead of purchasing expensive software applications, you can get most of what you need for free or in least prices on mobile devices even in rural area.

*D. Instant software updates*

Another software-related advantage in cloud computing is that we are no longer faced with choosing between obsolete software and high upgrade costs. When the app is web-based, updates happen automatically and are available the next time you log on to the cloud. When you access a web-based application, you get the latest version without needing to pay for or download an upgrade with our mobile device.

*E.   Improved document format compatibility*

In mobile cloud computing, we have more compatibility for opening the files, applications easily with installation of several software's on mobile device.

*F.   Increased data reliability*

In desktop computing, in which a hard disk crash can destroy all your valuable data, a computer crashing in the cloud should not affect the storage of your data. Even, if your personal computer crashes, all your data is still out there in the cloud, still accessible. Hence, cloud computing is the ultimate in data-safe computing.

*G.   Universal document access*

All your documents are instantly available from wherever you are and there is simply no need to take your documents with you.

*H.   Device independence.*

Finally, here's the ultimate cloud computing advantage: You're no longer tethered to a single computer or network. Change computers, and your existing applications and documents follow you through the cloud. Move to a portable device which is mobile phone, and your applications and documents are still available. There is no need to buy a special version of a program for a particular device, or to save your document in a device-specific format.

VI.     CONCLUSIONS

The cloud computing has the significant scope to change the whole education system. In present scenario the e-learning is getting the popularity and this application in cloud computing will surely help in the development of the

education offered to poor people which will increase the quality of education offered to them. Cloud based education will help the students, staff, Trainers, Institutions and also the learners to a very high extent and mainly students from rural parts of the world will get an opportunity to get the knowledge shared by the professor on other part of the world. Even governments can take initiatives to implement this system in schools and colleges in future and we believe that this will happen soon.

## REFERENCES

[1] Bacigalupo, David; Wills, Gary; De Roure, David; *Victor, A Categorization of Cloud Computing Business Models: IEEE/ACM May 2010.*

[2] Minutoli, G. Fazio, M. Paone, M. Puliafito, A. Engineering Fac, Univ. of Messina, Messina, Italy *Virtual Business Networks With Cloud Computing and Virtual Machines: IEEE/ICUMT Oct 2010.*

[3] Paul Hofmann, *SAP Labs,* Dan Woods, *CITO Research:* The Limits of Public Clouds for Business Applications: *Digital Library* November/December 2010.

[4] Uhlig,R., Neiger, G. Rodgers, D. S.M. Kagi, A.Leung, F.H. Smith : *Intel Corp., USA : Intel visualization technology IEEE Computer Society :* May 2005.

[5] Perez, R., van Doom, L., Sailer, R. *IBM T.J. Watson Res. Center, Yorktown Heights, NY : Visualization and Hardware-Based Security* -October 2008.

[6] GregorPetri: The Data Center is dead; *long live the Virtual Data Center? Join the Efficient Data Center* Nov 2010.

[7] The Independent Cloud Computing and Security *http://cloudsecurity.org/forum/stats/-Augast* 2010.

[8] *Dr. Rao Mikkilineni & Vijay Sarathy, Kawa Objects, Inc.* Cloud Computing and Lessons: *IEEE —June 2009*

[9] *Wenke Ji, Jiangbo Ma, XiaoYong Ji - A* Reference Model of Cloud Operating and Open Source Software Implementation Mapping. - *IEEE -June 2009*

[10] *Pankaj Goyal, Senior Member, IEEE, Rao Mikkilineni, Murthy Ganti:* FCAPS in the Business Services Fabric Model. - *IEEE-June 2009*

[11] *Tarry Singh, VMware vExpert and Cloud Technologist, Pavan Yara, Researcher and Cloud Techno:* Smart Metering the Clouds: *IEEE — June 2009*

[12] *Pankaj Goyal, Senior Member, IEEE, Rao Mikkilineni, Murthy Ganti:* Manageability and Operability in the Business Services Fabric. - *IEEE -June 2009*

[13] *Miyuki Sato, Fujitsu Co. Ltd Japan -* Creating Next Generation Cloud Computing based Network Services and The Contributions of Social Cloud Operation Support System (OSS) to Society Fabric. - *IEEE -June 2009*

[14] Borje Ohlman Anders Eriksson Rene Rembarz Borje. Ohlman@ericsson. Com Anders.E.Eriksson@ericsson. *comRene.Rembarz@ericsson.com Ericsson Research* - What Networking of Information Can Do for Cloud Com. - *IEEE-June 2009*

[15] N.Mallikharjuna Rao, V.Sathyendra Kumar, D.Sudhakar and P.Seetharam, "Cloud computing approaches for educational institutions", International Journal of Computational intelligence and information and Security, ISSN: 2150-5570, Vol 1 No 7 page no: 20-28, IJCIIS September 2010.

[16] Roy Bragg. Cloud computing: When computers really rule. Tech News World, July 2008. Electronic Magazine, available at http://www.technewsworld.com/story/63954.html.

[17] Rajkumar Buyya, Chee Shin Yeo, and Srikumar Venugopal. Market-oriented cloud computing: Vision, hype, and reality for delivering it services as computing utilities. CoRR,(abs/0808.3558), 2008.

[18] Brian de Haaff. Cloud computing - the jargon is back! Cloud Computing Journal, August 2008. Electronic Magazine, article available at http://cloudcomputing.sys-con.com/node/613070

[19] Kemal A. Delic and Martin Anthony Walker. Emergence of the academic computing clouds. ACM Ubiquity, (31), 2008.

[20] Flexi scale web site. http://www.flexiscale.com, last visited: August 2008.

## AUTHORS PROFILE

N.Mallikharjuna Rao is presently working as Associate Professor in the department of Master of Computer Applications at Annamacharya PG college of Computer Studies, Rajampet and having more than 12 years of Experience in Teaching UG and PG courses. He received his B.Sc Computer Science from Andhra University in 1995, Master of Computer Applications (MCA) from Acharya Nagarjuna University in 1998, Master of Philosophy in Computer science from Madurai Kamarj University, Tamilnadu, in India and Master of Technology in Computer Science and Engineering from Allahabad Agricultural University, India. He is a life Member in ISTE and Member in IEEE, IACSIT. He is a research scholar in Acharya Nagarjuna University under the esteemed guidance of Dr.M.M.Naidu, Principal, SV.Univeristy, Tirupathi, and AP.

Mr.V.Sathyendra Kumar is working as Assistant Professor in the Department of MCA at Annamacharya Institute of Technology & Sciences, Rajampet. He has more than 4 years of Teaching Experience for PG. He is life member in ISTE.

Mr. C. Sasidhar is working as Assistant Professor in the Department of MCA, Annamacharya P.G College of Computer Studies, Rajampet. He has more than 5 years of Teaching Experience for PG Courses. He received his MCA and M.Tech (CS) from JNTUA in the year 2005, 2010 respectively. He is life member in ISTE.

# Robust R Peak and QRS detection in Electrocardiogram using Wavelet Transform

P. Sasikala

Research Scholar, AP/Dept. Of Mathematics
V.M.K.V. Engineering College
Salem, Tamilnadu, India
Rgsasi@Gmail.Com

Dr. R.S.D. Wahidabanu

Professor & Head/ Dept. Of Ece
Govt. College of Engineering
Salem, Tamilnadu, India
Drwahidabanu@Gmail.Com

*Abstract—* **In this paper a robust R Peak and QRS detection using Wavelet Transform has been developed. Wavelet Transform provides efficient localization in both time and frequency. Discrete Wavelet Transform (DWT) has been used to extract relevant information from the ECG signal in order to perform classification. Electrocardiogram (ECG) signal feature parameters are the basis for signal Analysis, Diagnosis, Authentication and Identification performance. These parameters can be extracted from the intervals and amplitudes of the signal. The first step in extracting ECG features starts from the exact detection of R Peak in the QRS Complex. The accuracy of the determined temporal locations of R Peak and QRS complex is essential for the performance of other ECG processing stages. Individuals can be identified once ECG signature is formulated. This is an initial work towards establishing that the ECG signal is a signature like fingerprint, retinal signature for any individual Identification. Analysis is carried out using MATLAB Software. The correct detection rate of the Peaks is up to 99% based on MIT-BIH ECG database.**

*Keywords- Electrocardiogram, Wavelet Transform, QRS complex, Filters, Thresholds*

## I. INTRODUCTION

The **Electrocardiogram** is the electrical manifestation of the contractile activity of the heart. It is a graphical record of the direction and magnitude of the electrical activity that is generated by depolarization and repolarization of the atria and ventricles. It provides information about the heart rate, rhythm, and morphology. The importance of the Electrocardiography is remarkable since heart diseases constitute one of the major causes of mortality in the world. ECG varies from person to person due to the difference in position, size, anatomy of the heart, age, relatively body weight, chest configuration and various other factors. There is strong evidence that heart's electrical activity embeds highly distinctive characteristics, suitable for various applications and diagnosis.

The ECG is characterized by a recurrent wave sequence of P, QRS, T and U wave associated with each beat. The QRS complex is the most striking waveform, caused by ventricular depolarization of the human heart. A typical ECG wave of a normal heartbeat consists of a *P* wave, a *QRS* complex, and a *T* wave. Fig. 1 depicts the basic shape of a healthy ECG heartbeat signal with P, Q, R, S, J, T and U characteristics and the standard ECG intervals QT interval, ST interval and PR interval.



Figure 1. An ECG waveform with the standard ECG intervals

QRS detection is one of the fundamental issue in the analysis of Electrocardiographic signal. The QRS complex consists of three characteristic points within one cardiac cycle denoted as Q, R and S. The QRS complex is considered as the most striking waveform of the electrocardiogram and hence used as a starting point for further analysis or compression schemes. The detection of a QRS complex seems not to be a very difficult problem. However, in case of noisy or pathological signals or in case of strong amplitude level variations, the detection quality and accuracy may decrease significantly.

Numerous QRS detection algorithms such as derivative based algorithms [1-4], wavelet transform [5], Filtering Techniques [6] artificial neural networks [7-9], genetic algorithms [10], syntactic methods [11], Hilbert transform [12], Markov models [13] etc. are reported in literature. Kohler et al [14] described and compared the performance of all these QRS detectors. Recently few other methods based on pattern recognition [15], moving- averaging [16] etc are proposed for the detection of QRS complex. Once the position of the QRS complex is obtained, the location of other components of ECG like P, T waves and ST segment etc. are found relative to the position of QRS, in order to analyze the complete cardiac period.

Recently Wavelet Transform has been proven to be useful tool for non-stationary signal analysis. Among the existing wavelet approaches, (continuous, dyadic, orthogonal, biorthogonal), we use real dyadic wavelet transform because of its good temporal localization properties and its fast calculations. Discrete Wavelet Transform can be used as a good tool for non-stationary ECG signal detection. DWT is a sampled version of the Continuous Wavelet Transform (CWT) in a dyadic grid.

## II. WAVELET TRANSFORM

The Wavelet Transform is a time-scale representation that has been used successfully in a broad range of applications, in particular signal compression. Recently, Wavelets have been applied to several problems in Electrocardiology, including data compression, analysis of ventricular late potentials, and the detection of ECG characteristic points. The Wavelet transformation is a linear operation that decomposes the signal into a number of scales related to frequency components and analyses each scale with a certain resolution [17-21].

The WT uses a short time interval for evaluating higher frequencies and a long time interval for lower frequencies. Due to this property, high frequency components of short duration can be observed successfully by Wavelet Transform. One of the advantage of the Wavelet Transform is that it is able to decompose signals at various resolutions, which allows accurate feature extraction from non-stationary signals like ECG. A family of analyzing wavelets in the time frequency domain is obtained by applying a scaling factor *and* a translation factor *to* the basic mother wavelet.

Wavelet Transform of a signal $f(t)$ is defined as the sum of over all time of the signal multiplied by scaled, shifted versions of the wavelet function ψ and is given by,

$$W(a,b) = \int_{-\infty}^{\infty} f(t)\psi_{a,b}(t)dt \qquad (1)$$

$$\psi_{a,b}(t) = \frac{1}{\sqrt{a}}\psi *\left(\frac{t-b}{a}\right) \qquad (2)$$

Where * denotes complex conjugation and $\psi_{a,b}(t)$ is a window function called the mother wavelet, '***a***' is a scale factor and '***b***' is a translation factor. Here $\psi\left(\dfrac{t-b}{a}\right)$ is a shifted and scaled version of a mother wavelet which is used as bases for wavelet decomposition of the input signal. One of the key criteria of a good mother wavelet is its ability to fully reconstruct the signal from the wavelet decompositions.

The DWT is sufficient for most practical applications and for the reconstruction of the signal [22]. The DWT provides enough information, and offers a significant reduction in the computation time. Here, we have discrete function $f(n)$ and the definition of DWT is given by:

$$W(a,b) = c(j,k) = \sum_{n \in z} f(n)\psi_{j,k}(n) \qquad (3)$$

where $\psi_{j,k}(n)$ is a discrete wavelet defined as

$$\psi_{j,k}(n) = 2^{-\frac{j}{2}}\psi(2^{-j}n-k) \qquad (4)$$

The parameters a, b are defined in such a way that $a = 2^j$, $b = 2^j k$. In the DWT analyses, the signal at different frequency bands and at different resolutions is decomposed into a 'coarse approximation' and 'detailed information'. Two sets of functions are employed by the DWT, the scaling functions (associated with the low pass filter) and the wavelet functions (associated with the high pass filter). The signal is filtered by passing it through successive high pass and low pass filters to obtain versions of the signal in different frequency bands.

The original signal $x(n)$ is passed through a half band low pass and high pass filter. With the signal highest frequency being π/2, half of the samples are eliminated adhering to the Nyquist criterion. Thus, the signal can be sub-sampled by 2 as shown in Equation (5). Thus mathematically, this can be written as:

$$y(n) = \sum_n h(k)x(2n-k) \qquad (5)$$

$$y_{high}(k) = \sum x(n)g(2k-n) \qquad (6)$$

$$y_{low}(k) = \sum x(n)h(2k-n) \qquad (7)$$

The decomposition perform halves the time resolution and at the same time doubles the frequency resolution. Thus, at every level, the filtering and sub-sampling will result in half the time resolution and double the frequency resolution. The successive Low Pass Filter (LPF) and High Pass Filter (HPF) of the discrete time-domain signal are called the Mallat algorithm or Mallat Tree Decomposition (MTD). The sequence $x(n)$ is passed through several levels made up of low pass $g(n)$ and high pass $h(n)$ analysis filters. At each level, 'detail information' $d_j[n]$ is produced by the high pass filter while the 'coarse approximations' $a_j[n]$ is produced by the low pass filter.

The maximum number of levels of decomposition depends upon the length of the signal as shown in Fig. 2. The Discrete Wavelet Transform of the original signal is obtained by concatenating all the coefficients, a_j[n] and d_j[n].

Figure 2.  Three level wavelet decomposition tree

The reconstruction process is the reverse of decomposition, where the approximation and detail coefficients at every level are up-sampled by 2 and passed through low-pass $g(n)$ and high pass $h(n)$ synthesis filters and finally added as shown in Fig. 3. The same number of levels is taken as in the case of decomposition.



Figure 3.  Three level wavelet reconstruction tree

Wavelet Transform is popular because it satisfies energy conservation law and original signal can be reconstructed. It is obvious that Wavelet Transform at scale '*a*' is proportional to the derivative of the filtered signal with a smoothing impulse response at scale '*a*'. Therefore, local maxima or minima of the smoothed signal will occur on the zero crossings of the Wavelet Transform at different scales. Maximum absolute values of the Wavelet Transform will show the maximum slopes in the filtered signal.

*A.  Wavelet Selection*

The use of the Wavelet Transform has gained popularity in time-frequency analysis because of the flexibility it offers in analyzing basis functions. The selection of relevant wavelet is an important task before starting the detection procedure. The choice of wavelet depends upon the type of signal to be analyzed. The wavelet similar to the signal is usually selected. The are several wavelet families like Harr, Daubechies, Biorthogonal, Coiflets, Symlets, Morlet, Mexican Hat, Meyer etc. and several other Real and Complex wavelets. However, Daubechies (Db4) Wavelet has been found to give details more accurately than others [23]. Moreover, this Wavelet

shows similarity with QRS complexes and energy spectrum is concentrated around low frequencies. Therefore, we have chosen Daubechies (Db4) Wavelet for extracting ECG features in our application [22]. The Daubechies Wavelet is shown in Fig. 4.



Figure 4.  Daubechies4 Wavelet

### III.  DATA

ECG signals required for analysis are collected from Physionet MIT-BIH arrhythmia database where annotated ECG signals are described by a text header file (.hea), a binary file (.dat) and a binary annotation file (.atr).  The database contains 48 records, each containing two-channel ECG signals for 30 min duration selected from 24-hr recordings of 47 different individuals. Header file consists of detailed information such as number of samples, sampling frequency, format of ECG signal, type of ECG leads and number of ECG leads, patient's history and the detailed clinical information. In binary data signal file, the signal is stored in 212 format which means each sample requires number of leads times 12 bits to be stored and the binary annotation file consists of beat annotations.  Signals were sampled using a 12-bit analog-to-digital converter board (National Instruments, PCI-6071E). Matlab and its wavelet toolbox were used for ECG Signal processing and Analysis. Analysis was performed on the PQRST waveform.

### IV.  METHODOLOGY

In order to extract useful information from the ECG signal, the raw ECG signal should be processed. ECG signal processing can be roughly divided into two stages by functionality: Preprocessing and Feature Extraction as shown in Fig. 5.

Figure 5. Structure of ECG Signal Processing.

Feature Extraction is performed to form distinctive personalized signatures for every subject. The purpose of the feature extraction process is to select and retain relevant information from original signal. The Feature Extraction stage extracts diagnostic information from the ECG signal. The preprocessing stage removes or suppresses noise from the raw ECG signal. A Feature Extraction method using Discrete Wavelet Transform was proposed by Emran et al [22].

### A. Preprocessing

ECG signal mainly contains noises of different types, namely frequency interference, baseline drift, electrode contact noise, polarization noise, muscle noise, the internal amplifier noise and motor artifacts. Artifacts are the noise induced to ECG signals that result from movements of electrodes. One of the common problems in ECG signal processing is baseline wander removal and noise suppression.

#### 1) Removal of the baseline drift

Baseline wandering is one of the noise artifacts that affect ECG signals. Removal of baseline wander is therefore required in the analysis of the ECG signal to minimize the changes in beat morphology with no physiological counterpart. Respiration and electrode impedance changes due to perspiration are important sources of baseline wander in most types of ECG recordings. The frequency content of the baseline wander is usually in a range well below 0.5Hz. This baseline drift can be eliminated without changing or disturbing the characteristics of the waveform.

We use the median filters (200-ms and 600-ms) [24] to eliminate baseline drift of ECG signal. The process is as follows

a) The original ECG signal is processed with a median filter of 200-ms width to remove QRS complexes and P waves.

b) The resulting signal is then processed with a median filter of 600-ms width to remove T waves. The signal resulting from the second filter operation contains the baseline of the ECG signal.

c) By subtracting the filtered signal from the original signal, a signal with baseline drift elimination can be obtained.

#### 2) Removal of the NOISE

After removing baseline wander, the resulting ECG signal is more stationary and explicit than the original signal. However, some other types of noise might still affect feature extraction of the ECG signal. In order to reduce the noise many techniques are available like Digital filters, Adaptive method and Wavelet Transform thresholding methods. Digital filters and Adaptive methods can be applied to signal whose statistical characteristics are stationary in many cases. However, for non-stationary signals it is not adequate to use Digital filters or Adaptive method because of loss of information. To remove the noise, we use Discrete Wavelet transform.

This first decomposes the ECG signal into several subbands by applying the Wavelet Transform, and then modifies each wavelet coefficient by applying a threshold function, and finally reconstructs the denoised signal. The high frequency components of the ECG signal decreases as lower details are removed from the original signal. As the lower details are removed, the signal becomes smoother and the noise disappears since noises are marked by high frequency components picked up along the ways of transmission. This is the contribution of the discrete Wavelet Transform where noise filtration is performed implicitly. The preprocessed signal using DWT is shown in Fig. 6.



Figure 6. Baseline removed and Denoised Signal

## B. Detection of R peak and QRS

In order to detect the peaks, specific details of the signal are selected. The detection of R peak is the first step of feature extraction. The R peak in the signal from the Modified Lead II (MLII) lead has the largest amplitude among all the waves compared to other leads. The QRS complex detection consists of determining the R point of the heartbeat, which is in general the point where the heartbeat has the highest amplitude. A normal QRS complex indicates that the electrical impulse has progressed normally from the bundle of His to the Purkinje network through the right and left bundle branches and that normal depolarization of the right and left ventricles has occurred.

Most of the energy of the QRS complex lies between 3 Hz and 40 Hz [25]. The 3-dB frequencies of the Fourier Transform of the wavelets indicate that most of the energy of the QRS complex lies between scales of $2^3$ and $2^4$, with the largest at $2^4$. The energy decreases if the scale is larger then $2^4$. The energy of motion artifacts and baseline wander (i.e., noise) increases for scales greater then $2^5$. Therefore, we choose to use characteristic scales of $2^1$ to $2^4$ for the wavelet.

The detection of the QRS complex is based on modulus maxima of the Wavelet Transform. This is because modulus maxima and zero crossings of the Wavelet Transform correspond to the sharp edges in the signal. The QRS complex produces two modulus maxima with opposite signs, with a zero crossing between them shown in Fig. 7. Therefore, detection rules (thresholds) are applied to the Wavelet Transform of the ECG signal. The Q and S point occurs about the R Peak with in 0.1second. The left point denoted the Q point and the right one denotes the S point. Calculating the distance from zero point or close to zero of left side of R Peak within the threshold limit denotes Q point. Similarly the right side denotes the S point.



Figure 7. Maxima, Minima, and Zero crossing of Wavelet Transform at scale $2^4$

QRS width is calculated from the onset and the offset of the QRS complex. The onset is the beginning of the Q wave and the offset is the ending of the S wave. Normally, the onset of the QRS complex contains the high-frequency components, which are detected at finer scales. To identify the onset and offset of the wave, the wave is made to zero base. The onset is the beginning and the offset is the ending of the first modulus

maxima pair. Once this QRS complex is located the next step is to determine the onset and offset points for each QRS complex and to identify the component waves of the QRS complex. The R peak and QRS complex is shown in Fig. 8.





Figure 8. a) R-Peak  b) QRS Complex

## V. CONCLUSION

An algorithm for R Peak and QRS complex detection using Wavelet Transform technique has been developed. Table - 1 shows the detection results on the whole database. The information about the R Peak and QRS complex obtained is very useful for ECG Classification, Analysis, Diagnosis, Authentication and Identification performance. The QRS complex is also used for beat detection and the determination of heart rate through R-R interval estimation. This information can also serve as an input to a system that allows automatic cardiac diagnosis. The overall sensitivity of the detector improves. The main advantage of this kind of detection is less time consuming for long time ECG signal.

TABLE - 1: TEST RESULTS SHOW THE DETECTION RESULTS.

| Record | Total beats | FP | FN | FP + FN | Detection Error Rate | Sensitivity |
|--------|-------------|----|----|---------|----------------------|-------------|
| 100 | 2272 | 2 | 0 | 2 | 0.09 | 99.96 |
| 105 | 2543 | 18 | 11 | 29 | 1.14 | 100.00 |
| 108 | 1775 | 27 | 35 | 62 | 3.49 | 99.83 |
| 115 | 1953 | 0 | 0 | 0 | 0.00 | 99.85 |
| 118 | 2278 | 10 | 3 | 13 | 0.57 | 99.88 |
| 124 | 1473 | 1 | 2 | 3 | 0.20 | 99.75 |
| 200 | 2601 | 0 | 1 | 1 | 0.04 | 99.96 |
| 202 | 2136 | 8 | 0 | 8 | 0.37 | 100.00 |
| 208 | 2956 | 0 | 5 | 5 | 0.17 | 99.83 |
| 210 | 2647 | 0 | 4 | 4 | 0.15 | 99.85 |
| 215 | 3256 | 0 | 4 | 4 | 0.12 | 99.88 |
| Average | 2354 | 6 | 6 | 12 | 0.58 | 99.89 |

Further, ECG signal is a life indicator, and can be used as a tool for liveness detection. The physiological and geometrical differences of the heart in different individuals display certain uniqueness in their ECG signals. Hence ECG can be used as a Biometric tool for Identification and Verification of Individuals. The advantage of ECGs in biometric systems is their robust nature against the application of falsified credentials. Amplitude of P wave remains constant throughout the life and other amplitude features are changes on small scale. Future work is to calculate the amplitude distance between ECG features and comparison will be made for Identification. Further verification of individuals can be done by using statistical theory of Sequential Probability procedures.

## REFERENCES

[1] J. Pan and W.J. Tompkins, "A real-time QRS detection algorithm", *IEEE Trans. Biomed. Eng.*, vol. 32, pp. 230–236, 1985.

[2] V.X. Afonso, W.J. Tompkins, T.Q. Nguyen, and S. Luo, "ECG beat detection using filter banks", *IEEE Trans. Biomed. Eng.*, vol. 46, pp. 192–202, 1999.

[3] J. Fraden and M.R. Neumann, "QRS wave detection", *Med. Biol. Eng. Comput.*, vol. 18, pp. 125–132, 1980.

[4] W.P. Holsinger, K.M. Kempner, and M.H. Miller, "A QRS preprocessor based on digital differentiation", *IEEE Trans. Biomed. Eng.*, vol. 18, pp. 121–217, May 1971.

[5] M. Bahoura, M. Hassani, and M. Hubin, "DSP implementation of wavelet transform for real time ECG wave forms detection and heart rate analysis", *Comput. Methods Programs Biomed.* vol. 52, no. 1, pp. 35–44, 1997.

[6] S. A. Israel et al. "ECG to identify individuals". Pat. Rec. 38:133-142, 2005.

[7] Y.H. Hu, W.J. Tompkins, J.L. Urrusti, and V.X. Afonso, "Applications of artificial neural networks for ECG signal detection and classification", J. *Electrocardiology*, vol. 26 (Suppl.), pp. 66-73, 1993.

[8] G. Vijaya, V. Kumar, and H.K. Verma, "ANN-based QRS-complex analysis of ECG," *J. Med. Eng. Technol.*, vol. 22, no. 4, pp. 160-167, 1998.

[9] Q. Xue, Y. H. Hu, and W. J. Tompkins, "Neural-network-based adaptive matched filtering for QRS detection", *IEEE Trans. Biomed. Eng.*, vol. 39, pp. 317-329, 1992.

[10] R. Poli, S. Cagnoni, and G. Valli, "Genetic design of optimum linear and nonlinear QRS detectors", *IEEE Trans. Biomed. Eng.*, vol. 42, pp. 1137-1141, 1995.

[11] E. Skordalakis, "Syntactic ECG processing: a review", *Pattern Recognition.*, vol. 19, no. 4, pp. 305–313, 1986.

[12] S. K. Zhou, J.-T. Wang and J.-R. Xu, "The real-time detection of QRS-complex using the envelop of ECG", in *Proc. 10th Annu. Int. Conf., IEEE Engineering in Medicine and Biology Society,* New Orleans, LA, 1988, p. 38.

[13] R. V. Andreao, B. Dorizzi, and J. Boudy, "ECG signal analysis through hidden Markov models", *IEEE Transactions on Biomedical Engineering*, vol. 53, no. 8, pp. 1541–1549, 2006.

[14] B.U. Kohler, C. Hennig, and R. Orglmeister, "The principles of software QRS detection", Engineering *in Medicine and Biology Magazine, IEEE*, vol. 21, pp. 42 – 57, Jan.-Feb. 2002

[15] S. S. Mehta et al. "Computer-aided interpretation of ECG for diagnostics", Int. Journal of System Science, 43-58, 1996.

[16] S. W. Chen et al. "A real time QRS detection method based on moving-averaging incorporating with wavelet denoising", Comp. Methods and Progs. in Biomed. 82:187-195, 2006.

[17] Mallat S: "Multiresolution frequency channel decomposition of images and wavelet models", *IEEE Transactions on Acoust, Speech Signal Processing,* 37, No. 12: 2091-2110, 1989.

[18] Chui K.: "An Introduction to Wavelets", Academic Press, Inc.1992.

[19] S. Kadambe, R. Murray, and G.F. Boudreaux -Bartels, "Wavelet transform - based QRS complex detector", *IEEE Trans. Biomed.Eng.*, vol. 46, pp. 838–848, 1999.

[20] Cuiwei Li, Chongxun Zheng, and Changfeng Tai, "Detection of ECG Characteristic Points using Wavelet Transforms", IEEE Transactions on Biomedical Engineering, Vol. 42, No. 1, pp. 21-28, 1995.

[21] Saritha, V. Sukanya, and Y. Narasimha Murthy, "ECG Signal Analysis Using Wavelet Transforms", Bulgarian Journal of Physics, vol. 35, pp. 68-77, 2008.

[22] E. M. Tamil, N. H. Kamarudin, R. Salleh, M. Yamani Idna Idris, M. N. Noorzaily, and A. M. Tamil, (2008) Heartbeat electrocardiogram (ECG) signal feature extraction using discrete wavelet transforms (DWT), in Proceedings of CSPA, 1112–1117.

[23] S. Z. Mahmoodabadi, A. Ahmadian, and M. D. Abolhasani, "ECG Feature Extraction using Daubechies Wavelets", Proceedings of the fifth IASTED International conference on Visualization, Imaging and Image Processing, pp. 343-348, 2005.

[24] P. de Chazal, C. Heneghan, E. Sheridan, R.Reilly, P. Nolan, M. O'Malley, "Automated Processing of the Single-Lead Electrocardiogram for the Detection of Obstructive Sleep Apnoea", *IEEE Trans. Biomed. Eng.,* 50( 6): 686-689, 2003.

[25] N.V.Thakor, J.G.Webster and W.J.Tompkins, "Estimation of QRS complex power spectra for design of a QRS filter", IEEE Transactions on Biomedical Engineering, vol. BME-31, no. 11,pp. 702-706, 1986, pp. 702-706, Nov,1986.

[26] S.S.Mehta 1, N.S. Lingayat, "Identification of QRS complexes in 12-lead electrocardiogram", Science Direct, Expert Systems with Applications May 2007.

[27] Gordan Cornelia, Reiz Romulus , "ECG signals processing using Wavelets", IEEE, proceedings of the fifth IASTED International conference May 2005.

# Performance Evaluation of Node Failure Prediction QoS Routing Protocol (NFPQR) in Ad Hoc Networks

Dr D Srinivas Rao

Professor,  Dept of ECE, CE

Jawaharlal Nehru Technological University

HYDERABAD - 500 085, INDIA

dsraoece@gmail.com

Sake Pothalaiah

Project Associate, Dept of ECE, CE

Jawaharlal Nehru Technological University

HYDERABAD - 500 085, INDIA

pothalaiahs@gmail.com

*Abstract*—**The characteristics of ad hoc networks make the QoS support a very complex process unlike traditional networks. The nodes in ad hoc wireless networks have limited power capabilities. The node failure in the network leads to different problems such as network topology changes, network partitions, packet losses and low signal quality. Many QoS routing protocols like  Predictive location based QoS routing protocol (PLBQR), Ticket based QoS routing, Trigger based  distributed QoS routing (TDR) protocol ,Bandwidth routing(BR) protocol, Core extracted distributed routing (CEDAR) protocol have been proposed. However these algorithms do not consider the node failures and their consequences in the routing. Thus most of the routing protocols do not perform well in frequent or unpredictable node failure conditions.  Node Failure Predication QoS Routing" (NFPQR) scheme provides an optimal route selection by predicting the possibility of failure of a node through its power level. The NFPQR protocol has been modified as C-NFPQR (Clustered NFPQR) in order to provide power optimization using clustered based approach. The performance of the NFPQR and C-NFPQR is evaluated through the same QoS parameters.**

*Keywords- NFPQ;, C-NFPQR; CEDAR; PLBQR; TDR; QoS; Routing Protocols;*

## I.    INTRODUCTION

Ad hoc Network is a wireless network consisting of a collection of mobile nodes with no fixed infrastructure. Networking infrastructure refers to the facility of which sole purpose is network management and routing. The network nodes communicate with one another over scarce wireless channel in multi hop fashion.  Each node behaves as a router and it takes part in discovery and maintenance of routes to other nodes in the network. The main characteristics of this network are dynamic topology, bandwidth constraint, variable capacity links, power constrained operation, limited physical security and quickly deployable. One of the major research issues is related to these networks is quality of service (QoS) routing.

## II.    RELATED PROTOCOLS

Quality of service is set of service requirements provided to certain traffic by the network to meet the satisfaction of the user of that traffic. It has been investigated by different researchers and several proposals have been published to address how the QoS can be supported in MANETs [2].  QoS routing support in MANET still remains as an open problem. In this paper we discuss the related protocols and compare them with our work. The predictive location-based QoS routing protocol (PLBQR)  is based on the prediction of the location of nodes in Ad hoc networks. PLBQR protocol uses location and delay prediction schemes which reduce to some extent the problem arising due to the presence of stale routing information. PLBQR has no resources are reserved along the route from the source to the destination, it is not possible to provide hard QoS guarantees using this protocol. Even soft QoS guarantees may be broken in cases when the network load is high[2].

In the trigger-based (on-demand) distributed QoS routing (TDR) protocol [2] due to small-scale fading, the received power level may vary rapidly over short periods of time or distance traveled. Some of the factors that influence fading are multi-path propagation, velocity of the nodes, and bandwidth of the channel.

The bandwidth routing (BR) protocol [2] consists of an end-to-end path bandwidth calculation algorithm to inform the source node of the available bandwidth to any destination in the ad hoc network, a bandwidth reservation algorithm to reserve sufficient number of free slots for the QoS flow, and a standby routing algorithm to reestablish the QoS flow in case of path breaks. The CDMA-over-TDMA channel model that is used in this protocol requires assigning a unique control slot in the control phase of super-frame for each node present in the network. This assignment has to be done statically before commissioning the network. Due to this, it is not possible for a new node to enter into the network at a later point of time. If a particular node leaves the network, the corresponding control slot remains unused and there is no way to reuse such slot(s). and also network needs to be fully synchronized.

CEDAR: the core-extracted distributed routing (CEDAR) protocol has been proposed as a QOS routing protocol. CEADR dynamically establishes a core of the network, and then incrementally propagates link state of stable high bandwidth to the nodes of the core to identify and avoid using congested parts of the network. The core nodes are elected by approximating a minimum dominating set of the ad hoc network. However, it is overhead for CEDAR

### III. NODE FAILURE PREDICTION QoS ROUTING PROTOCOL (NFPQR)

Most of the QoS routing protocols proposed previously do not perform well in frequent and /or unpredictable node failure conditions. So QoS routing need a relatively accurate prediction of network's future conditions that is not included in the previous works. This work addresses a new routing algorithm NFPQR NFPQR decreases end to end packet delay and some extent of packet loss by predicting the future power level of a node. It calculates the future condition of a node to make it as next relay node in the path discovery. The estimation of future condition of a node depends on the power level of the node at a particular time [2].

The devices generally are dependent on finite battery sources. Once the battery power is completely consumed then the device will go down i.e., the device is considered as under failure .If the radio interface of the mobile device is not functioning, then all the communications from this device will be stopped .A prediction on node failure helps us in providing better QoS routing for ad hoc networks. Suppose a node, having high probability of failure in the near future due to the lack of sufficient power in the battery and the node is selected as a router to forward the packets, after a few seconds the node will be failed and its communication links with its neighbor will be broken. So it cannot forward the packets and many packets will be lost and these packets have to be regenerated and retransmitted then. Another penalty due to the node failure is that the route discovery process should be performed once again to establish new path that may take few more seconds .During this entire process, the packets will be queued up at the downstream node until a new path is set up.

During certain lower power levels, the signal strength is reduced and delay is increased at the MAC level. Sometimes the network partitioning also occurs due the node failures, where the packet may never reach the destination, if the source and the destination nodes are not in the same partition leads to node failure in the network and causes QoS violation which increases end to end delay, packet loss and network throughput.

In order to solve the problems due to node failure and to support QoS ,a new approach is proposed which predicts whether a node will be failed in the near future or not .Before the upstream node is selected as a router to forward the packets ,the downstream nodes predicts whether the upstream node will be failed in the near future or not. The heuristic which is used here is based on the power levels in the battery. Power is consumed during communication and processing or computing. Communication power is much higher than the computing power. In communication power, the transmission power, the power needed to transmit a packet, is much higher than others like receiving power, idle power etc.

If the transmission power is $C_t$ and overhead energy is $C_o$ then the total power needed to transmit the entire buffer is:

$$((B_f. C_t)/ Ps) + C_o \qquad (1)$$

Here $B_f$ is buffer capacity and $Ps$ is packet size.

The threshold power level is based on the packet size, buffer capacity and the packet transfer rate of the node. If t1 is the present time, then the maximum power consumption at a particular node after time t2 is given by:

$$P_{12} = (t2 - t1). (tr. C_t + C_o) \qquad (2)$$

### IV. NFPQR ALGORITHM

When a node j receives a route request message from a node i, then the node j predicts its future condition by considering power level of node j.

If the power level is above the threshold which is 80% of the power of the maximum power node in the random network then the node j will forward the RREQ to the next hop; otherwise it will drop the route request message. The same procedure is repeated for all the nodes till the destination is reached.

Definition (HEAD):

HEAD nodes are nodes such that all non-HEAD nodes (nodes within the transmission range of that HEAD node) are connected to any one of the HEAD node and route packets for any other nodes with the help of Mobile Agents. The route consists of Source node, Corresponding HEAD node, Gateway nodes and intermediate HEAD and Gateway nodes and destination node.

In this paper our topology management scheme HEAD nodes are selected such that HEAD nodes have maximum power level among their one hop neighbors and all non- HEAD nodes are within the transmission range of HEAD nodes. These HEAD nodes have the routing intelligence i.e. they make all decisions related to routing. The gateway nodes are selected which are having enough/high power so that they can forward packets between HEAD nodes and they don't have routing intelligence.

#### A. Head Placement

HEAD nodes along with gateways confirm a path in the virtual backbone, which is used for routing and there is demands for additional power for transmission, reception and processing of packets. Thus the HEAD nodes should be selected in such a way that they have enough/higher power level.

Undecided nodes periodically checks if it has a maximum POWER level among its one-hop neighbors, which have not, joined to any HEAD node (i.e. undecided neighbors). If a node has maximum POWER level among such one hop neighbors, it becomes a HEAD node and declares itself as a HEAD node in the status field of next LINK message and communicates to all its neighbors.

If undecided node knows that its neighbor node has become HEAD node from received LINK message, it changes status to member. It declares its status as member and it is current HEAD node in next LINK message. If more than one neighbors of an undecided node became HEAD, undecided node select its HEAD node from which it has received the

LINK packet earlier There may be undecided nodes whose one hop neighbors with power level more than the undecided node chose to join HEAD nodes, as the HEAD nodes have more power level than its one hop neighbors. Such undecided nodes with maximum power level among one hop undecided neighbors declares themselves as HEAD nodes in the next LINK message.

A HEAD node prepares a list of its member nodes, which are joined to the HEAD node, form the broadcast of LINK messages received from one hop neighbors. This information in the table is periodically changes as a new LINK packet is received.



Figure 1. Flowchart for head placement



Figure 2. Illustration of HEAD node selection in a Random Ad Hoc Network

## B. Head Node Withdrawal

The HEAD node will drain its energy more rapidly, as compared to member nodes. Before the HEAD node loses its major part of its power, the responsibilities of the HEAD node should be transferred to other node with sufficient power level. Also RIMA nodes should not be changed frequently which will increase the overhead.

When a HEAD node observes that its POWER level is gone below a threshold, it will withdraw its status of HEAD node. The withdrawal of HEAD node is declared to its member nodes in the next WAKEUP message as a undecided node. The threshold can be set to 80% of HEAD level when the node decided to become a HEAD node.

When a gate way or member node comes to know that it cannot contact its HEAD node, it changes its status to undecided and starts HEAD node placement procedure.

## C. Gateway Selection

The maximum number of hops between any two close HEAD nodes is two; hence gateways are required and are used to forward the packets between the HEAD nodes. The gateway nodes must have sufficient amount of power, to transmit and receive the packets to and from the HEAD nodes.



Figure 3. Gateway selection and flow in Ad Hoc Network

M – HEAD Nodes; m – Member Nodes; G – Gateways

In NFPQR algorithm, more stable paths are formed during route discovery .Here, the stable path means the packets, which traverses on these paths, will not experience long delays and improves throughput .Also it increases the network life time of the ad hoc networks

## D. Scheduling of Sleep Cycle

The POWER saving features to 802.11 CSMA/CA to make the MAC layer power efficient by using randomized wake up time for member nodes in ad hoc network. HEAD nodes and Gateways continuously stay awake to forward packets of other nodes

### CNFPQR

Minimizing of power consumption is an important challenge in mobile networking. The requirement of co-operation between power saving and routing protocols is particularly acute in the case of multi-hop Ad hoc wireless networks where nodes must forward packets for each other. Although the quality of routing is improved using NFPQR protocol but it failed to impart power saving mechanism. In

this paper implementing C-NFPQR (Clustered NFPQR) protocol is implemented over random ad hoc wireless network in which route selection is done using NFPQR protocol. Proposed approach not only optimizes power consumption along with better quality routing leading to better network performances in terms of various QoS parameters

TABLE I.    NETWORK SCENARIO

| No .of nodes | 20 |
|---|---|
| Source node id | 14 |
| Destination node id | 8 |
| Network Area | 20X20units |
| Node density | Average |
| Mobility | Random |
| Nodes | Static |
| Power Allocation | Random |
| Communication Standard | IEEE 802.11b |

## V.    RESULTS AND DISCUSSION

The mentioned QoS routing approach has been implemented using MATLAB 7.2. The following simulation results obtained during implementation.



Figure 4.    Overhead for DSR, NFPQR and C-NFPQR

The Fig 4 the overhead per node increases with the passage of time for DSR as the nodes drain out of power and the packet needs to be retransmitted through new route. In case of the NFPQR overhead ha reduced a lot as node failure due on the power level basis has been predicted at time of route discovery. This has been further improved using C-NFPQR.

TABLE II.    OVERHEAD VALUES

| Routing protocol | Minimum overhead | Maximum overhead |
|---|---|---|
| DSR | 204 | 223 |
| NFPQR | 203 | 206 |
| C-NFPQR | 201 | 204 |



Figure 5.    Power consumption for DSR, NFPQR and C-NFPQR

TABLE III.    POWER CONSUMPTION RANGE

| Routing protocol | Power consumption range ( mW) |
|---|---|
| DSR | 12-580 |
| NFPQR | 10-300 |
| C-NFPQR | 8-60 |

As shown in Fig 5 maximum power is consumed using DSR protocol. Power consumption has been reduced further through NFPQR. Though NFPQR does provide power optimization still lesser power consumption accounts to stable routes formed through it. C- NFPQR consumes least power due to power optimization routing



Figure 6.    Throughput plot for DSR, NFPQR and C-NFPQR

The above figure Fig 6 shows maximum throughput can be obtained through C-NFPQR protocol. For NFPQR the upper limit for throughput is much less and is least for DSR. Thus in terms of throughput also the performance of C-NFPQR supercedes other two protocols

TABLE IV.    THROUGHPUT RANGE

| Routing Protocol | Min. Throughput % | Max. Throughput % |
|---|---|---|
| DSR | 18 | 42 |
| NFPQR | 32 | 68 |
| C-NFPQR | 34 | 98 |



Figure 7. Delay performance plot for DSR, NFPQR and C-NFPQR

From Fig.7 the end to end delay offered to packet communicate from respective source to destination using DSR is almost two times the delay offered using NFPQR, and it comes out to be six times in comparison to C-NFPQR

TABLE V.    DELAY RANGE

| Routing protocol | Minimum delay(ms) | Maximum Delay(ms) |
|---|---|---|
| DSR | 54 | 132 |
| NFPQR | 30 | 70 |
| C-NFPQR | 20 | 24 |



Figure 8.    Network plot for DSR, NFPQR and C-NFPQR

From Fig 8 the network life of the ad hoc network has been improved almost ten times through C-NFPQR in comparison to DSR. Though not much significant was improvement reflected in comparison with NFPQR.

TABLE VI.    NETWORK LIFE IMPROVEMENT

| Routing protocol | Network life Improvement (%) |
|---|---|
| DSR | 1.6 |
| NFPQR | 10.5 |
| C-NFPQR | 13.8 |

## VI.    CONCLUSION

In this paper performance of C-NFPQR, NFPQR and DSR protocols for random ad hoc network has been evaluated and compared for various QoS parameters. The power level of each node and the respective geographical position is randomly defined in the network. Simulations have been run for 3 seconds considering almost no mobility of nodes during routing and communication. Better QoS routing is provided through NFPQR protocol as it provided stable routes in comparison to DSR due to node failure prediction based on power to some extent improvement in QoS but it fails to impart any power optimization facility. The C-NFPQR overcomes the limitation of both the mentioned approaches by enhancing the performance of the network in terms of evaluated QoS parameters.

REFERENCES

[1] C.S.R. Murthy and B.S. Manoj, "Ad hoc Wireless Networks: Architectures and Protocols", Prentice Hall, 2004 J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp. 68–73.

[2] D.Satynarayana, S.Sathyashree, "Node Failure Predication QoS routing Protocol for ad hoc sensor networks", 2nd International Conference on wireless communication & sensor networks, December 17-19, 2006.

[3] Tony Larsson and Nicklas Hedman, "Routing Protocols in Wireless Ad hoc Networks -A Simulation Study", Master's thesis in Computer Science and Engineering, Luleå University of Technology Stockholm, 1998.

[4] Mehran Abolhasan,Tadeusz A.Wysocki, and Eryk Dutkiewicz, "A review of routing protocols for mobile ad hoc networks " protocols for mobile ad hoc networks "

[5] David B. Johnson, David A. Maltz, Josh Broch, "DSR: The Dynamic Source Routing Protocol for Multi-Hop Wireless Ad Hoc Networks", Computer Science Department Carnegie Mellon University Pittsburgh.

[6] David B. Johnson, David A. Maltz, "Dynamic Source Routing in Ad Hoc Wireless Networks", Computer Science Department Carnegie Mellon University 5000 Forbes Avenue Pitts burgh.

[7] C.R. Lin and J.S .Liu, "QoS routing in ad hoc wireless networks", IEEE Journal on selected areas in Commuinications",17: 1426-1438, 1999.

[8] T. Bheemarjuna Reddy, I. Karthigeyan, B.S. Manoj, C. Siva Ram Murthy, "Quality of service provisioning in ad hoc wireless networks: a survey of issues and solutions", Department of Computer Science and Engineering, Indian Institute of Technology, Madras 600036, India , April 14,2004.

[9]  Vikas Kawadia, P. R. Kumar*,* "Principles and Protocols for Power Control in Wireless Ad Hoc Networks", IEEE journal on selected areas in Communications, vol. 23, no. 1, January 2005.

[10] Shigang Chen and Klara Nahrstedt, "A distributed quality of service routing in ad hoc networks", IEEE Journal on Selected Areas in Commuinications.17 (8), August 1999.

[11] V. Kawadia and P. R. Kumar, "Power control and clustering in ad hoc networks," in Proc. IEEE INFOCOM, 2003, pp. 459–469.

AUTHORS PROFILE

Sake Pothalaiah, graduated from the Department of ECE in National Institute of Technology Waragal in 2006, he obtained his M.E. degree from the department ECE, University College of Engineeering ,OU in2008. He is now working as project Associate, Deprtment of ECE, College of Engineering, JNTU . Kukatpalli

D.Srinivasa Rao, Received the B-Tech degree in Electronics & Communication Engg from Nagarjuna University, India in 1986 and M.E degree in DigitalSystems, Osmania University, India in 1994, and PhD degree in Computer Networks from University of Hyderabad, India in 2004. He is now working as Professor, Deprtment of ECE, College of Engineering, JNTU, Kukatpalli, India

# Microcontroller Based Home Automation System With Security

Inderpreet Kaur  (Asstt. Prof.)

Rayat and Bahra Institute of Engineering and Bio-technology, Mohali, India

Email: inder_preet74@yahoo.com

*Abstract*— **With advancement of technology things are becoming simpler and easier for us. Automatic systems are being preferred over manual system.  This unit talks about the basic definitions needed to understand the Project better and further defines the technical criteria to be implemented as a part of this project.**

***Keywords-component; Automation, 8051 microcontroller, LDR, LED, ADC, Relays, LCD display, Sensors, Stepper motor***

## I.  INTRODUCTION

With advancement of technology things are becoming simpler and easier for us. Automation is the use of control systems and information technologies to reduce the need for human work in the production of goods and services. In the scope of industrialization, automation is a step beyond mechanization. Whereas mechanization provided human operators with machinery to assist them with the muscular requirements of work, automation greatly decreases the need for human sensory and mental requirements as well. Automation plays an increasingly important role in the world economy and in daily experience.

Automatic systems are being preferred over manual system. Through this project we have tried to show automatic control of a house as a result of which power is saved to some extent.

## II.  HOME AUTOMATION

Home/office automation is the control of any or all electrical devices in our home or office, whether we are there or away. Home/office automation is one of the most exciting developments in technology for the home that has come along in decades. There are hundreds of products available today that allow us control over the devices automatically, either by remote control; or even by voice command.

Home automation (also called domotics) is the residential extension of "building automation". It is automation of the home, housework or household activity. Home automation may include centralized control of lighting, HVAC (heating, ventilation and air conditioning), appliances, and other systems, to provide improved convenience, comfort, energy efficiency and security. Disabled can provide increased quality of life for persons who might otherwise require caregivers or institutional care.

A home automation system integrates electrical devices in a house with each other. The techniques employed in home automation include those in building automation as well as the control of domestic activities, such as home entertainment systems, houseplant and yard watering, pet feeding, changing the ambiance "scenes" for different events (such as dinners or parties), and the use of domestic robots. Devices may be connected through a computer network to allow control by a personal computer, and may allow remote access from the internet.

Typically, a new home is outfitted for home automation during construction, due to the accessibility of the walls, outlets, and storage rooms, and the ability to make design changes specifically to accommodate certain technologies. Wireless systems are commonly installed when outfitting a pre-existing house, as they reduce wiring changes. These communicate through the existing power wiring, radio, or infrared signals with a central controller. Network sockets may be installed in every room like AC power receptacles.

Although automated homes of the future have been staple exhibits for World's Fairs and popular backgrounds in science fiction, complexity, competition between vendors, multiple incompatible standards  and the resulting expense have limited the penetration of home automation to homes of the wealthy or ambitious hobbyists.

## III.  NEED OF AUTOMATION

Earlier, we looked into the face of future when we talked about automated devices, which could do anything on instigation of a controller, but today it has become a reality.

a) An automated device can replace good amount of human working force, moreover humans are more prone to errors and in intensive conditions the probability of error increases whereas, an automated device can work with diligence, versatility and with almost zero error.

   Replacing human operators in tasks that involve hard physical or monotonous work.

   Replacing humans in tasks done in dangerous environments (i.e. fire, space, volcanoes, nuclear facilities, underwater, etc)

   Performing tasks that are beyond human capabilities of size, weight, speed, endurance, etc.

☐ Economy improvement. Automation may improve in economy of enterprises, society or most of humankind. For example, when an enterprise that has invested in automation technology recovers its investment, or when a state or country increases its income due to automation like Germany or Japan in the 20th Century.

b) This is why this project looks into construction and implementation of a system involving hardware to control a variety of electrical and electronics system.

## IV. SUPPLY UNIT

Initial stage of every electronic circuit is power supply system which provides required power to drive the whole system. The specification of power supply depends on the power requirement and this requirement is determined by its rating. The main components used in supply system are:

- transformer
- rectifier
- input filter
- regulator
- output filter
- output indication

### A. Transformer:

The main source of power(Fig 1) supply is a transformer. The maximum output power of power supply is dependent on maximum output power of transformer .We determine power from its current and voltage rating. e.g.: if there is a transformer of 12V, 500mA then maximum power delivered by transformer is 6Watt.

It means we can drive a load from this transformer up to 6w. In our project our maximum power requirement is 1watt. So to provide this power we use 12V/250mA transformer. The maximum output power of this transformer is 4watt.it means it can easily drive load up to 4 watt.

### B. Rectifier

Rectifier is a circuit which is used to convert ac to dc. Every electronic circuit requires a dc power supply for rectification. We have used four diodes.

### C. Input filter:

After rectification we obtain dc supply from ac but it is not pure dc it may have some ac ripples .To reduce these ripples we use filters. It comprises of two filters –low frequency ripple filter and high frequency ripple filter.

To reduce low frequency ripples we use electrolytic capacitor. The voltage rating of capacitor must be double from incoming dc supply. It blocks dc and passes ripples to ground.

### D. Regulator:

Regulator is a device which provides constant output voltage with varying input voltage. There are two types of regulators-

(a) Fixed voltage regulator

(b) Adjustable regulator

We have used fixed voltage regulator LM78XX last two digits signify output voltage. The voltage for our system is 5V that is why we have used 7805 regulator which provides 5V from 12V dc.

### E. Output filter:

It is used to filter out output ripple if any.

### F. Output indication

We use LED to observe the functioning of our system. If the LED glows it confirms proper functioning of our supply. We have used four power supply units.



Fig 1-volt DC supply

This supply is for the microcontroller, display and relay unit. . The microcontroller requires 5 volt supply to perform any desired task.

### G. Control Unit

Two control units were used one for internal system and one for external system and these control unit based on ATMEL'sAT89S52 microcontroller(Fig 2). The given capture shows the pins and basic requirement of microcontroller to make it functional. Detailed description of the controller is

given                                        as                                    follow:



Fig 2-Chip Board

AT89S52 is an ATMEL controller with the core of Intel MCS-51. It has same pin configuration as give above. The AT89S52 is a low-power, high-performance CMOS 8-bit microcomputer with 8K bytes of Downloadable Flash programmable and erasable read only memory and 2K bytes of EEPROM. The device is manufactured using Atmel's high density nonvolatile memory technology and is compatible with the industry standard 80C51 instruction set and pin out.

The on-chip Downloadable Flash allows the program memory to be reprogrammed in-system through an SPI serial interface or by a conventional nonvolatile memory programmer. By combining a versatile 8-bit CPU with Downloadable Flash on a monolithic chip, the Atmel AT89S52 is a powerful microcomputer which provides a highly flexible and cost effective solution to many embedded control applications.

The AT89S52 provides the following standard features: 8K bytes of Downloadable Flash, 2K bytes of EEPROM, 256 bytes of RAM, 32 I/O lines, programmable watchdog timer, two Data Pointers, three 16-bit timer/counters, a six-vector

two-level interrupt, a full duplex serial port, on-chip oscillator, and clock circuitry.

In addition, the AT89S52 is designed with static logic for operation down to zero frequency and supports two software selectable power saving modes. The Idle Mode stops the CPU while allowing the RAM, timer/counters, serial port, and interrupt system to continue functioning. The Power down Mode saves the RAM contents but freezes the oscillator, disabling all other chip functions until the next interrupt or hardware reset.

The Downloadable Flash can be changed a single byte at a time and is accessible through the SPI serial interface. Holding RESET active forces the SPI bus into a serial programming interface and allows the program memory to be written to or read from unless Lock Bit 2 has been activated.

*H. Features*

➢ Compatible with MCS-51™Products

➢ 8K bytes of In-System Reprogrammable Downloadable Flash Memory

➢ SPI Serial Interface for Program Downloading

➢ Endurance: 1,000 Write/Erase Cycles

➢ 4.0V to 5.5V Operating Range

➢ Fully Static Operation: 0 Hz to 33 MHz

➢ 256 x 8 bit Internal RAM

➢ 32 Programmable I/O Lines

➢ Three 16 bit Timer/Counters

➢ Eight Interrupt Sources

➢ Full Duplex UART Serial Channel

➢ Low Power Idle and Power Down Modes

➢ Interrupt Recovery from Power Down Mode

*I. Advantages*

➢ Less power consumption

➢ Low cost

➢ Less space required

➢ High speed

*J. Pin Description*

VCC: Supply voltage. GND: Ground., Port 0: Port 0 is an 8-bit open drain bidirectional I/O port. As an output port, each pin can sink eight TTL inputs. When 1s are written to port 0 pins, the pins can be used as high impedance inputs. Port 0 can also be configured to be the multiplexed low-order address/data bus during accesses to external program and data memory. In this mode, P0 has internal pull-ups. Port 0 also receives the code bytes during Flash programming and outputs the code bytes during program verification. External pull-ups are required during program verification.

Port 1: Port 1 is an 8-bit bidirectional I/O port with internal pull-ups. The Port 1 output buffers can sink/source four TTL inputs. When 1s are written to Port 1 pins, they are pulled high by the internal pull-ups and can be used as inputs. As inputs, Port 1 pins that are externally being pulled low will source current (IIL) because of the internal pull-ups.

In addition, P1.0 and P1.1 can be configured to be the timer/counter 2 external count input (P1.0/T2) and the timer/counter 2 trigger input (P1.1/T2EX), respectively, as shown in the following table .Port 1 also receives the low-order address bytes during Flash programming and verification.

Port 2: Port 2 is an 8-bit bidirectional I/O port with internal pull-ups. The Port 2 output buffers can sink/source four TTL inputs. When 1s are written to Port 2 pins, they are pulled high by the internal pull-ups and can be used as inputs. As inputs, Port 2 pins that are externally being pulled low will source current (IIL) because of the internal pull-ups.

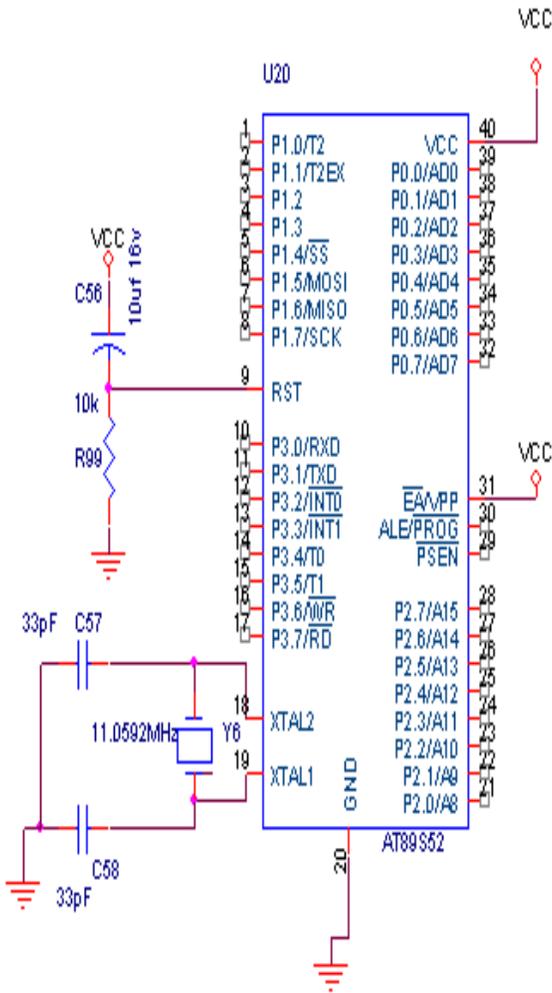Port 2 emits the high-order address byte during fetches from external program memory and during accesses to external data memory that uses 16-bit addresses (MOVX @ DPTR). In this application, Port 2 uses strong internal pull-ups when emitting 1s. During accesses to external data memory that uses 8-bit addresses (MOVX @ RI), Port 2 emits the contents of the P2 Special Function Register. Port 2 also receives the high-order address bits and some control signals during Flash programming and verification.

Port 3 Port 3 is an 8-bit bidirectional I/O port with internal pull-ups. The Port 3 output buffers can sink/source four TTL inputs. When 1s are written to Port 3 pins, they are pulled high by the internal pull-ups and can be used as inputs. As inputs, Port 3 pins that are externally being pulled low will source current (IIL) because of the pull-ups.

Port 3 receives some control signals for Flash programming and verification. Port 3 also serves the functions of various special features of the AT89S52.

RST: Reset input. A high on this pin for two machine cycles while the oscillator is running resets the device. This pin drives high for 98 oscillator periods after the Watchdog times out. The DISRTO bit in SFR AUXR (address 8EH) can be used to disable this feature. In the default state of bit DISRTO, the RESET HIGH out feature is enabled.

ALE/PROG: Address Latch Enable (ALE) is an output pulse for latching the low byte of the address during accesses to external memory. This pin is also the program pulse input (PROG) during Flash programming. In normal operation, ALE is emitted at a constant rate of 1/6 the oscillator frequency and may be used for external timing or clocking purposes. Note, however, that one ALE pulse is skipped during each access to external data memory.

If desired, ALE operation can be disabled by setting bit 0 of SFR location 8EH. With the bit set, ALE is active only during a MOVX or MOVC instruction. Otherwise, the pin is weakly pulled high. Setting the ALE-disable bit has no effect if the microcontroller is in external execution mode.

PSEN: Program Store Enable (PSEN) is the read strobe to external program memory. When the AT89S52 is executing code from external program memory, PSEN is activated twice each machine cycle, except that two PSEN activations are skipped during each access to external data memory.

EA/VPP: External Access Enable. EA must be strapped to GND in order to enable the device to fetch code from external program memory locations starting at 0000H up to FFFFH. Note, however, that if lock bit 1 is programmed, EA will be internally latched on reset. EA should be strapped to VCC for internal program executions.

This pin also receives the 12-volt programming enable voltage (VPP) during Flash programming.

XTAL1: Input to the inverting oscillator amplifier and input to the internal clock operating circuit.

XTAL2: Output from the inverting oscillator amplifier.

H. Display Unit



Fig 3-Display Unit

Liquid crystal displays (LCD) is an alphanumeric display and widely used in recent years as compared to LEDs. This is due to the declining prices of LCD, the ability to display numbers, characters and graphics, incorporation of a refreshing controller into the LCD, their by relieving the CPU of the task of refreshing the LCD and also the ease of programming for characters and graphics. We have used JHD162A advanced version of HD44780 based LCDs.

## V. WHAT CAN BE AUTOMATED

Virtually anything in the home/office that is powered by electricity can be automated and/or controlled. We can control our electrical devices. We can turn our porch lights on

automatically(Fig 4) at dark or when someone approaches and can see

| LCD Display | TEMPRATURE MONITORING | FIRE SENSOR |
|---|---|---|

| Led Light | LDR |
|---|---|

8051 MICROCNTROLLER

| Led Light | LDR |
|---|---|

PERSON COUNTING

RELAY DRIVE UINIT

| DEVICE 1 | DEVICE 2 |
|---|---|

Fig 4-Automatic System

who is at the front door from any nearby television, and talk to them or unlock the door from any nearby telephone. Have the security system turn off lights, close drapes and setback the temperature when we leave and turn on the alarm system. The possibilities are only limited by our imagination.

## VI. FEATURES

### A. *Password Based Locking System*

In this system we have ensured a safe locking system. On seeing from outside the lock would not be visible but this inbuilt locking system ensures security. This lock can be opened and closed with the help of a password which we will give using a keypad. The door will only open or close only if the password is correct else it will remain in its original state. The lock cannot be broken because to the person standing outside can just see the closed door and not the lock as it is inbuilt. The password is given with the help of controller and can be changed by simply making a small change in the program and then burning the program in the controller.

### B. *Counter dependent automatic switching system of room*

After opening the lock when the person enters the room the counter gets incremented. Now if it is a day then the lights would not be switched on but if it is dark then the lights will automatically switch on. Now whatever may be the number of people entering the room the counter will automatically get incremented by itself and on leaving the room the counter will get decremented but the system will keep on working .Once the counter is zero in other words once everyone leaves the room the switching system will automatically stop working.

### C. *Temperature controlled cooling system*

Once the person has entered the room he would not require to switch on anything everything will just happen automatically. Like if the temperature is high then the fan will switch on, on its own. Else it will remain in off state. This temperature is predefined by us in the controller. But this system will only work if there is a person in the room in other words if the counter is not zero.

### D. *Light saving system*

This light saving system is used in two places for internal section and external section. If a person is not at home or sitting inside the room and it is dark outside then the lights will automatically get switched on and when its day the light will get switched off. This ensures power saving.

### E. *Fire and Smoke sensor*

This part detects any fire or smoke from a fire and set an alarm or an indication.

## VII. CONCLUSION

An automated home can be a very simple grouping of controls, or it can be heavily automated where any appliance that is plugged into electrical power is remotely controlled. Costs mainly include equipment, components, furniture, and custom installation.

Ongoing costs include electricity to run the control systems, maintenance costs for the control and networking systems, including troubleshooting, and eventual cost of upgrading as standards change. Increased complexity may also increase maintenance costs for networked devices.

Learning to use a complex system effectively may take significant time and training.

Control system security may be difficult and costly to maintain, especially if the control system extends beyond the home, for instance by wireless or by connection to the internet or other networks.

Future of Automation: Future will be of Automation of all products. Each and every product will be smart devices that we use daily and that will be controlled through a smart chip called microcontrollers. Each and Every home appliances will be controlled either by PC or hand held devices like PDA or mobile handsets. Some examples of it are when you want you can switch on/off Fan of your home by mobile handset or PC.

Smart Grid: Home automation technologies are viewed as integral additions to the Smart grid. The ability to control lighting, appliances, HVAC as well as Smart applications (load shedding, demand response, real-time power usage and price reporting) will become vital as Smart Grid initiatives are rolled out.

REFERENCES

[1] http://www.smartcomputing.com/editorial/article.asp?article=articles%2F1995%2Fmar95%2Fpcn0323%2Fpcn0323.asp retrieved 2010 09 02

[2] "U.S. Patent 613809: Method of and apparatus for controlling mechanism of moving vessels and vehicles". United States Patent and Trademark Office. 1898-11-08. Retrieved 2010-06-16.

[3] William C. Mann (ed.) Smart technology for aging, disability and independence : the state of the science, John Wiley and Sons, 2005 0-471-69694-3, pp. 34-66

[4] http://www.drdobbs.com/184404040;jsessionid=IM5NJPJYWXAOFQE 1GHPCKH4ATMY32JVN Dag Spicer, If You Can't Stand the Coding, Stay Out of the Kitchen, Dr. Dobb's Journal, August 2000 , retrieved 2010 Sept 2

[5] "Home automation costs". Totalavcontrol.co.uk. Retrieved 2010-02-18.

[6] "About Us". InsteonSmartGrid.com. Retrieved 2009-11-20.

[7] Worlds Open Protocol ISO/IEC 14543-3 KNX www.knx.org

AUTHORS PROFILE

Inderpreet Kaur received the Master of Engineering from PEC University of Technology (Formerly Punjab Engineering College), Chandigarh. She is pursuing her Ph.D in Electronics and Communication Engineering. Her area of research is optical fiber communication. She has total of 14 years of experience. She is life member of ISTE, OSI, IEI. She is in the committee of Reviewer in national and international Journals. She regularly contributes in various Journals, Magazines, and Conferences. She can be contacted at inder_preet74@yahoo.com

# Characterization and Architecture of Component Based Models

Er.Iqbaldeep kaur *(Author)*
Assistant Professor ,Department of
Computer Science and Engineering,
Rayat and Bahra Institute of
Engineering and Bio-Technology
,Kharar, India
(eriqbaldeepkaur@gmail.com)

Dr.P.K.Suri
Dean and Professor, Computer
Science and Application Department,
Kurukshetra University,
Kurukshetra, India
( pksuritf25@yahoo.com)

Er.Amit Verma
Assistant Professor , Department of
Electronics and Communication
Engineering, Rayat and Bahra
Institute of  Engineering and Bio-
Technology,Kharar,India
(ervermaamit@gmail.com)

*Abstract*—**Component based Software Engineering is the most common term nowadays in the field of software development. The CBSE approach is actually based on the principle of 'Select and Use' rather than 'Design and Test' as in traditional software development methods. Since this trend of using and 'reusing' components is in its developing stage, there are many advantages and problems as well that occur while use of components. Here is presented a series of papers that cover various important and integral issues in the field concerned. This paper is an introductory research on the essential concepts, principles and steps that underlie the available commercialized models in CBD. This research work has a scope extending to Component retrieval in repositories and their management and implementing the results verification**.

*Keywords- Components, CBSD, CORBA, KOAYLA, EJB, Component retrieval, repositories etc.*

## I.  INTRODUCTION

The advantages of component based development include lesser development time, lower costs, reusability and better modification. A component is the basic building block of an application or system created with CBD. Generally, a component can be defined as an independent and replaceable part of a system that fulfills a clear function. It works in the context of a well defined architecture and can communicate with other components through its interfaces (Fig. 1). Although the basic principle of 'Plug and play' is very promising, but it also brings in some practical difficulties faced by the stakeholders involved. For instance, when we buy a component, we do not know exactly about its maintenance, the security arrangements and the most important its behavior when integrated with other components. There exist some models in the market that, to an extent, provide us with some standards and interfaces to aid the intercommunication process of components within integration. The models enable the independently designed components to be deployed and ease the communication between them. Rightly stated, it can be said that a component model supports components by forcing them to conform to certain standards and allows instances of these components to cooperate with other components in this

model (fig. 2). In the absence of component models, there would be obvious non-cooperation among independently developed components, so the aim of 'independent deployment and assembled integration' of components



Fig. 1 Component and its features

Fig. 2 Component Models provide interface to components

would not be realized. Thus, these models play a significant role in making the real goal[26-28] of CBD achieved. In the next sections, a detailed characteristic listing has been done for the main component models in market.

## II. EXISTING COMPONENT MODEL(BACKGROUND & RELATED WORK)

The cornerstone of any CBD methodology is its underlying component model which defines what components are, how they can be constructed, how they can be assembled[1]. Component-based approach has shown considerable successes in recent years in many application domains like Distributed and web-based systems, desktop and graphical applications etc. In these domains the general-purpose component technologies, such as COM, .NET, EJB, J2EE are used [12]. According to [5], there are some commercial players involved in the software component revolution, such as BEA, Microsoft, IBM and Sun. [5]  also states that among the component infrastructure technologies that have been developed, three have become somewhat standardized: OMG's CORBA, Microsoft's Component Object Model (COM) and Distributed COM (DCOM), and Sun's JavaBeans and Enterprise JavaBeans .

Most of the literature contains description about three major component models viz, OMG's CORBA, SUN's EJB and Microsoft's COM. The present work includes these three and some other less known models that are still maturing. At present there are various component models that are being used. These are shown pictorially in figure 3. Some approaches, such as Visual Basic Controls (VBX), ActiveX controls, class libraries, and JavaBeans, make it possible for their related languages, such as Visual Basic, C++, Java and the supporting tools to share and distribute application pieces. But all of these approaches rely on certain underlying services to provide the communication and coordination necessary for the application. The infrastructure of components, called a component model, in fact, acts as the "plumbing" that allows communication among components [9].



Fig.3 Component Models in Market

Generally Component Models work in three different service categories as follows: Basic, Distributed & Enterprise For example, the basic services include the simple component model version like COM, CORBA or EJB. Similarly, Distribution is provided[32-33] with a communication protocol that has been added to the basic component model.

## III. COMPONENT OBJECT MODEL

It provides platform-dependent, based on Windows and Windows NT, and language-independent component based applications. COM defines how components and their clients interact. This interaction is defined such that the client and the component can connect without the need of any intermediate system component. Specially, COM provides a binary standard that components and their clients must follow to ensure dynamic interoperability. This enables on-line software update and cross-language software reuse [7].The following features characterize COM model:

• A model for designing components that have multiple interfaces with dynamic binding

• COM is an open standard, with main platform as Microsoft Windows

• Interfaces are the only means for components to expose themselves

• The interfaces are binary which provide the obvious ease to implement the component in multiple programming languages such as C++, Visual Basic and Java.

• A COM component can implement and expose multiple interfaces

- COM helps client to locate server components and desired interfaces by establishing connection between client and server.
- Interfaces [35-36] are defined as unchangeable units (A basic COM rule is that one cannot change an interface when it has been released), hence solving the interface versioning problem. Each time a new version of the interface is created a new interface will be added instead of changing the older version.

DCOM is the protocol that is used to make COM location transparent. A client talks to a proxy, which looks like the server and manages the real communication with the server. [3] has stated on DCOM , the extension of the Component Object Model (COM) as follows. Distributed COM (DCOM) is a protocol that enables software components to communicate directly over a network in a reliable, secure, and efficient manner.

DCOM is designed for use across multiple network transports, including Internet protocols such as HTTP. When a client and its component reside on different machines, DCOM simply replaces the local inter-process communication with a network protocol. Neither the client nor the component is aware the changes of the physical connections. COM+ is an extension to COM with technologies that supports various additional services like transactions, directory service, load balancing and message queuing. COM+ is implemented to connect the clients to the business logic, through an Internet Information Server (IIS) or DCOM, as shown in figure 4.

The business logic uses ActiveX Data Objects (ADOs) to access the data in the databases.

## IV. ENTERPRISE JAVA BEANS(EJB)

U In accordance with [3], Java platform offers an efficient solution to the portability and security problems through the use of portable Java byte codes and the concept of trusted and non-trusted Java applets. Java provides a universal integration and enabling technology for enterprise application development, which includes:

➢ Interoperating across multivendor servers;
➢ Propagating transaction and security contexts;
➢ Servicing multilingual clients; and
➢ d)Supporting ActiveX via DCOM/CORBA bridges.

[8] has mentioned that the JavaBeans component architecture supports applications of multiple platforms, as well as reusable[24],[35], client-side and server-side components. JavaBeans and EJB extend all native strengths of Java including portability and security into the area of component-based development. The portability, security, and reliability of Java are well suited for developing robust server objects independent of operating systems, Web servers and database management servers. Sun's Java-based component model consists of two parts:

- JavaBeans for client-side component development
- Enterprise JavaBeans (EJB) for the server-side component development.

The following are the main features of EJB: EJB is part of the Java 2 Platform Enterprise Edition (J2EE) which includes remote method invocation (RMI), naming and directory interface (JNDI), database connectivity (JDBC), Server Pages (JSPs) and Messaging services (JMS). Fig. 5 shows the architectural style of EJB used in a three-tier application. EJB is designed so it can run together



Fig.4 COM Architecture



Fig 5 EJB Architecture

with CORBA and access CORBA objects easily.

## V. COMMON OBJECT REQUEST BROKER ARCHITECURE (CORBA)(CURRENT TECHNOLOGY USED)

The Common Object[25],[29] Request Broker Architecture[30],[34] (CORBA) is a standard that has been developed by the Object Management Group (OMG) in early nineties. The OMG provides industry guidelines and object management[31]specifications to supply a common framework for integrating application development. Primary requirements for these specifications are reusability[21-24],[35] portability and interoperability of object based software components in a distributed environment. CORBA is part of the Object Management Architecture (OMA)[31] which covers object services, common facilities and definitions of terms. Object services include naming, persistency, events, transactions and relationships.

The following are the primary working principle of OMG's CORBA:

• The most important part of a CORBA system is the Object Request Broker (ORB).

• An object request broker (ORB) provides the basic mechanism for transparently

• Requests can be made through the ORB without regard to the service location or implementation.

• Objects publish their interfaces using the Interface Definition Language (IDL)

• Objects are stored in an interface repository where they can be found and activated on demand from the clients.

• The stubs and proxies are generated from the IDL specification

According to [3], CORBA manages details of component interoperability. Also CORBA is widely used in Object-Oriented distributed systems[6].

## VI. LESS POPULAR COMPONENT MODEL TECHNOLOGIES SOFA SOFA 2 KOALA KOBRA

As stated by [13], the component model SOFA is a part of SOFA project (Software Appliances). It is a software system is described as a hierarchical composition of primitive and composite components. A component is an instance of a template, which is described by its frame and architecture. The frame is a "black-box" specification view of the component defining its provided and required interfaces. Primitive components are directly implemented by described software system they have a primitive architecture[37]. The architecture of a composed component is a "grey-box" implementation view, which defines first level of nesting in the component. It describes direct subcomponents and their interconnections via interfaces. The connections of the interfaces can be realized[38] via connectors, implicitly for simply connections or explicitly. Explicit connectors are described in a similar

way as the components, by a frame and architecture. The connector frame is a set of roles, i.e. interfaces, which are compatible with interfaces of components.

SOFA 2 is a component system employing hierarchically composed components. It is a direct successor of the SOFA component model.

### KOALA

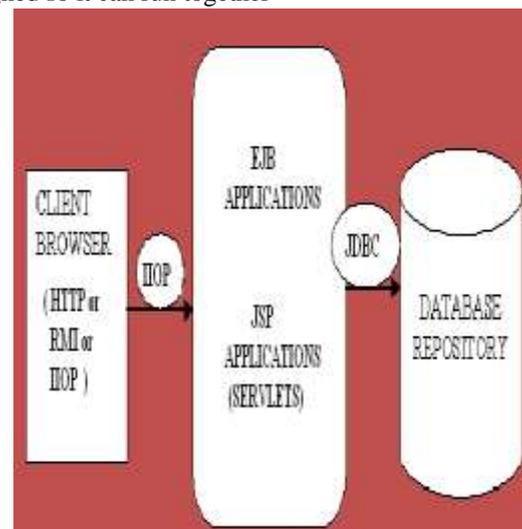Having most of its uses within Philips, Koala [14] offers explicit management of a special graphical notation that is very helpful in design discussions, and an elegant parameterization mechanism. Its partial evaluation techniques can calculate part of the configuration at compile time while generating code for that part that must be determined at runtime. In designing Koala, a strict separation is sought between component and configuration development.

• Koala components are units of design, development, and – more importantly – reuse.

• As in COM and Java, a Koala interface is a small set of semantically related functions.

• Koala components access all external functionality through requires interfaces which provides the architects with a clear view of the of the system's resources use.

• Koala components are designed independently of each other. They have interfaces to connect to other components, but this binding is late – at configuration time.

Koala has some extra features that are aimed at handling diversity efficiently: interface compatibility, function binding, partial evaluation, diversity interfaces, diversity spreadsheets, switches, optional interfaces, and Connected interfaces.

### KOBRA

[16] states that KobrA is a UML-based method for describing components and component-based systems developed at the Fraunhofer Institute for Experimental Software Engineering at the beginning of the decade. The acronym stands for the term "Komponenten basierte Anwendungsentwicklung" – German for "Component-based Application Development". KobrA has been successfully used by a number of companies in industrial settings and has given rise to numerous specializations and offshoots . The original version of the method was developed for the UML 1.x flavor of the UML.

## VII. CONCLUSION AND COMPARISON

Component-based systems result from adopting a component-based design with strategy, and software component technology includes the products and concepts that support this design strategy. By design strategy we mean something almost near to architectural style—a high-level design pattern and system described by the types of components in a system and their patterns of interaction [20]. Component based software development (CBSD) refers to the development of software component systems making considerable use of software components. Component based software development can help the software industry to realize

productivity and quality gains similar to those achieved in hardware and manufacturing organizations. A detailed characterization of known component model technologies has been done in the present research work. The difference in all the model with respect to properties as shown in table 1 is illustrated. Some models like COM, CORBA and EJB are very well known among users and developers, whereas some other quite effective model technologies for component

Table 1: Comparative Study

| Property | CORBA | EJB | COM/DCOM | KOALA | SOFA |
|---|---|---|---|---|---|
| ENVIRONMENT | Not well developed | new | Widely supported | New | new |
| BINARY INTERFACE | ---------- | Java based | Key to work | Java based | --------- |
| COMPATIBILITY AND PORTABILITY | Strong in binding poor in portability | Portable but not compatible | Not having source level concept | Portable but less compatible | Less portable |
| MODIFICATION AND MAINTENANCE | Hard to do the same | Easier in the similar task | Need extra modification | -------- | Need Modification |
| SERVICES PROVIDED | Poor in implementation | Nor implemented | Supplemented by number of services | Supplements by number of services | ------- ----- |
| PLATFORM DEPENDENCY | independent | independent | Dependent | Independent | dependent |
| LANGUAGE DEPENDENCY | Independent | Dependent | Independent | Independent | dependent |
| APPLICATION | Traditional computing | Used in Web client | Traditional desktop application | diversity spreadsheets, switches | For black-box specification |

Based software development are less popular as compared to these. Since the CBSE is a new discipline and is still maturing, a lot has to be done to find solutions to its associated problems which remain unsolved.

REFERENCES

[1]  Kung –Kiu and Zheng Wang, "A survey of Software Component Models", School of computer Science, University of Manchester, April, 2005, available at http://www.cs.man.ac.uk/preprints/index.htm

[2]  A. Campbell. A Quality of Service Architecture. PhD Thesis, Lancaster University,1996

[3]  Component-Based Software Engineering: Technologies, Development Frameworks, and Quality Assurance Schemes, Xia Cai, Michael R. Lyu, Kam-Fai Wong Roy Ko, The Chinese University of Hong Kong Hong Kong Productivity Council {xcai@cse, lyu@cse, kfwong@se}.cuhk.edu.hk roy@hkpc.org

[4]  http://www.omg.org/corba/whatiscorba.html, Mar, 2000.

[5]  W. Kozaczynski, G. Booch, "Component-Based Software Engineering," IEEE Software Volume: 155, Sept.-Oct. 1998, pp. 34–36.

[6]  S.S.Yau, B. Xia, "Object-Oriented Distributed Component Software Development based on CORBA," Proceedings of COMPSAC'98. The Twenty-Second Annual International, 1998, pp. 246-251

[7]  Y.M. Wang, O.P. Damani, W.J. Lee, "Reliability and Availability Issues in Distributed Component Ojbect Model (DCOM)," Fourth International Workshop on Community Networking Proceedings, 1997, pp. 59 –63.

[8]  http://developer.java.sun.com/developer, Mar. 2000

[9]  A.W.Brown, K.C. Wallnau, "The Current State of CBSE,"IEEE Software, Volume: 15 , Sept.-Oct. 1998, pp. 37- 46

[10] C.Szyperski, "Component Software: Beyond Object-Oriented Programming," Addison-Wesley, New York, 1998.

[11] G. Pour, "Enterprise JavaBeans, JavaBeans & XML Expanding the Possibilities for Web-Based Enterprise Application Development," Proceedings Technology of Object-Oriented Languages and Systems, 1999, TOOLS 31, pp.282-291.

[12] Component-based Development Process and Component Lifecycle by Ivica Crnkovic, Stig Larsson, Michel Chaudron

[13] 'Component Model with Support of Mobile Architectures' by Marek Rychllli, Brno University of Tcchnology, Czech Republic

[14] The Koala Component Model for Consumer Electronics Software, Rob van Ommering, Frank van der Linden, Jeff Kramer, Jeff Magee, IEEE Computer, March 2000, p78-85

[15] Enterprise Distributed Object Computing Conference, 2001. EDOC '01. Proceedings. Fifth IEEE International Publication Date: 200, Pages 212 - 223 , Seattle, WA.

[16] 'Modeling Components and Component-Based Systems in KobrA' by Colin Atkinson

[17] Enterprise JavaBeans Specification. Version 2.0. Sun Microsystems. 2001.

[18] Object Management Group. The Common Object Request Broker: Architecture and Specification. version 3.0./02-06-33. 2002.

[19] F. E. Redmond III. DCOM: Microsoft Distributed Component Object Model. [sofa] Frantisek Plasil, Dusan Balek, and Radovan Janecek. SOFA/DCUP: Architecture for Component Trading and Dynamic Updating. Proceedings of ICCDS 98, May 4-6, 1998, Annapolis, Maryland, USA. IEEE CS Press. 1998.

[20] Bass, L; Clements, P.; & Kazman R. Software Architecture in Practice. Boston, Ma.: Addison Wesley, March 1998.

[21] zyperski, C, Component Software: Beyond Object-Oriented Programming, Addison Wesley, 1999.

[22] Cecilia Albert and Lisa Brownsword, Evolutionary Process for Integrating COTS-Based Systems (EPIC): An overview, Technical Report CMU/SEI-2002-TR-009 ESC-TR-2002-009, July, 2002.

[23] Jerry Zeyu Gao, Jacob Tsao, Ye Wu, Testing and Quality Assurance for Component Based Software, Artech House Publishers, 2003.

[24] David Garlan et al, "Architecural Mismatch: Why Reuse is so Hard", IEEE software, 1995.

[25] Ian graham, Object Oriented Methods, - Principles and practice, 3rd Edition, Addison Wesley, Object Technology Series.

[26] Ian Sommervilee, Software Engineering, 7th Edition, Pearson Education.

[27] R.S.Pressman, Software Engineering – A Practioners Approach, Fourth Edition, McGraw Hill International Series.

[28] Hafedh Mili et al, Reuse Based Software Engineering, Techniques, organization and Controls, John Wiley and Sons, 2002.

[29] Alencar, A. & Goguen, J. "OOZE," Stepney, S.; Barden, R.; & Cooper, D., ed. Object Orientation in Z, Workshops in Computing. Los Angeles, Ca.: Springer-Verlag, 1992.

[30] Allen, R.; Douence, R.; & Garlan, D. "Specifying Dynamism in Software Architectures," Proceedings of the 1st Workshop on Component-Based Systems. Zurich, Switzerland, 1997, in con-junction with European Software Engineering Conference (ESEC) and ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE), 1997

[31] Baggiolini, V. & Harms, J. "Toward Automatic, Run-Time Fault Management for Component-Based Applications," Proceedings of the 2nd International Workshop on Component-Oriented Pro-gramming (WCOP97), in conjunction with the European Confer-ence on Object-Oriented Programming (ECOOP98, Brussels, Belgium, July 1998.

[32] Barnes, J. High Integrity Ada: the SPARK Approach. Boston, Ma.: Addison-Wesley, 1997

[33] Box, D. Essential COM. Boston, Ma.: Addison-Wesley, 1998.

[34] Ciancarani, P. & Cimato, S. "Specifying Component-Based Soft-ware Architectures," 60–70. Proceedings of the ESEC/FSE-Workshop on Foundations of Component-Based Systems (FoCBS), Zürich, Sep. 1997.

[35] Deline, R. "Avoiding Packaging Mismatch with Flexible Packaging," Proceedings of the 21st International Conference on Soft-ware Engineering. Los Angeles, Ca., May 1999.

[36] Della, C.; Cicalese, T.; & Rotenstreich, S. "Behavioral Specification of Distributed Software Component Interfaces." IEEE Computer (Jul. 1998): 46-53

[37] Dowson, M. "ISTAR and the Contractual Approach," 287-288. Proceedings of the 9th International Conference on Software En-gineering. Monterey, Ca, March 30-April 2, 1987. Washington DC, Baltimore, Md.: IEEE Computer Society and the Association for Computing Machinery, April 1987.

[38] Hissam, S & Carney, D. "Isolating Faults in Complex COTS-based Systems." Journal of Software Maintenance: Research and Practice, No. 11 (1999): 183-1999

# Performance Improvement by Changing Modulation Methods for Software Defined Radios

Bhalchandra B. Godbole
Karmaveer Bhaurao Patil College of Engineering and
Polytechnic, Satara - 415001
bbgodbole@rediffmail.com

Dilip S. Aldar
Karmaveer Bhaurao Patil College of Engineering and
Polytechnic, Satara - 415001
dilip_aldar@rediffmail.com

*Abstract-* **This paper describes an automatic switching of modulation method to reconfigure transceivers of Software Defined Radio (SDR) based wireless communication system. The programmable architecture of Software Radio promotes a flexible implementation of modulation methods. This flexibility also translates into adaptively, which is used here to optimize the throughput of a wireless network, operating under varying channel conditions.**
**It is robust and efficient with processing time overhead that still allows the SDR to maintain its real-time operating objectives. This technique is studied for digital wireless communication systems.**
**Tests and simulations using an AWGN channel show that the SNR threshold is 5dB for the case study.**

*Keywords- Wireless mobile communication, SDR, reconfigurability, modulation switching.*

## I. INTRODUCTION

Reconfiguration, dynamic/static, partial/complete is an essential part of software radio technology. Thanks to SDR, so that the systems can be designed for change and evolution. In other words "*change*" becomes part of mainstream system operation. Recent work in European, Asian, and American R&D projects and the SDR Forum, lately WWRF, clearly shows that the concept of reconfigurability especially in the context of mobile cellular networks is not very complicated business. Reconfiguration still raises question on required *system-level support* both at the reconfigured devices and the network side [1], [2].

Over the past decade, previous work has correctly demonstrated the technical feasibility of the SDR approach in the design of radio equipment. This potential can be concretely exploited through equipment reconfiguration. In the evolution path towards 4G and beyond, this potential can be useful in many technically challenging as well as commercially attractive scenarios. Optimization in the QoS will necessitate considering reconfiguration as part of the mainstream operation [3].

This paper describes one application that exploits the flexibility of a software radio. The ability to select automatically, the correct modulation scheme used in an unknown received signal is a major advantage in a wireless network. As a channel capacity varies, modulation scheme switching enables the baud rate to be increased or decreased thus maximizing channel capacity usage. However the finite processor computing power limits the complexity of software radio if real time constraints are to be met [4], [5]. This paper proposes an automatic modulation scheme-switching algorithm based on SDR for *General-Purpose Processor* (GPP), *Digital Signal Processing* (DSP) and /or VLSI chips like FPGA. In this paper, we have simulated the algorithm on GPP.

The GPP approach makes use of P-IV system for rapid application development, large amount of program memory, and relatively inexpensive compared to the inflexible dedicated hardware. Conventional DSP processors rely on assembly language optimization for maximization of application efficiency, but with the P-IV, the degree of optimization possible using high-level languages such as C++ and MATLAB is much greater [6]. Software Defined Radio system closely tracks the advances in new high-speed processor technology allowing the addition of more complex signal processing techniques to the SDR system, while maintaining the real-time objective.

While switching the modulation method, it should be ensured that both transmitter and receiver are in tandem. To do so, any of two existing techniques can be used:

1. Automatic Modulation Recognition Technique,

2. Cross Layer Adaptation Technique.

Work in automatic modulation recognition is interesting and extensive work has been carried out for a number of years producing processor intensive techniques mainly restricted to non real-time operation. Recently published modulation recognition algorithms include a decision theoretic approach and pattern recognition approach used to discriminate between digitally modulated signals [7]. Modulation recognition scheme by the signal envelope extraction method and a digital modulation scheme classifier based on a pattern recognition technique generalises the moment matrix technique to gray scale images used in binary image word-spotting problems. A successful algorithm for an SDR implementation must be robust and efficient but the processing overhead must not stop the software radio from maintaining its real-time operating objectives. As the techniques described above are in general very processor intensive, they are not presently suitable for the

SDR systems [8]. Here, we are using the physical and MAC layer signaling and control channel for the transmitter and receiver handshaking to signal the switched modulation. But it is also possible to use modulation recognition instead of cross layer signaling and use synchronous communication without adding delays.

With an attempt to introduce briefly Software Defined Radios in Section I, we explain SDR terminal hardware and software architecture in Section II. Section III discusses Cross Layer Adaptation technique to handshake the modulation method switching. Section IV is on the Case Study in details with the algorithm for modulation switching and the simulation results are presented in Section V. The work is concluded in Section VI.

II. SDR TERMINAL: ARCHITECTURE AND CONFIGURATION

*A. Hardware Architecture:*

The future mobile terminals are designed to support high bit traffic and they will be certainly equipped with high computational capabilities [4], [9]. Fourth generation mobile terminals, as shown in figure-1, will use

more than one processor to increase the global computational capacity of the device. The Operating System (OS) performs all the control functions. In a SDR mobile terminal, the OS can also be used to generate the firmware to be inserted into the reconfigurable strata of the hardware.

A feasible architecture of baseband digital signal processing section is constituted by the following blocks:

- A general purpose microprocessor;

- Multiple DSP/FPGA blocks, depending on functional requirements of device;

- FIFO memories for both the RX and TX branches;

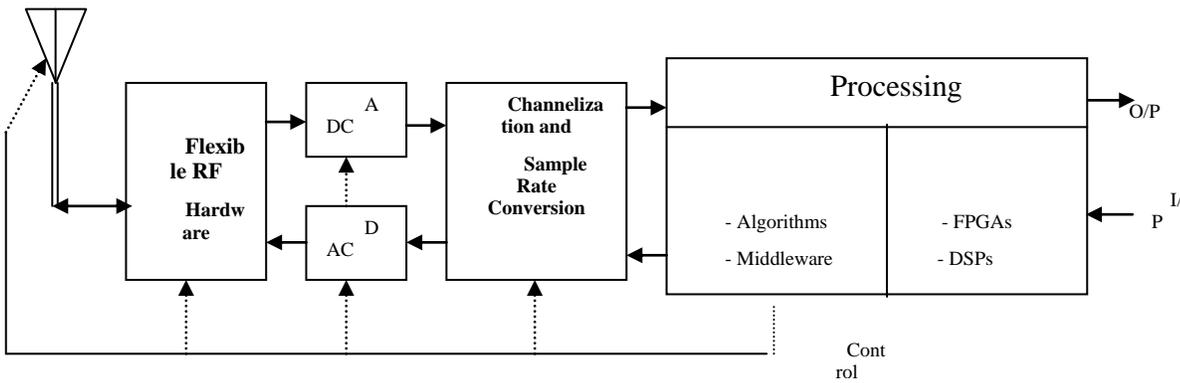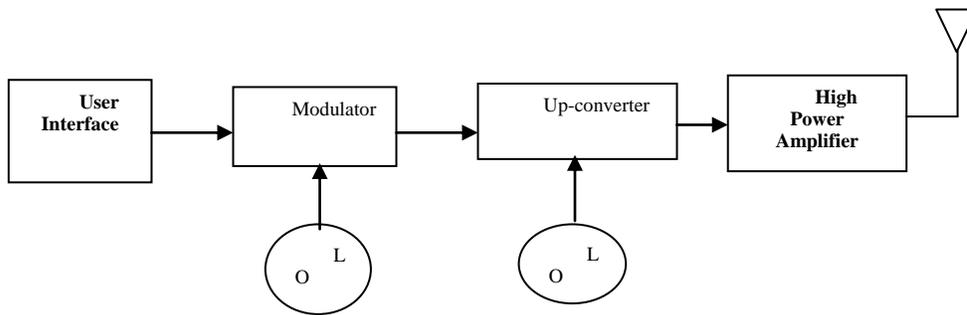- ROM/FLASH memories, for storing the resident parts of the OS;

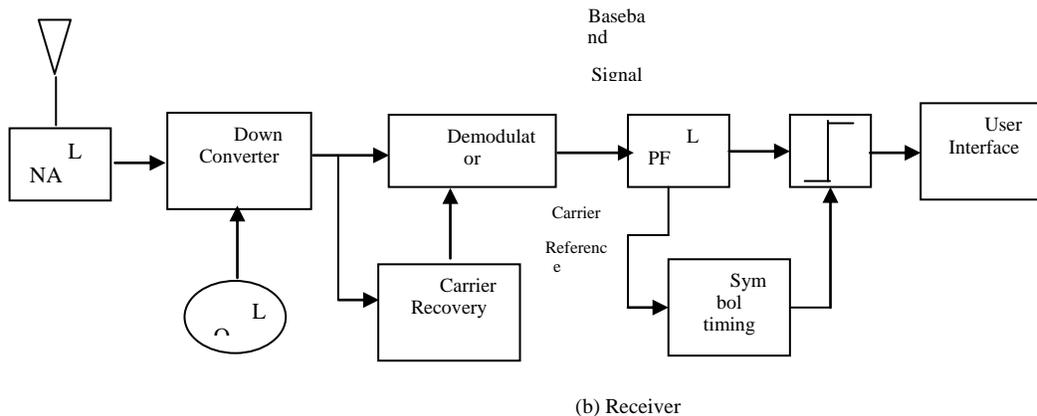Figure 1. Model of a Software Radio

(a)Transmitter

(b) Receiver

Figure-2. General Architecture of a transceiver in a software radio based wireless communication system.

The reconfigurable hardware, the FIFO Memories and the general purpose CPU are interconnected by high-speed bus. Alternately, the transceiver in a software radio based wireless communication system for GPP hardware will be as depicted in figure-2.

### B. Software Architecture

The device operating in SDR technology acts as a node in the hosting networks architecture both for control functions (control plane) and for the communication functions (traffic plane). The relevant entities of SDR network architecture are the configuration manager and the SDR mobile terminal. The physical layer of the mobile terminal is split into two sub-layers: the reconfigurable hardware and the macrocode (firmware), which implements the target communication protocol.

The key process for remote terminal configuration is the radio software download. The download process is constituted by phases like, pre-download, download and post-download [10]. The radio download data is intended platform dependent. In [11] the remote physical layer definition is operated through a high level descriptive language, called RADL (Radio Access Definition Language). As an example, functions like digital filtering, coding and decoding, modulation, pulse shaping, carrier recovery, timing acquisition are defined by a set of parameters provided by the SDR terminal using its libraries.

We simulate these functions using MATLAB and perform the sequence of operations as per requirement.

### III. CROSS LAYER APPROACH

Any innovation adopted at physical layer will not obtain the maximum performance if it does not interact with the upper MAC layer. This concept, opposed to the historical separation of functions between layers, has demonstrated significant results in terms of efficiency in the use of the communication resources. Cross-layer design, whichever is its target, requires a physical layer in communication with MAC not only giving information about the physical status of the channel, but also accepting real-time reconfiguration of

transmission parameters. Also in this case, the SDR acts as the enabling technology, providing a full logical control over the signal processing functions located at physical layer.

The so-called cross-layer approach, depicted in figure-3, is the mutual exchange of information between physical & upper layers and is naturally enabled by Software Radio [12]. Several objects were proposed for cross-layer design. Some research has addressed improving QoS improvement and power management, proposing strategies based on modifying different layers of the communication system and also show that adapting transmission methodology to channel fading significantly improves link efficiency [13].

Therefore, assuming that all second and third generation systems can be simulated on the reconfigurable multi-processor platform, we propose modulation-switching algorithm to improve performance.

### IV. CASE STUDY

In order to achieve high data rate transmission under a target BER, to attain so-called Quality of Service (QoS), we propose modulation-switching algorithm to change the number of bits for each symbol. Computing receiver noise and SNR at a receiver is important for determining coverage and QoS in a wireless communication system. SNR determines the link quality and impacts the probability of error in a wireless communication system. Thus, the ability to estimate SNR is important for determining suitable transmitter powers or received signal levels in various propagation conditions. SNR at a receiver is a convenient metric that allows a designer to factor in the noise induced by the channel.

Modulation in channel has a significant effect on the quality of the information transfer measured in BER and on the complexity of the receiver. Receiver complexity dominates the complexity of SDR.

A receiver is typically four times more complex than a transmitter in terms of MIPS required to implement the baseband and the IF processing in the software. The digital modulation_demodulation algorithm topics include AGC,

Channel waveform coherence, coding/decoding and spreading /dispreading of the spectrum used [14]. In the present case study, we will be focusing on the modulation switching, therefore, AGC, coding and spreading are not discussed.

As shown in Figure-4, the probability of bit error is a function of channel modulation. There are broadly three main digital modulation types: ASK, FSK and PSK. To approach channel capacity upper bound, adaptive modulation with continuous modulation order have to be used. However, in reality, the continuous modulation order is very complex, therefore, discrete modulation orders, M = 2, 4, 16, 64 are used here, where M = 2 is the BPSK modulation, M = 4 is the

QPSK modulation, and 16, 64 are M-QAM modulations. When the full CSI is known at the transmitter side, it is possible to estimate the channel quality by SNRs. A channel with better quality can be assigned a larger number of bits and a higher order modulation, whereas a channel with poorer quality has to be assigned fewer bits or no bit when the channel quality is too bad to transmit even BPSK modulated signals.



Figure-3 The cross-layer concept

1. Speech Coder =47.2 %

2. Channel Encoder, interleaver, cipher = 2.7 %

3. Speech Decoder= 16.8 %

4. Decipher & deinterleaver =0.73 %

5. Channel Decoder =12.0 %

6. Modulator = 13.6 %



Figure-4 CPU requirement for EIGHT logical channels on Pentium-IV: 2 GHz [14].

For a certain modulation type, the relationship of modulation order M, SNR and baud, $P_b$ performance can be expressed as:

$$P_b = f(SNR, M, baud) \qquad (1)$$

Some of the F functions can be solved only numerically, but a number of them have a closed-form solution. For example, in flat Rayleigh fading channels, The BER for coherently detected M-QAM with Gray bit mapping is approximated. This paper presents an SR system physical layer, which is less complex, to meet constraints set by cellular mobile communication system depending on the specific application and type of data. A communication system must meet various QoS constraints e.g. BER, data rate, or energy dissipation at different times. For example, a cellular system may require a low BER for control data, a high data rate, or minimal energy dissipation for regular updates.

Typical communication systems are built to meet these QoS constraints ever under the worst calculation. However, a communication system operates in the worst condition infrequently. So, it wastes valuable resource as a result existing approaches to meet QoS requirements in cellular systems focus on MAC layer and above, but these optimization are often application dependent. Adaptive technique is 17 dB more power efficient than non-adaptive modulation in fading [11].

In the Physical Layer, the transmitter is the energy efficient pulse generator. The transmitter modulates an input bit stream into a wave of sinusoids with center frequency of 900 MHz and a bandwidth of 2.5 MHz.

$$S(t) = \sum_{i=-\infty}^{\infty} A\, p(t - iTs - \delta * di(t)) \qquad (2)$$

Where, *di(t)* shifts the phase in time by some multiple *δ,* which is larger than the Ts to assure an orthogonal signal set.

The channel model may include AWGN, multipath effects, and also sources of co-channel interference. But, the model under study is AWGN. Note that techniques such as coding and spreading would achieve better performance in a practical system, but they are omitted from the model to focus on the effects of the modulation scheme [12].

In relation to the modulation switching scheme, the function of the MAC protocol is inter-node communication of

the current conditions, resources, and, QoS requirements. It communicates these directly to and from the Physical Layer as header information. The MAC is also responsible for tracking such Physical Layer conditions as channel conditions and link distances.

The Application Layer provides pre-defined QoS constraints depending on the data type. For example, control data may require the lowest possible BER because it is typically not resilient to errors. The error rate and data rate could be sacrificed to operate with minimal energy dissipation. Still other applications may desire a trade-off that results in a compromise between data rate, BER, and energy dissipation [13].

TABLE I. COMPARISON OF MODULATION METHODS WITH DIFFERENT DATA RATES

| Data Rate (Mbps) | Modulation | $N_d$ bps | Data transfer* Duration (μS) |
|---|---|---|---|
| 3 | BPSK | 24 | 2012 |
| 6 | QPSK | 48 | 1008 |
| 12 | 16-QAM | 96 | 504 |
| 24 | 64-QAM | 192 | 252 |
| *for the processor Pentium-IV: 2 GHz. | | | |

The overall strategy for choosing a modulation is motivated by the behavior of MQAM in different environmental conditions. This section characterizes some of this behavior. Configurations with high M and short Ts are more susceptible to interference, which eventually limits the performance regardless of the Eb/N0. In contrast, configurations with low *M* and along Ts are less susceptible to interference, and therefore the performance is more limited by the Eb/No ratio [14]. Therefore, for low Eb/No ratios, modulation schemes with high *M* are preferable. However, for higher Eb/No ratios, the interference dominates the noise, and modulation schemes with low values of *M* are preferable.

The *local cost function (γ)* defines the cost of a transaction.

$$\gamma = \frac{(\alpha \cdot BER) \cdot (\beta \cdot Eb/No)}{(\chi \cdot Data\ Rate)} \qquad (3)$$

In the local function, each performance parameter includes a weighting function that shows the relative contribution of BER, energy, and data rate to the overall cost. In (2), α ☐is the weighting function of the BER, β☐ is the weighting function of the energy, and χ ☐is the weighting function of the data rate. The BER is a function *W* of M, *Ts, Eb/No*, channel impulse response *h(x)*, and interference *i(x)*. The *Eb/No* is a function *Y* of the transmitter-receiver distance, the channel impulse response, and the maximum radiated energy allowed by the ITU-R for a given data rate. The data rate depends only on the Ts and M.

$$BER = W(Eb/No, Ts, m, h(x), i(x)) \qquad (4)$$

$$Eb/No = Y(D, m, Ts, h(x)) \qquad (5)$$

$$Bit\ Rate = \log_2 M / T_s \qquad (6)$$

Thus, from (2) - (5), γ (z) defines the local cost function, where z is a set of environmental and QoS 6-tuples [Eb/N0, m, Ts, h(x), i(x), dist] within **Z**, the set of all possible 6-tuples.

$$\gamma(z) = \frac{\chi}{\alpha\beta} \cdot \frac{Ts \cdot W(Eb/No, Ts, m, h(x) \cdot i(x))\, Y(Dis\tan ce, m, Ts, h(x))}{\log_2 m}$$

$$(7)$$

When environmental conditions vary over several transactions, the local cost function is described in terms of the *local expected cost function*

$$\overline{\gamma} = \sum_{\forall \in Z} \gamma(z) \cdot p(z) \qquad (8)$$

Where *p(z)* is the probability distribution of z $\in$☐Z, and p*(z)* is determined from observation. The system chooses a modulation scheme by optimizing the local expected cost function. A configuration is said to be optimal if it costs less than any other possible configuration for the given environmental and preset parameters. Numerical gradient-based optimization techniques can be used to minimize or

maximize the cost function. However, gradient calculations can be costly for implementation, and hence, calculation may consume more resources than it saves. To reduce implementation complexity, the nodes use linear approximations to choose the optimal modulation scheme for the current operating conditions [15], [16], [17].

## V. SIMULATION RESULTS

The first case is minimum BER, has a fixed data rate, and radiates the maximum allowed Eb. The distances of each transaction determine the Eb/No such that it varies from 5 dB to 12 dB in integer steps, and the distribution is shown in (9).

$$P(E_b/N_o) = \{0.0525, 0dB \leq Eb/No \leq 25dB$$
$$0 \quad else \tag{9}$$

The channel impulse response is a random instance of the model for each transaction, and considers M and Ts values that result in the minimum BER without regard without unacceptably high-energy dissipation or unacceptably low data rate.

Next, the system maximizes the data rate for a target BER in various channel conditions. The environmental and QoS parameters are the same as above, but now the system increases the data rate by choosing an appropriate modulation method. Any non-adaptive system must operate at the fixed data rate of 25 Mbps to meet the target BER even in the best channel conditions.

Third, the system minimizes the radiated energy over various QoS constraints. The first QoS constraint requires a BER of $5 \times 10^{-5}$ and a data rate of 3 Mbps, and this represents the case of distributing microcode over several hops. The second case requires a BER of $2 \times 10^{-4}$ and a data rate of 12 Mbps, and this represents the case of video data.

For different scenarios [17], [18] the system always performs best, and different fixed systems performed second best depending on the QoS goal. For example, the fixed system with $M=16$ achieved the second best BER when BER is the goal, and the fixed system with $M=4$ achieved the best energy efficiency when energy efficiency is goal. Thus, the modulation-switching algorithm should be especially suitable for systems that have changing QoS goals. For example, a sensor network may wish to minimize energy dissipation during normal operation, minimize BER for control data, and maximize data rate when it detects an event. We now compare the performance of the proposed system to the fixed modulation systems as all three of the QoS goals change.

The following results consider that the goals of minimum BER, minimum energy, and maximum data rate all occur with equal probability. Further, the network designer weights each performance parameter such that $\alpha = \beta = \chi = 1$. Note that designers are free to change the weighting function to result in optimal performance for specific applications. The results in Table-2 are normalized and obtained from the cost function of (8). The proposed system performs much better each of the fixed modulation systems. This is in contrast to the previous results that concentrate on just one QoS parameter. In the previous results, one system may have similar performance to the proposed system.

Simulations show that the algorithm improves performance significantly as compared to a conventional, fixed modulation system under variable environmental and QoS requirements. The proposed system improves BER up to 50%, data rate up to 100 %, without sacrificing performance of any other parameter. When the QoS goals change dynamically, the modulation switching system performs significantly better than for a static QoS goal. The proposed Cellular Mobile Software Radios system improves performance for ad hoc and sensor networks by supplementing the MAC protocols in the previous chapters to improve BER, energy efficiency, or data rate. Further, it does not increase the hardware cost or complexity of the Cellular Mobile Software radio transceiver architecture, because it requires changes only to the control logic [19]. Next, a cross-level optimization scheme adapts the Cellular Mobile Software radio to meet various application level QoS constraints as channel conditions change

Cross layer optimization further improves performance for any protocol. The protocols and the optimization scheme are custom tailored to meet the requirements of both Cellular Mobile Software radio. Thus, they significantly outperform more general approaches.

The cpu-time computed using profiler with MATLAB-7 shows that the time required for fixed modulation and that for the modulation switching is approximately same. Also, the data rate is improved. Therefore, on account of negligible complexity the proposed scheme of modulation switching to improve performance of mobile communication system is efficient.

## VI. CONCLUSION

Most of the technological efforts in the wireless communication area are devoted to increase the rational us of resources and devices with two main objectives:

- To increase of efficiency, measured in terms o radio coverage, number of served users, power consumption, spectrum usage, biological impact, short time-to-market and fast network (re)-planning,

- To provide a good degree of services for next generation systems, possibly with a contained investment.

The success of modulation switching algorithm integration into commercial standards such as 3G, WLAN, and beyond (4G, short range communications, etc.) will rely on a fine compromise between rate maximization (Layered type) and diversity (space-time coding) solutions, also including the ability to adapt to the time changing nature of the wireless channel using some form of feedback.

This paper proposes an adaptive system that adapts its resources to efficiently meet QoS requirements in dynamic channel conditions. The system is particularly suitable for, which have demanding QoS requirements that change for

various types of data. The CA block works across the Application, MAC and Physical Layers, which is a departure from established network design techniques. Although this complicates the design process, it meets the special demands of Cellular Mobile Software Radios.

Finally upcoming trials and performance measurements in specific deployment conditions will be required in order to evaluate precisely the overall benefits of SDR systems in real-time world wireless scenarios (for the actual and future services).

TABLE II. BER OF FIXED MODULATION VERSUS MODULATION SWITCHING SYSTEM

| Data Rate (Mbps) | Fixed modulation BER | | | | | | Proposed System BER | Percentage Decrease in BER |
|---|---|---|---|---|---|---|---|---|
| | M=2 | M=4 | M=8 | M=16 | M=32 | M=64 | | |
| 12 | 0.0474 | 0.0310 | 0.0303 | 0.0289 | 0.0355 | 0.0365 | 0.0263 | 44.7% |
| 3 | 0.0418 | 0.0239 | 0.0233 | 0.0252 | 0.0322 | 0.0332 | 0.0212 | 49.3% |

TABLE III. DATA RATE OF FIXED MODULATION VERSUS MODULATION SWITCHING SYSTEM

| Fixed modulation System Data Rate | Modulation switching System Data Rate | Percentage Increase in Data Rate |
|---|---|---|
| 12Mbps | 18.5 Mbps | 68.2% |
| 3 Mbps | 6.1 Mbps | 100% |

.



Figure-6. BER versus $E_b/No$ for modulation schemes used for simulation

REFERENCES

[1] J. Mitalo,"The Software Radio Architecture," IEEE Commun. Mag., vol.33, no.5, Feb. 1995, pp. 26-38.

[2] K C. Zangi and R. D. Koilpillai, "Software Radio issues in cellular Base Station," IEEE JSAC, Vol.1, No.4, pp. 561- 573, April 1999.

[3] J. E. Gunn, & et al., "A Low-Power DSP Core – based Software Radio Architecture", IEEE JSAC, Vol.1, No.4, April 1999, pp.573-590..

[4] B. Hedberg, "Technical challenges in introducing Software radio for mobile telephony base station'" in ACTS software radio workshop, Belgium, May 1997.

[5] Lidtke F. F., "Computer simulation of an automatic classification procedure for digitally modulated communication signals with unknown parameters", Signal Process.1984. 6(4), pp. 311-323.

[6] Jondral F., "Foundations of automatic modulation classification", ITC-Fachberh.1989.107pp.201-206.

[7] Nolan K.E. & et. Al., "Modulation Scheme Classification for 4G Software Radio Wireless Networks", IEEE WCNC, March 2002, Orlando, Florida.

[8] E Del R "Software Radio Technologies and services", Springer Ed., September 2000.

[9] J. Hoffmeyer, II-Pyung Park, M. Majmundar, S. Blust, "Radio software download for Commercial Wireless Reconfigurable devices", pp. s26-s32.

[10] J. J. Patti, R. M. Husnay, "A Smart Software Radio: Concept, development and Demonstration," IEEE JSAC, Vol.1, No.4, April 1999, pp. 631-649.

[11] Oetting J., "Cellular Mobile Radio-An Emerging Technology," IEEE Communications Magazine, Vol.21, No.8, Nov.1983.

[12] Rohde, U. L. & J. C. Whitaker," Communication Receivers: DESP, Software Radios and Design," 3$^{rd}$ Edition, McGraw Hill – N.Y. pp. 321-574, 2001.

[13] A. Wiesler and F. Jondral,o "A Software Radio for second and third Generation Mobile Systems," IEEE Trans. on Vehicular Tech. Vol. 51, No.4, July 2002, pp.738-748.

[14] I. Oka and M.P.C. Fossorier, "A General Orthogonal Modulation Model for Software Radio", IEEE Trans., Comm., vol. 54, No.1, Jan. 2006, pp. 7-12.

[15] Texas Instruments Website at http:// www.ti.com /sc /docs /news /1996 /96037a .htm # support.

[16] Proakis, J. G. & M. Selehi, "Contemporary Communication Systems using MATLAB," Brooks/Cole Thomson Learning, U.S., 1995, pp. 286 – 335.

[17] T. Keller and L. Hanzo, "Adaptive multicarrier modulation: a convenient framework for time frequency professing in wireless communications," *IEEE Proc.*, vol. 88, pp. 611- 640, May 2000.

[18] S. M. Alamouti and S. G.Kallel, "Adaptive trellis coded multiple-PSK FOR Rayleigh fading channels," *IEEE Trans. Commun .*, vol. 42, pp. 2305-2314, June 1994.

[19] A.J.Goldsmith and S.G.Chua, "Variable-rate variable power MQAM for fading channels," *IEEE Trans. Commun .*, vol. 45, pp. 1218-1230, Oct. 1997.

# Randomized Algorithmic Approach for Biclustering of Gene Expression Data

Sradhanjali Nayak[1], Debahuti Mishra[2], Satyabrata Das[3] and   Amiya Kumar Rath[4]

[1,3,4] Department of Computer Science and Engineering, College of Engineering Bhubaneswar, Odisha, INDIA

[2] Institute of Technical Education and Research, Siksha O Anusandhan University, Bhubaneswar, Odisha, INDIA

sradha.mtech09@gmail.com, debahuti@iter.ac.in, satya.das73@gmail.com and amiyamaiya@rediffmail.com

*Abstract*—**Microarray data processing revolves around the pivotal issue of locating genes altering their expression in response to pathogens, other organisms or other multiple environmental conditions resulted out of a comparison between infected and uninfected cells or tissues. To have a comprehensive analysis of the corollaries of certain treatments, deseases and developmental stages embodied as a data matrix on gene expression data is possible through simultaneous observation and monitoring of the expression levels of multiple genes. Clustering is the mechanism of grouping genes into clusters based on different parameters. Clustering is the process of grouping genes into clusters either considering row at a time(row clustering) or considering column at a time(column clustering). The application of clustering approach is crippled by conditions which are unrelated to genes. To get better of these problems a unique form of clustering technique has evolved which offers simultaneous clustering (both rows and columns) which is known as biclustering. A bicluster is deemed to be a sub matrix consisting data values. A bicluster is resulted out of the removal of some of the rows as well as some of the columns of given data matrix in such a fashion that each row of what is left reads the same string. A fast, simple and efficient randomized algorithm is explored in this paper, which discovers the largest bicluster by random projections**.

*Keywords: Bicluster; microarray data; gene expression; randomized algorithm*

## I.    INTRODUCTION

Gene expression data is typically arranged in the form of a matrix with rows corresponding to genes, and columns corresponding to patients, tissues, time points, etc.  Gene expression data are being generated by DNA chip and other microarray technology and they are presented as matrices where each entry in the matrix represents the expression levels of genes under various conditions including environments, individuals and tissues. Each of the $N$ rows represents a gene (or a clone, ORF, etc.) and each of the $M$ columns represents a condition (a sample, a time point, etc.) [8]. It can either be an absolute value (e.g. Affymetrix GeneChip) or a relative expression ratio (e.g. cDNA microarrays). A row/column is sometimes referred to as the "expression profile" of the gene/condition [4]. Due to complex procedure of microarray experiment, gene expression data contains a huge amount of data. Clustering is applied to extract useful information from the gene expression data matrix. The process of grouping data

objects into a set of disjoint class clusters, so that objects within a class have high similarity to each other, while objects in separate classes are more dissimilar [1]. Clustering can be applied either conditions (column clustering). Table 1 show the row clustering where, all the columns for the rows G2, G3 and G4 is selected and table 2 shows column clustering, where C3, C4 and C5 column is clustered with all the rows/genes.

TABLE 1: Row Clustering

|  | $C_1$ | $C_2$ | $C_3$ | $C_4$ | $C_5$ | $C_6$ | $C_7$ |
|---|---|---|---|---|---|---|---|
| $G_1$ | $a_{11}$ | $a_{12}$ | $a_{13}$ | $a_{14}$ | $a_{15}$ | $a_{16}$ | $a_{17}$ |
| $G_2$ | $a_{21}$ | $a_{21}$ | $a_{23}$ | $a_{24}$ | $a_{25}$ | $a_{26}$ | $a_{27}$ |
| $G_3$ | $a_{31}$ | $a_{32}$ | $a_{33}$ | $a_{34}$ | $a_{35}$ | $a_{36}$ | $a_{37}$ |
| $G_4$ | $a_{41}$ | $a_{42}$ | $a_{43}$ | $a_{44}$ | $a_{45}$ | $a_{46}$ | $a_{47}$ |
| $G_5$ | $a_{51}$ | $a_{52}$ | $a_{53}$ | $a_{54}$ | $a_{55}$ | $a_{56}$ | $a_{57}$ |
| $G_6$ | $a_{61}$ | $a_{62}$ | $a_{63}$ | $a_{64}$ | $a_{65}$ | $a_{66}$ | $a_{67}$ |

The classical approach to analyze microarray data is clustering. The process of clustering partitions genes into mutually exclusive clusters under the assumption that genes that are involved in the same genetic pathway behave similarly across all the testing conditions. The assumption might be true when the testing conditions are associated with time points. However, when the testing conditions are heterogeneous, such as patients or tissues, the clustering can be proven as the method of extraction information [6].

TABLE 2: Column Clustering

|  | $C_1$ | $C_2$ | $C_3$ | $C_4$ | $C_5$ | $C_6$ | $C_7$ |
|---|---|---|---|---|---|---|---|
| $G_1$ | $a_{11}$ | $a_{12}$ | $a_{13}$ | $a_{14}$ | $a_{15}$ | $a_{16}$ | $a_{17}$ |
| $G_2$ | $a_{21}$ | $a_{21}$ | $a_{23}$ | $a_{24}$ | $a_{25}$ | $a_{26}$ | $a_{27}$ |
| $G_3$ | $a_{31}$ | $a_{32}$ | $a_{33}$ | $a_{34}$ | $a_{35}$ | $a_{36}$ | $a_{37}$ |
| $G_4$ | $a_{41}$ | $a_{42}$ | $a_{43}$ | $a_{44}$ | $a_{45}$ | $a_{46}$ | $a_{47}$ |

| $G_5$ | $a_{51}$ | $a_{52}$ | $a_{53}$ | $a_{54}$ | $a_{55}$ | $a_{56}$ | $a_{57}$ |
|---|---|---|---|---|---|---|---|
| $G_6$ | $a_{61}$ | $a_{62}$ | $a_{63}$ | $a_{64}$ | $a_{65}$ | $a_{66}$ | $a_{67}$ |

However clustering has got its own limitations. Clustering is based on the assumption that all the related genes behave similarly across all the measured conditions. It may reveal the genes which are very closely co-regulated along the entire column. Based on a general understanding of the cellular process, the subsets of genes are co-regulated and co-expressed under certain experimental conditions. But they behave almost independently under other conditions. Moreover, clustering partitions the genes into disjoint sets i.e. each gene is associated with a single biological function, which is in contradiction to the biological system [8]. In order to make the clustering model more flexible and to overcome the difficulties associated with clustering the concept of biclustering was introduced (see table 3). Biclustering is clustering applied in two dimensions, i.e. along the row and column, simultaneously. This approach identifies the genes which show similar expression levels under a specific subset of experimental conditions. The objective is to discover maximal subgroups of genes and subgroups of conditions. Such genes express highly correlated [18] activities over a range of conditions.

One would expect that a group of genes would exhibit similar expression patterns only in a subset of conditions, such as the subset of patients suffering from the same type of disease. Under this circumstance, biclustering becomes the alternative to the traditional clustering paradigm. Biclustering is a process which performs clustering in two dimensions simultaneously. Clustering method derives a global model while biclustering produces a local model. Biclustering enables one to discover hidden structures in gene expression data in which many genetic pathways might be embedded [2]. It might also allow one to uncover unknown genetic pathways, or to assign functions to unknown genes in already known genetic pathways, while clustering technique is applied a given gene cluster is defined using all the conditions ,similarly each condition cluster is defined for all genes. But each gene in a bicluster is selected using only a subset of the conditions and each condition in a bicluster is selected using only a subset of genes [2]. The goal of biclustering is to identify subgroups of genes and subgroups of conditions by performing simultaneous clustering of both the rows and columns instead of in two dimensions separately as in clustering [2].

Randomized algorithm approach is based on the idea of randomly selecting a set of columns and rows [6]. It is a very simple, effective method to find bicluster on both the aspect of time complexity and space complexity. The sub matrix produced by the biclustering has the property that each row reads the same string, so such a sub matrix would therefore correspond to a group of genes that exhibit a coherent pattern of states over a subset of conditions. [3].

## A. Proposed Model



Figure 1: Our Proposed Model

Our proposed work is to find the biclusters from gene expression data using randomized algorithm. First, we have used a synthetic data set, and then we have validated our work with Yeast data set [20]. Second, we have pre-processed our data set using Z-score method to put the attribute values in a standard range of values. Finally, we validate our randomized model by comparing our model with existing biclustering models by considering various parameters. Our model (See figure 1) outperforms the existing model of Cheng and Church [15]on the basis of run time for finding number of patterns and also the scalability issues have been found to be improved significantly considering both the attributes and objects as they increases.

## B. Paper Layout

This paper is arranged in the following manner, section I gives the introduction as well as our proposed model is also outlined, section II deals with related work on biclustering models. In section III the preliminary information about gene expression data, bicluster, randomized approach, problem statement and algorithms are described. Section VI describes our proposed algorithm. Section V gives the analysis of our work and shows its significance over the Cheng and Church[15 ] algorithm. Finally, section VI gives the conclusion and future directions of our work.

## II. RELATED WORK

Shyama Das et al [13] proposed a greedy randomized adaptive search procedure to find the biclusters. The bicluster seeds are generated using k-means algorithm, and then these seed are enlarged using GRASP. GRASP happens in two phases i.e construction and local search. In the construction phase a feasible solution is developed iteratively by adding one element each time which will generate a feasible solution whose neighborhood will be searched until a local minimum is identified during the local search phase. The best solution is stored as the result.

In this study GRASP is applied for the first time to identify biclusters from Human Lymphoma dataset. In this paper the GRASP meta heuristics is used for finding biclusters in gene expression data. In the first step K-Means algorithm is used to group rows and columns of the data matrix separately. Then they are combined to produce small biclusters.

Bing Liu et al [7] proposed an efficient semi-unsupervised gene Selection method via spectral biclustering. From biological and clinical point of view finding smaller number of important genes help the doctor to concentrate on these genes and investigating the mechanism for cancer causes and its remedies.

Haider Banka et al. [8] give an evolutionary biclustering of gene expression data. They have proposed to uncover genetic pathways (or chains of genetic interactions) which is equivalent to generating clusters of genes with expression levels that evolve coherently under subsets of conditions, *i.e.*, discovering biclusters where a subset of genes are co-expressed under a subset of conditions. Such pathways can provide clues genes that contribute towards a disease. This emphasizes the possibilities and challenges posed by biclustering. The objective here is to find sub matrices or maximal subgroups of conditions where the genes exhibit highly co-related activities over a range of conditions.

Stefano Lonardi et al. [6] find biclusters by random projection. From a given matrix *X* composed of symbols, a bicluster is a sub matrix of *X* obtained by removing some of the rows and columns, so that each row left will read the same string. An efficient randomized approach is used to find largest bicluster which is probabilistic that is each entry of the matrix is associated with the probability.

Daxin Jiang et al. [11] proposed an interactive exploration of gene expression patterns from a gene expression data set. Analyzing coherent gene expression patterns is an important task in bioinformatics research and biomedical applications. The development of microarray technology provides a great opportunity for functional genomics. Identifying co-expressed genes and coherent expression patterns in gene expression data can help biologists understand the molecular functions of the genes and the regulatory network between the genes. However, due to the distinct characteristics of gene expression data and the special requirements from the biology domain, mining coherent patterns from gene expression data presents several challenges, which cannot be solved by traditional clustering algorithms.

### III. PRELIMINARIES

#### A. Microarray or Gene Expression Data

Microarrays is a small chip made of chemically coated glass , nylon membrane or silicon onto which thousands of DNA molecules are attached in fixed grids[19]. Microarray is used in the medical domain to produce molecular profiles of diseased and normal tissues of patients. Microarray captures the expression level of thousands of genes under one experiment. Microarray operations are done under different condition to have parallel comparison between the experimental levels of gene. The relative abundance of mRNA of a gene is called the expression level of a gene [9].This is measured using DNA microarray technology which revolutionized the gene expression study by simultaneously measuring the expression levels of thousands of genes in a single experiment [8][13].

The data generated by these experiments high dimensional matrix contain thousands of rows (genes) and hundreds of conditions. The experimental conditions can be patients, tissue types, different time points etc. Each entry in this matrix is a real number which denotes the expression level of a gene. Genes participating in the same biological process will have similar expression patterns. Clustering is the suitable mining method for identifying these patterns [1][13]. The ability of arrays to monitor thousands of separate but unrelated events simultaneously has captured the thoughts of scientists practicing in both basic and applied research [8][16][17].

The process of microarray formation experiment is associated with a collection of experimental factors describing the variables under study, e.g. "disease state","gender state". Each microarray in an experiment takes on a specific value for each of the experimental factors, e.g. "disease state = normal" and "gender = male" [9][19]. In the very first stage the mRNA (messenger RNA) of normal male cell and a cancer male cell is obtained through RNA isolation process. Then by the reverse transcriptage enzyme the cDNA is obtained from mRNA. The cDNA (complimentary DNA) of the diseased cell is labeled with red color and the normal cell cDNA is labeled with green color. Then by the hybridization process the diseased cell and normal cell is hybridized to a small chip made of chemically coated glass, nylon membrane or silicon called as microarray in a fixed form (grids). Gene expression data are being generated by DNA chip and other microarray technology and they are presented as matrices of expression levels of genes under various conditions including environments, individuals and tissues. Gene expressions provide a fundamental link between genotypes and phenotypes, and play a major role in biological processes [18][19] and systems including gene regulation, evolution, development and disease mechanism.

A gene expression data from microarray experiment is represented by a real valued matrix. $M = \{ A_{ij} | 1 \leq i \leq n , 1 \leq j \leq m \}$ where rows ,$G = \{ g_1, g_2, g_3, \ldots \ldots g_r \}$ represents the expression pattern of the genes and the column, $S = \{ s_1, s_2, s_3, \ldots \ldots s_c \}$ represents expression profiles for samples and each element $w_{ij}$ is measured expression level of gene $i$ in sample $j_1$ which is shown in the below table 3.

TABLE 3: Gene Expression Data

| Gene | Condition 1 | ... | Condition j | ... | Condition c |
|------|-------------|-----|-------------|-----|-------------|
| Gene₁ | $a_{11}$ | .... | $a_{1j}$ | ... | $a_{1c}$ |
| Gene... | ... | .... | .... | ... | ... |
| Geneᵢ | $a_{i1}$ | ... | $a_{ij}$ | ... | $a_{ic}$ |
| Gene... | ... | .... | ... | ... | ... |
| Geneᵣ | $a_{r1}$ | ... | $a_{rj}$ | ... | $a_{rc}$ |

Where, $r$ = no of genes, $c$ = no of samples, $M$ = gene expression data matrix, $a_{ij}$ = element in the gene expression matrix, gene = different whose expression levels are taken in the row and condition = genes are studied under different conditions which are taken in the column.

Gene expression data set contains thousands of genes while the no. of tissue sample ranges from tens to hundreds, while analyzing expression profiles, a major issue is gene selection for target phenotype [5][19]. For example Cancer is a disease that begins in the cells of the body. Cancer is ultimately the result of cells that uncontrollably grow and don't die. Cancer occurs when cells become abnormal and keep dividing and forming more cells without order or control. From biological and clinical point of view finding the small number of important genes can help medical researchers to concentrate on these gens and investigating the mechanism for cancer development. Clustering is a reputed algorithmic technique that partitions a set of input data (vectors) into subsets such that data in the same subset are close to one another in some metric [3]. Recent developments require finding the largest bicluster satisfying some additional property with the largest area. For a given matrix of size $n \times m$ over a alphabet set $\sum$, a bicluster is a sub matrix composed of selected columns and rows satisfying a certain property [3]. Bicluster is a subset of genes that jointly respond across a subset of conditions, where a gene is termed responding under some condition if its expression level changes significantly under that condition with respect to its normal level. A bicluster of a gene expression data is a local pattern such that the gene in the bicluster exhibit similar expression patterns through a subset of conditions [6].

Each bicluster is represented as a tightly co-regulated sub matrix of the gene expression matrix. A $(X, Y)$ is a matrix, $I =$ subset of rows, $J =$ subset of columns and $(I, Y) =$ a subset of rows that exhibits similar behavior across the set of columns $=$ cluster of rows. $(X, J) =$ a subset of columns that exhibit similar behavior across set of all rows $=$ cluster of columns. $(I, J) =$ is a bicluster i.e. subset of genes and subsets of conditions, where the genes exhibit similar behavior across the conditions and vice versa. Cluster of columns $(X, J) = (C_3, C_4, C_5)$, Cluster of Rows $(I, Y) = (G_2, G_3, G_4)$, Bicluster $(I, J) = \{ (G_2, G_3, G_4), (C_3, C_4, C_5)\}$. (See table 4)

TABLE 4: Bicluster

| | $C_1$ | $C_2$ | $C_3$ | $C_4$ | $C_5$ | $C_6$ | $C_7$ |
|---|---|---|---|---|---|---|---|
| $G_1$ | $a_{11}$ | $a_{12}$ | $a_{13}$ | $a_{14}$ | $a_{15}$ | $a_{16}$ | $a17$ |
| $G_2$ | $a_{21}$ | $a_{21}$ | $a_{23}$ | $a_{24}$ | $a_{25}$ | $a_{26}$ | $a_{27}$ |
| $G_3$ | $a_{31}$ | $a_{32}$ | $a_{33}$ | $a_{34}$ | $a_{35}$ | $a_{36}$ | $a_{37}$ |
| $G_4$ | $a_{41}$ | $a_{42}$ | $a_{43}$ | $a_{44}$ | $a_{45}$ | $a_{46}$ | $a_{47}$ |
| $G_5$ | $a_{51}$ | $a_{52}$ | $a_{53}$ | $a_{54}$ | $a_{55}$ | $a_{56}$ | $a_{57}$ |
| $G_6$ | $a_{61}$ | $a_{62}$ | $a_{63}$ | $a_{64}$ | $a_{65}$ | $a_{66}$ | $a_{67}$ |

The basic goal of biclustering is to identify subgroups of genes and subgroups of conditions, where the genes exhibit highly correlated activities for every condition, Identify sub-matrices with interesting properties and to perform simultaneous clustering on the rows and column dimensions of

the genes. The underlying bases for using bi-clustering in the analysis of gene expression data are similar genes may exhibit similar behaviors only under a subset of conditions, not all conditions, genes may participate in more than one function resulting in one regulation pattern in one context and a different pattern in another.

*B. Randomized Approach for finding Biclusters*

Biclustering algorithms may have two different objectives: to identify one or to identify a given number of biclusters. Randomized algorithm is an approach to find one bicluster at a time which is very easy to understand and implement [6].

Let's assume that, given a large set of a data matrix, $X \in \sum_{n \times m}$ from which a sub matrix, $x_{(r*, c*)}$ has to be discovered where the sub matrix $x_{(r*, c*)}$ is the largest one from the data matrix set. For the simplicity $r^* = |R^*|$ and $c^* = |C^*|$. The concept of the algorithm owes its origin to the following simple observation. It is analyzed that if we can know what is the value of $R^*$ then we can easily determine $C^*$ by selecting the clean columns with respect to $R^*$ or if instead we know $C^*$, then, $R^*$ could be obtained by taking the maximal set of rows which read the same string. Unfortunately, if neither $R^*$ nor $C^*$ is known then the approach is to "sample" the matrix by random algorithm, with the expectation that at least some of the projections will overlap with the solution $(R^*, C^*)$, one can focus to either rows or columns, but here, in this algorithm, it is described how to retrieve the solution by sampling columns.

The steps for the algorithm are as described below:

1. Select a random subset of columns as $S$ of size $k$ uniformly from the set of columns $\{1, 2, \ldots, m\}$.
2. Lets assume that for the instant that $S \cap C^* \neq \Phi$. If we know $S \ C^*$, then $(R^*, C^*)$ could be determined by the following three steps:

   a. select the string(s) $w$ that appear exactly $r^*$ *times* in the rows of $X [1:n.S \cap C^*]$
   b. set $R^*$ to be the set of rows in which $w$ appears and
   c. set $C^*$ to be the set of clean columns corresponding to $R^*$.

Given a selection of rows $R$, we say that a column $j$, $1 \leq j \leq m$, is *clean* with respect to $R$ if the symbols in the $j$th column of $X$ restricted to the rows $R$, are identical. In general, a solution of the largest bicluster can contain a column of zeros, as long as they appear in all rows of the sub matrix [6].

*C. Problem Statement*

The main problem behind this algorithm is to find the largest bicluster from the given data matrix.

Largest Biclsuter$(f)$ problem
Instance: let $\sum$ denotes the set of nonempty symbols.
Let $X$ be a gene expression data matrix asdefined over the alphabet $\sum^{n \times m}$ of symbol.

$n$ = no of rows or genes
$m$ = no of columns or conditions

The set ∑ denotes a non-empty *alphabet* of *symbols* and a *string* over ∑ *is* an ordered sequence of symbol Largest Biclsuter (*f*) problem.

*Objective:* To find a row selection $R$ and a column selection $C$ such that the rows of $X(R,C)$ are identical strings and the objective function $f(X(R,C))$ is maximized from the alphabet set.

Assume that we are given a large matrix $X \in \sum^{n \times m}$ in which a sub matrix $X \in (R^*, C^*)$ is to be selected. Assume also that the sub matrix $X(R^*, C^*)$ is maximal. To simplify, let the notations are $r^* = |R^*|$ = set of rows and $C^* = |C^*|$ = set of columns. Let the examples of objective functions which are used as a basis to find the bicluster are as follows:

- ○ $f1(X(R,C)) = |R| + |C|$;
- ○ $f2(X(R,C)) = |R|$ provided that $|C| = |R|$; and
- ○ $f3(X(R,C)) = |R|\|C|$.
- ○ $f4(x_{(r^*,c^*)}) = |R^*| \cap |C^*|$

## IV. OUR PROPOSED ALGORITHM

*Randomized search (step 1):* Select a random subset $S$ of size $k$ uniformly from the set of columns $\{1, 2, …, m\}$ ;

Example: Let us take an example of data matrix as follows:

TABLE 5: Example Data Matrix

| 2 | 1 | 0 | 1 | 1 | 2 |
|---|---|---|---|---|---|
|   |   |   |   |   |   |
| 0 | 0 | 1 | 1 | 0 | 2 |
| 0 | 1 | 2 | 0 | 1 | 1 |
| 2 | 0 | 0 | 1 | 0 | 2 |
| 1 | 2 | 1 | 2 | 1 | 2 |
| 0 | 0 | 2 | 1 | 0 | 1 |

Let us randomly select 3 columns as: $C_1^* = (1,2,3,4)$ Let us randomly select another 4 columns as: $C_2^* = (2,4,5,6)$

*Randomized search (Step-2):* From the selected columns take the common columns which are the subset of the given matrix $X$. As per our example the common column is : $C_1^* \cap C_2^* = (2,4)$. The common column is shown in green color in the below table 6.

*Randomized search (Step-3):* For all the subset of $S$, find the occurrences of string $w$ that appears at least $r$ times in each subset of $S$.

As per our example,
   The String    11   appears = 1
   The String    01   appears = 3
   The string    12   appears = 1
   The string    02   appears = 1

TABLE 6: Example

| 2 | 1 | 0 | 1 | 1 | 2 |
|---|---|---|---|---|---|
|   |   |   |   |   |   |
| 0 | 0 | 1 | 1 | 0 | 2 |
| 0 | 1 | 2 | 0 | 1 | 1 |
| 2 | 0 | 0 | 1 | 0 | 2 |
| 1 | 2 | 1 | 2 | 1 | 2 |
| 0 | 0 | 2 | 1 | 0 | 2 |

*Randomized search (Step-4):* Record the maximum no string which appears in the subset and record the corresponding rows. As per our example, the maximum string which appears is 01 and the corresponding rows are rows are (2,4,6)

*Randomized search (step-5):*

- Select the set of clean columns $C$ with size at least ' $c$' corresponding to each $R$
- A column $j$ is clean with respect to $R$ if the symbols in the $j^{th}$ column of $X$ restricted to the rows $R$, are identical.

As per our example, the clean column with respect to the rows are (5,6).

*Randomized search (step-6):* Save the solutions and repeat step 1 to 4 for $t$ iterations. As per our example, the largest bicluster is:

$$ X'' = \begin{pmatrix} 0 & 1 & 0 & 2 \\ 0 & 1 & 0 & 2 \\ 0 & 1 & 0 & 2 \\ 0 & 1 & 0 & 2 \end{pmatrix} $$

Parameters used in the algorithms are as follows:

- Projection size $k$ ($k_{min}$)
- Column threshold $c$
- Row threshold $r$
- Number of iterations $t$
- In our example the clean column w.r.t the rows are (5,6)

## V. RESULT ANALYSIS

In this paper, we have simulated the randomized biclustering algorithm to find the maximal bicluster embedded in the data matrix using the synthetic data set as well as Yeast data set [20]. We have also implemented the the Minimum Square Residue (MSR) approach of Cheng and Church [15] to find the biclusers.

We have tested both the approaches in Intel Dual Core machine with 2GB HDD. The OS used is Microsoft XP and all programs are written in C. We have observed the similar trends on runtime versus number of biclusters found in both MSR based approach and our proposed randomized approach, the figure 2 shows the running time is significantly less as

compared to MSR approach. Table 7 shows the comparative study on both the approaches.



Figure 3: Performance Analysis

TABLE 7: Comparative Analysis

| Model | Run Time (ms) For synthetic data set | Run Time (ms) for Yeast data set | No.of Biclusters found for synthetic data set | No.of Biclusters found for Yeast data set |
|---|---|---|---|---|
| MSR based Method | 4000 | 24000 | 18 | 286 |
| Our Randomized Approach | 1500 | 20000 | 24 | 788 |

## VI. CONCLUSION

The simultaneous clustering of the rows and columns of a matrix falls under diversified names such as biclustering, co-clustering or two mode clustering. The unique features of gene expression data and the specific demands from the domain of biology, propels different challenges on the front of coherent patterns from gene expression data which is a taxing task for traditional clustering algorithms. The underlying basis for using biclustering in the analysis of gene expression data are similar genes may exhibit similar behaviors only under a subset of conditions, not all conditions. Genes are regulated by multiple factors/processes concurrently. Genes may participate in more than one function resulting in one regulation pattern in one context and a different pattern in another. Using biclustering algorithms, one can obtain sets of genes that are co-regulated under subsets of conditions. Here, we have presented a rather simple algorithm based on random projections. We have presented a probabilistic analysis of the largest bicluster problem, which allows one to determine the statistical significance of a solution. In future we plan to extend our system to the following aspects randomized approach provides a flexible and consistent model to organize the expression patterns in individual data sets. In future this approach can be extended to the application of various soft computing techniques as genetic algorithm, pattern matching, unsupervised learning algorithm which can able to find more than one bicluster from the single data set. We are hopeful that this concept of biclustering will meet the future challenges and will prove itself as more effective and result oriented.

## REFERENCES

[1] Daxin jiang, Chun Tang, and Aidong Zhang ,"Cluster analysis for gene expression data: a survey", *IEEE Transaction On Knowledge and Data Engineering,* Vol. 16, no.11. November 2004.

[2] S.C. Madeira and A.L. Oliveira, "Biclustering Algorithms for Biological Data Analysis: A Survey," IEEE/ACM Trans. Computational Biology and Bioinformatics, Vol. 1, No. 1, pp. 24-45, 2004.

[3] Gahyun Park and Wojciech Szpankowsk, "Analysis of biclusters with application to gene expression data". International Conference on Analysis of Algorithms, pp.267–274, 2005

[4] Mel N Kronick," Creation of the whole human genome microarray, *Technology Profile, Future Drugs Limited*, pp.19-28,www.futute-drugs.com

[5] L. Lazzeroni, A. Owen, "Plaid models for gene expression data", *Statistica Sinica* 12 (1), pp.61–86, 2002

[6] Stefano Lonardia,,Wojciech Szpankowskib, QiaofengYanga," Finding biclusters by random projections", *Theoretical Computer Science* , 368, pp. 217 – 230, 2006

[7] Bing Liu, Chunru Wan, and Lipo Wang, *"*An Efficient Semi-Unsupervised Gene Selection Method via Spectral Biclustering". *IEEE Transactions on Nano Bioscience*, Vol. 5, No. 2, 2006.

[8] Alain B. Tchnag and Ahmed H.Tewfik," DNA Microarray Data Analysis: A Novel Biclustering Algorithm Approach". Volume 2006, pp. 1- 12 , DOI 10.1155/ASP/2006/59809,2006

[9] Jos´e Caldas, Nils Gehlenborg , Ali Faisal , Alvis Brazma and Samuel Kaski, *"*Probabilistic retrieval and visualization of biologically relevant microarray experiments*", BMC Bioinformatics.*, 10(Suppl 13), pp.1-9, 2009.

[10] Arifa Nisar, Waseem Ahmady Wei-keng Liao, Alok Choudhary,"High Performance Parallel/Distributed Biclustering Using Barycenter Heuristic". *SIAM,* pp. 1050-1061, 2009.

[11] R. Sharan and R. Shamir. Click: A clustering algorithm with applications to gene expression analysis. In ISMB, pages 307-216, 2000.

[12] Q. Sheng, Y. Moreau, B.D. Moor, "Biclustering Microarray data by Gibbs sampling", *Proceeding of European Conf. on Computational Biology". (ECCB'03)*,pp. 196–205,2003.

[13] Shyama Das, Sumam and Mary Idicula," Application of Greedy Randomized Adaptive Search Procedure to the Biclustering of Gene Expression Data", International Journal of Computer Applications, Volume 2 – No.3,pp. 0975 – 8887, 2010.

[14] Haider Banka, Sushmita Mitra, "Evolutionary Biclustering of Gene Expressions", *ACM Ubiquity*, Volume 7, Issue 42 , 2006

[15] Yizong Cheng and George M. Church. "Biclustering of expression data". *In proceedings of the 8th International conference on intelligent systems for molecular Biology* (ISMB '00'), pages 93-103, 2000.

[16] Keisuke Lida, Ichiro Nishimura, "Gene expression profiling by DNA Microarray Technology", *Critical Reviews in Oral Biology and Medicine*, Vol. 13 no. 1,pp. 35-50, 2002.

[17] Daxin Jiang ,Jian Pei, Aidong Zhang ,"Towards Interactive Exploration of Gene Expression Patterns ", *SIGKDD Explorations*, 5(2), pp.79-90 ,2003.

[18] Haider Banka, Sushmita Mitra " Evolutionary  biclustering of gene expression data", *Proceedings of the 2nd international conference on Rough sets and knowledge technology*, Pages: 284-291, 2007.

[19] Madan Babu,M., Luscombe, N., Aravind, L., Gerstein, M., Teichmann, S.A. , Structure and evolution of transcriptional regulatory networks. *Curr. Opin. Struct. Biol*. ,2004

[20] UCI Repository for Machine Learning Data bases retrieved from the *World Wide Web: http://www.ics.uci.edu*

[21] Uetz  P., et al." A Comprehensive analysis of protein protein interaction in saccharomyces cerevisiae", *Nature*, 403(6770): 601-3, Feb-2000.

[22] Gavin A.C., et. al. "Functional organization of yeast proteome by systematic analysis of protein complexes". *Nature* 415(6868) :13-4, Jan-2002.

AUTHORS PROFILE

Sradhanjali Nayak is a scholar of M.Tech (CSE) at College of Engineering, Biju Pattanaik University, Bhubaneswar, Odisha, INDIA. Her research areas includes Data mining, Soft Computing Techniques etc.

Debahuti Mishra is an Assistant Professor and research scholar in the department of Computer Sc. & Engg, Institute of Technical Education & Research (ITER) under Siksha O Anusandhan University, Bhubaneswar, Odisha, INDIA. She received her Masters degree from KIIT University, Bhubaneswar. Her research areas include Datamining, Bio-informatics, Software Engineering, Soft computing. Many publications are there to her credit in many International and National level journal and proceedings. She is member of OITS, IAENG and AICSIT. She is an author of a book Aotumata Theory and Computation by Sun India Publication (2008).

Satyabrata Das is as Assistant Professor and Head in the department of Computer Sc. & Engineering, College of Engineering Bhubaneswar (CEB). He received his Masters degree from Siksha O Anusandhan  University, Bhubaneswar. His research area includes Data Mining, Adho-network etc.

Dr.Amiya Kumar Rath obtained Ph.D in Computer Science in the year 2005 from Utkal University for the work in the field of Embedded system. Presently working with College of Engineering Bhubaneswar (CEB) as Professor of Computer Science & Engg. Cum Director (A&R) and is actively engaged in conducting Academic, Research and development programs in the field of Computer Science and IT Engg. Contributed more than 30 research level papers to many national and International journals. and conferences Besides this, published 4 books by reputed publishers. Having research interests include Embedded System, Adhoc Network, Sensor Network, Power Minimization, Biclustering, Evolutionary Computation and Data Mining.

# A Method of Genetic Algorithm (GA) for FIR Filter Construction: Design and Development with Newer Approaches in Neural Network Platform

Ajoy Kumar Dey, Susmita Saha
Department of Signal Processing
Blekinge Tekniska Hogskola (BTH)
Karlskrona, Sweden
e-mail: dey.ajoykumar@yahoo.com

Avijit Saha, Shibani Ghosh
Department of Electrical and Electronic Engineering
Bangladesh University of Engineering and Technology
Dhaka, Bangladesh
e-mail: avijit003@gmail.com

*Abstract*—**The main focus of this paper is to describe a developed and dynamic method of designing finite impulse response filters with automatic, rapid and less computational complexity by an efficient Genetic approach. To obtain such efficiency, specific filter coefficient coding scheme has been studied and implemented. The algorithm generates a population of genomes that represents the filter coefficient where new genomes are generated by crossover, mutation operations methods. Our proposed genetic technique has able to give better result compare to other method.**

*Keywords-Genetic Algorithm; FIR: filter design; optimization; neural network.*

## I. INTRODUCTION

This paper represent a developed and dynamic method for genetic algorithm for to design a FIR filter in Neural network platform where the FIR filter has certain kinds of finite impulse response and genetic algorithm provide a automatic, efficient and less complex design method [6].

Select the standard polynomial transfer function by satisfied the response specification, followed by the implementation of the transfer function in one of the standard circuit structures is the conventional approach to filter design where optimization approach is required in many case of this approach [7]. The non standard response specification, the computation complexity for digital filters and many other ways we need to use this approach [3], [4].

We know that Genetic algorithm can be successfully employed for minimizing or maximizing a cost function based on the genetic and natural selection. The design of quantized digital FIR filters, being an optimization problem over a discrete coefficient space, can therefore be faced using a genetic approach [3]. Genetic algorithm optimization methods have emerged as a powerful approach to solving the more difficult optimization problems. This paper gives the attention to theoretical analysis of genetic algorithm and a efficient and developed method analysis for genetic algorithm for FIR filter design [1].

The paper organization is as follows: Section II describes the Genetic algorithm and its implementation. Section III discusses Research rationale. Section IV describes about the design method of FIR filter and Section V discusses about the Analysis of findings. Finally, Section VI presents a conclusion and an indication towards the future scope of this work.

## II. GENETIC ALGORITHM (GA)

### A. General Idea of GA

Genetic algorithms are stochastic strategies for optimization. Evolutionary computation is used by the genetic algorithms which one is one of the major features. As an optimization tool, it can be used as traning algorithm for any supervised Neural Network [5].

Natural selection and Natural genetics mechanics approaches first time introduced by the genetic algorithm where this two mechanism is the based part for the GA. Genetic algorithms may be differentiated from more conventional techniques by these characteristics, a) direct manipulation of the encoded representation of variables at the string level, rather than manipulation of the variables themselves, b) search form a population of points rather than form a single point, thus reducing the probability of reaching a false peak, c) blind search by sampling, ignoring all information except the outcome of the sample, d) use of stochastic rather than deterministic operators [6].

Effectiveness in searching large, noisy, multimodal problem spaces and versatility of genetic algorithm gave by these characteristics where smooth and differentiable surfaces are to be searched, calculus based methods are likely to do better and where the problem space is small, genetic algorithms may show no advantage over enumerative and random search methods .
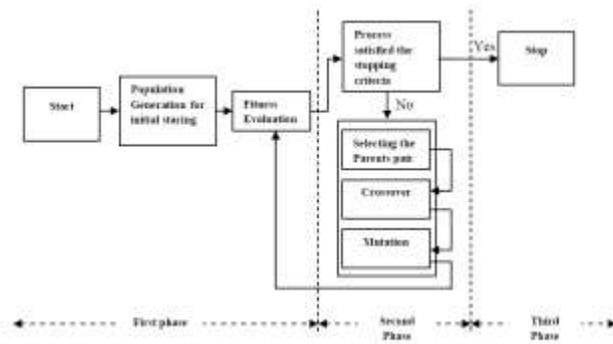
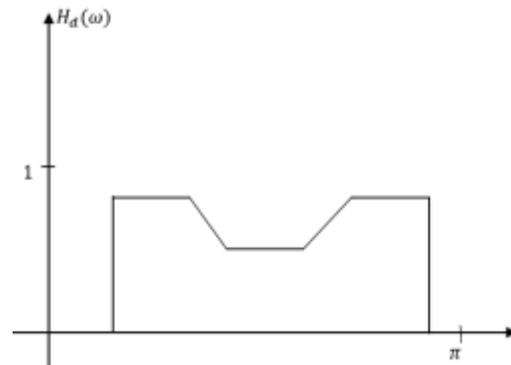Figure 1.   Natural Flow of Genetic algorithm.



Figure 2.   Observation of desired frequency response.

Inherently parallel search performed by the genetic algorithms and make effective use of historical information in their search for improved areas of performance as well [2]. They are not the only class of optimization technique to have been inspired by evolutionary process where evolutionary strategies based on selection according to fitness, mutation and reproduction.

There are many advantages and drawback in genetic algorithm. i) It can escape local minimums in the error function, ii) any error function can be used easily (not only the least square error function), iii) genetic algorithm can be easily implemented in parallel computers; those are consider

as the benefits with genetic algorithm. As well so many disadvantages with genetic algorithm like, a) the

computational cost is high, b) the convergence speed varies strongly with different tasks, c) some optimization problems are not solvable with genetic algorithms.

The genetic algorithm described in this paper in a simple way only the process of selection according to fitness, crossover, mutation and reproduction.

### B.   Structure of Genetic Algorithm

There are many ways to implements a GA but most of genetic algorithm consists of an iteration of 3 steps (from figure 1). Like, 1) selection of parent (we do not really need two parents even though it is normally used. In real life the use of two parents is one way to make the iteration process continue.), 2) Creation of new population, 3) Mutation of the population.

Now we consider the multilayer feed forward neural network to use genetic algorithm in order to train the network. Assume we have a training sequence, $(x_1,t_1),(x_2,t_2),.....,(x_N,t_N)$ and we have initialized the weights $W_1,W_2,W_3$ (parents)

As a step 1 we have to Create P number of new sets of weights

$$\begin{array}{ll} 1: & W_1^1, W_2^1, W_3^1 \\ 2: & W_1^2, W_2^2, W_3^2 \\ \vdots & \vdots \quad \vdots \quad \vdots \\ P: & W_1^P, W_2^P, W_3^P \end{array} \Bigg\} \text{New Population}$$

By changing each entry in each parent matrix, with probability $p_i$, by adding a small random value [4].

As a step 2, now we have to pick a random input target pair from the training set, say $(x_i,t_i)$. Now we have to calculate the output from all children and create the error signal.

As a step 3, We have to recorder the children so that the best children can get the place in top of the position and the rest children are in decreasing order. As a step 4, if the best child

is worse than the parent go to step 1. As a step 5, we should iterate it until no changes in parent [11].

### III.   RESEARCH RATIONALE

Genetic Algorithm is optimization tools. As such we can use them in any supervised neural network and also use in unsupervised neural networks, as long as we have error criteria [12]. In this part, we wish to create an FIR filter with an arbitrary transfer function by the use of Genetic Algorithms (GA:s). The FIR filter should be a linear phase filter of odd length. This will give symmetry in the filter taps [3].

We will use a Genetic Algorithm where each parents gives birth to 10 children as a new population. We use the exponential decease function to weight the suitability of becoming a new parent. When initiating the first parent use a pure delay of half the FIR filter length an impulse at position $L - 1/2$ of the FIR filter taps and the rest of the taps at zero [9], [10].

Now we have to create a function where a desired frequency response and the length of the FIR filter are taken as input parameters. The function should then find the least

square or the mini max approximation to the desired transfer function with Genetic Algorithm [7]. When finding the least square or the mini max filter, the only thing that has to be modified is the definition of the error function [8].

We have to remember that when we using the mini max error criteria one should allow for transition regions where the desired frequency response has abrupt changes. Omitting some points around these regions when sampling the frequency does this. Otherwise the maximum error will always appear at these regions.  Now we have to use the FIR filter lengths: 33, 65, and 129. The step size and the probability of change, $p_i$ , should be decreased during learning. The constant in the exponentially weighing can be set to 0.1. Sample the frequency region at 1024 points in order to calculate the error function. Now we have to use an appropriate stopping criterion.

## IV.    DESIGN METHOD OF FIR FILTER

Consider a linear phase FIR filter to approximate any given transfer function. From the (fig. 2 ) we want to find a L tap FIR filter which have the following transfer function. L- Odd integer number.

We denote the filter taps as $h(n), n = 0, 1, ..., L - 1$. The Fourier transfer of this filter is,

$$H(\omega) = \sum_{n=0}^{L-1} h(n) e^{-j\omega n}$$

(1)

$$H(\omega) = e^{-j\omega \frac{L-1}{2}} \left( h\left(\frac{L-1}{2}\right) + \sum_{n=\frac{L-1}{2}+1}^{L-1} 2h(n) \cos\left(\omega\left(n - \frac{L-1}{2}\right)\right) \right)$$

(2)

It has linear phase and it is symmetric. The absolute value of $H(\omega)$ is,

$$|H(\omega)| = h\left(\frac{L-1}{2}\right) + 2 \sum_{n=\frac{L-1}{2}+1}^{L-1} h(n) \cos(\omega(n - \frac{L-1}{2}))$$

(3)

(since $|e^{jx}| = 1$). We sample the frequency in $[0, \pi]$ with N points,

$$H_d(\omega) = [H_d(\omega_1), H_d(\omega_2), ....., H_d(\omega_N)]^T$$

(4)

And

$$H(\omega) = [H(\omega_1), H(\omega_2), ....., H(\omega_N)]^T$$

(5)

We can now define an error function,

$$E(H) = \|H_d(\omega) - H(\omega)\|_2^2$$

(6)

TABLE I.        THE NUMBER OF THE FILTER ORDER IMPLEMENTED IN THE NETWORK.

| Experiment No. | No. of Filter Order for LS appro. | Implemented Corresponding Functions | | Another Observation | |
|---|---|---|---|---|---|
| | | | | Filter order for Mini max | Filter order for Band pass filter |
| 1 | 33 | Error Function | Frequency response | 33 | 33 |
| 2 | 65 | Error Function | Frequency response | 65 | 65 |
| 3 | 129 | Error Function | Frequency response | 129 | 129 |

Where,

$$h = \left[ h\left(\frac{L-1}{2}\right), h\left(\frac{L-1}{2} + 1\right), ...., h(L-1) \right]^T$$

(7)

or

$$E(h) = \max|H_d(\omega) - H(\omega)|$$

(8)

We now use the Genetic Algorithm to find the impulse response,

$$h = \left[ h\left(\frac{L-1}{2}\right), h\left(\frac{L-1}{2} + 1\right), ...., h(L-1) \right]^T$$

(9)

Initialize h to small random values.

## V.    ANALYSIS OF FINDINGS

According to the table 1 we implemented the number of filter order at LS approximation, Mini max approximation find out the corresponding error function and frequency response. As we mentioned above, we use in here the three filter order 33, 65 and 129 and observe the corresponding function. First we find out the LS approximation with this filter orders and find out the corresponding error function and frequency response function. After that we find out the Mini max approximation with this filter orders.

To choose the faster and efficient approximation between the LS and mini max, we implemented all the filter orders to LS approximation for a band pass filter and we can analysis all the output at the figure 3.

Basically the LS approximation is faster than Mini Max Approximation. But the Step size is also a Factor which can effect on the Convergence.

## VI.    CONCLUSION

For to do Filter Design and coefficient optimization, we always use the genetic algorithm as efficient and powerful tool. When genetic algorithm designing new filters, it directly optimizes the coefficients by considering the quantization effect. On the other hand, when the genetic algorithm mapping predesigned filter coefficients to integer

arithmetic, since hand quantization does not results in optimal results, it can quickly find a satisfactory result. This genetic algorithm can modify and reshaped as the specific needs and characteristics of the applications.

## REFERENCES

[1] Sabbir U. Ahmad and Andreas Antoniou , "Cascade-Form Multiplierless FIR Filter Design Using Orthogonal Genetic Algorithm", IEEE International Symposium on Signal Processing and Information Technology, 2006, pp. 932-937, Aug. 2006

[2] Khadijeh Khamei, Abdolreza Nabavi, Shaahin Hessabi ,"Design Of Variable Fractional Delay Fir Filteirs Using Genetic Algorithm", Proceedings of the 2003 10th IEEE International Conference on Electronics, Circuits and Systems, 2003. ICECS 2003, vol. 1, pp. 48-51, Dec. 2003.

[3] Mehmet Oner and Murat Agkar. "Incremental Design Of High Complexity Fir Filters By Genetic Algorithms", vol.2, pp. 1005-1008, 1999.

[4] Mehmet ONER, "A Genetic Algorithm for Optimisation of Linear Phase FIR Filter Coefficients", Conference on Signals, Systems & Computers, 1998, vol.2, pp. 1397-1400, Nov. 1998.

[5] Paolo Gentili, Fruncesco Piazza and Aurelio Uncini, "Efficient Genetic Algorithm Design For Power-Of-Two Fir Filters", International Conference of Acoustics, Speech, and Signal Processing, 1995. ICASSP-95., vol.2, pp. 1268-1271, May 1995.

[6] D. Suckley ,"Genetic algorithm in the design of FIR filters", IEE Proceedings Circuits, Devices and Systems, vol. 138, pp. 234-238, Apr. 1991

[7] .B. Deng, "Discretization-free design of variable fractional delay FIR filters," IEEE Trans. Circuits Syst. 11: Analog and Digital Signal Processing, June 2001, Vol. 48, No. 6, pp. 637- 644.

[8] C. K. **S.** Pun, **Y.** C. Wu, K. L. Ho, **"**An efficient design offractional delay digital FIR filters using the Farrow structure," Proc. of the 11th EEE Signal Processing Workshop on Statistical Signal Processing, 2001 ,pp. 595 -598.

[9] Y. W. Leung and Y. Wang, "An orthogonal genetic algorithm with quantization for global numerical optimization," IEEE Trans. Evolutionary Comp., vol. 5, no. 1, pp. 41 - 53, Feb. 2001.

[10] J. W. Adams, "FIR digital filters with least-squares stopbands subject to peak-gain constraints," IEEE Trans. Circuits Syst. vol. 39, pp. 376-388, Apr. 1991.

[11] M. Yagyu, A. Nishihara, and N. Fujii, "Fast FIR digital filter structures using minimal number of adders and its application to filter design," IEICE Trans. Fundamentals, vol. E79-A, no. 8, pp. 1120-1128, Aug. 1996.

[12] T.Arslan, D.H.Horrocks, "A genetic algorithm for the design of finite word length arbitrary response cascaded IIR digital filtersLEEE Int. Conf on GA in engineering systems, Sheffield, Sept. 1995, pp 276-281.

(a)  (c)  (e)  (g)



(b)  (d)  (f)  (h)

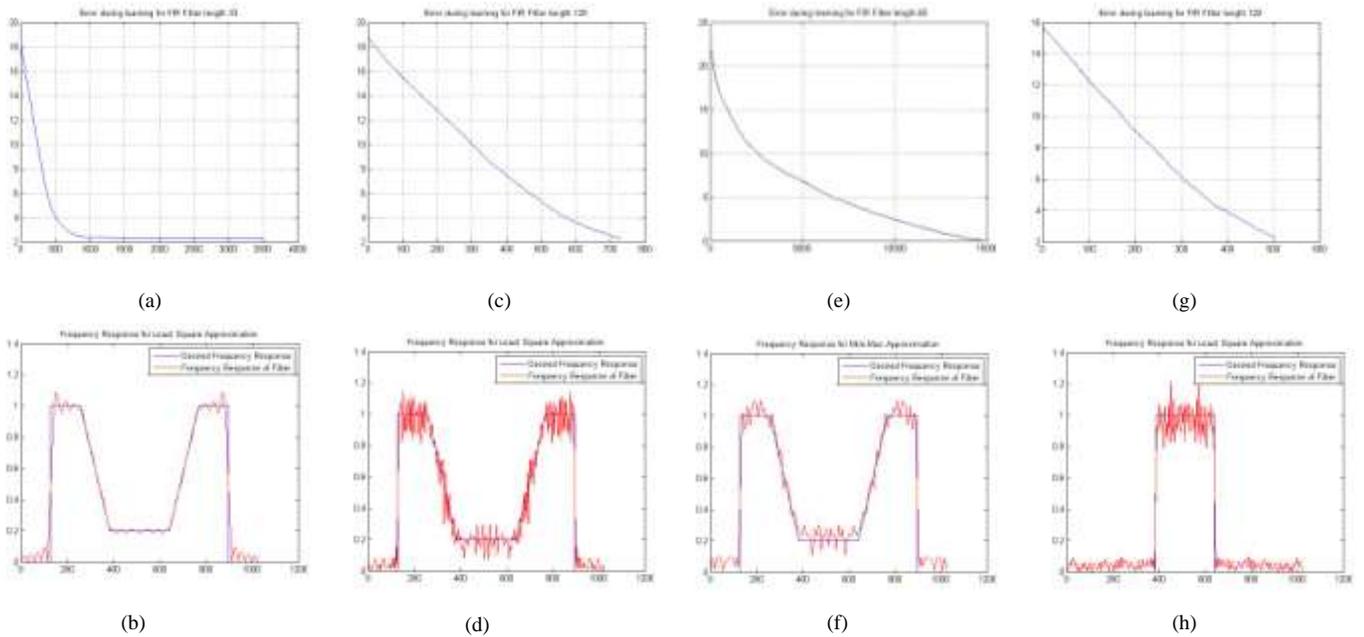Figure 3. (a) Learning Error for filter order 33, (b) Frequency response for LS approximation for filter order 33, (c) Learning Error for filter order 129, (d) Frequency response for LS approximation for filter order129, (e) Learning Error for filter order 65 in mini max approximation, (f) Frequency response for mini max approximation for filter order 65, (g) Learning Error for filter order 129 for bandpass filter, (h) Frequency response for LS approximation for filter order 129 for bandpass filter.

# Hybrid Technique for Human Face Emotion Detection

Renu Nagpal, Pooja Nagpal
M.Tech Student,CE deptt.
Yadavindra College of Engineering
Guru Khashi Campus
Talwandi Sabo, Bathinda, (Punjab)

Sumeet Kaur
Assistant Professor,CE deptt.
Yadavindra College of Engineering
Guru Khashi Campus
TalwandiSabo,Bathinda,(Punjab)

*Abstract*—**This paper presents a novel approach for the detection of emotions using the cascading of Mutation Bacteria Foraging optimization and Adaptive Median Filter in highly corrupted noisy environment. The approach involves removal of noise from the image by the combination of MBFO & AMF and then detects local, global and statistical feature form the image. The Bacterial Foraging Optimization Algorithm (BFOA), as it is called now, is currently gaining popularity in the community of researchers, for its effectiveness in solving certain difficult real-world optimization problems. Our results so far show the approach to have a promising success rate. An automatic system for the recognition of facial expressions is based on a representation of the expression, learned from a training set of pre-selected meaningful features. However, in reality the noises that may embed into an image document will affect the performance of face recognition algorithms. As a first we investigate the emotionally intelligent computers which can perceive human emotions. In this research paper four emotions namely anger, fear, happiness along with neutral is tested from database in noisy environment of salt and pepper. Very high recognition rate has been achieved for all emotions along with neutral on the training dataset as well as user defined dataset. The proposed method uses cascading of MBFO & AMF for the removal of noise and Neural Networks by which emotions are classified.**

*Keywords-biomertics;adaptive median filter; bacteria foraging optimization;feature detection;facial expression*

## I. FACIAL EXPRESSION RECOGNITION

Biometric is the science and technology of recording and authenticating identity using physiological or behavioral characteristics of the subject. A biometric representation of an individual. It is a measurable characteristic, whether physiological or behavioral, of a living organism that can be used to differentiate that organism as an individual. Biometric data is captured when the user makes an attempt to be authenticated by the system. This data is used by the biometric system for real-time comparison against biometric samples. Biometrics offers the identity of an individual may be viewed as the information associated with that person in a particular identity management system.

Automatic recognition of facial expressions may act as a component of natural human machine interfaces Such interfaces would enable the automated provision of services that require a good appreciation of the emotional state of the service user, as would be the case in transactions that involve negotiation[1], for example Some robots can also benefit from the ability to recognize expressions .

Noise is added as unwanted variations in the image when image is transmitted over the network [22]. It causes a wrong conclusion in the identification of images in authentication and also in pattern recognition process. Firstly there should be the removal of noise from the image then features are detected. Noise in imaging systems is usually either additive or multiplicative. In practice these basic types can be further classified into various forms such as amplifier noise or Gaussian noise, Impulsive noise or salt and pepper noise, quantization noise, shot noise, film grain noise and non-isotropic noise. However, in our experiments, we have considered only salt and pepper impulsive noise. .

## II. RELATED WORK

Automated analysis of facial expressions for behavioral science or medicine is another possible application domain. From the viewpoint of automatic recognition, a facial expression can be considered to consist of deformations of facial components and their spatial relations, or changes in the pigmentation of the face[1].

There is a vast body of literature on emotions. Recent discoveries suggest that emotions are intricately linked to other functions such as attention, perception, memory, decision making, and learning. This suggests that it may be beneficial for computers to recognize the human user's emotions and other related cognitive states and expressions. Ekman and Friesen [1] developed the Facial Action Coding System (FACS) to code facial expressions where movements on the face are described by a set of action units (AUs). Ekman's work inspired many researchers to analyze facial expressions by means of image and video processing.

The AAM approach is used in facial feature tracking due to its ability in detecting the desired features as the warped texture in each iteration of an AAM search approaches to the fitted image. Ahlberg [7] use AAM in their work. In addition, ASMs - which are the former version of the AAMs that only use shape information and the intensity values along the

profiles perpendicular to the shape surface are also used to extract features such as the work done by Votsis et al. [9].

Many algorithms have been developed to remove salt/pepper noise in document images with different performance in removing noise and retaining fine details of the image, like:

Simard and Malvar [10] shows image noise can originate in film grain, or in electronic noise in the input device such as scanner digital camera, sensor and circuitry, or in the unavoidable shot noise of an ideal photon detector. Beaurepaire et al. [2] tells the identification of the nature of the noise is an important part in determining the type of filtering that is needed for rectifying the noisy image. Noise Models from Wikipedia [11] shows the noise in imaging systems is usually either additive or multiplicative. Image Noise [12] shows in practice these basic types can be further classified into various forms such as amplifier noise or Gaussian noise, Impulsive noise or salt and pepper noise, quantization noise, shot noise, film grain noise and non-isotropic noise. Al-Khaffaf [13] proposes several noise removal filtering algorithms. Most of them assume certain statistical parameters and know the noise type a priori, which is not true in practical cases.

Prof. K. M. Passino [8] proposed an optimization technique known as Bacterial Foraging Optimization Algorithm (BFOA) based on the foraging strategies of the E. Coli bacterium cells. Until date there have been a few successful applications of the said algorithm in optimal control engineering, harmonic estimation in Ref [15], transmission loss reduction in Ref [16], machine learning in Ref [14] and so on. Its performance is also heavily affected with the growth of search space dimensionality Kim *et al* [17] proposed a hybrid approach involving GA and BFOA for function optimization. Biswas *et al* [18] proposed a hybrid optimization technique, which synergistically couples the BFOA with the PSO.

Up to our knowledge this is the first time we are using this hybridized technique for the detection of emotions in noisy environment. However, in our experiments, we have considered only salt and pepper impulsive noise.

### 2.1 Organization of the Paper

The rest of the paper is organized as follows: Section 3 describes general set up of proposed technique, experimental results have been discussed with respect to percentage of correct recognition considering JAFFE facial image database under salt and pepper noisy environment in section 4 and comparative analysis of the proposed technique with existing techniques at the end of section 4. The paper is concluded with some closing remarks and future scope in section 5.

### III. THE GENERAL SET UP

The design and implementation of the Facial Expression Recognition System can be subdivided into three main parts:

- The first part is Image Pre-processing.

- Second part is a Recognition technique, which includes Training of the images.
- Third part is testing and then there is result of classification of images.

### A. Image Preprocessing

The image processing part consists of image acquisition of noisy image. Filtering, Feature Extraction, Region of Interest clipping, Quality enhancement of image. This part consists of several image-processing techniques. First, noisy face's image acquisition is achieved by scanner or from JAFFE database by introducing the noise in the images then adaptive median filter is used to remove noise from the image then Mutation Bacteria Foraging Optimization Technique is used to remove noise that is remained there after using median filter and finally features are extracted. The region of interest is eyes, lips, mouth (eyes and lips) that are independently selected through the mouse for identification of feature extraction. Statistical analysis is also done by that mean, median and standard deviation of the noisy frame, restored frame, cropped frame and enhanced frame are calculated. The significance of these shows that mean and standard deviation should be less in enhanced frame as compared to noisy frame. These extracted features of image are then fed into Back-Propagation and Radial Basis Neural Network for training and then emotions are detected.

### B. Training of Neural Network

The Second part consists of the artificial intelligence, which is composed by Back Propagation Neural Network and Radial Basis Neural network. First training of the neurons is there and then testing is done. Back-propagation and RBF algorithms are used in this part. It consists of two layers. At input to hidden layer Back-propagation Neural Network is used which consists of feed forward and feed backward layers and at hidden to output layer, RBF Neural Network is used for classification of expressions.

- As the recognition machine of the system, a three layer neural network has been used that is trained with several times using hit and trial method on various input ideal and noisy images forced the network to learn how to deal with noise. The window size is of 9.We can increase the size of window with that computation time is also increased. The variation in the density of noise is taken from 0.05 to 0.9.
- **The combination of adaptive median filter and Mutation BFO removes noise up to 90% from the image with more accuracy** and the learning ranges from 0.1 to 0.9. Main accuracy or goal is 0.01, for that it takes more computation time. For single image total of 10 iterations are needed for zero error but for 21 images the error is zero in total of 15000 epochs because there are different images and different types of motions so more iterations are there.

### C. Testing

The third phase consists of testing of expressions that shows the percentage of accurate results and result of classification for different expression.

• A special advantage of the technique is that the expression is recognized even there is more noise density in the image up to 0.9. The range of density is taken from 0.05 to 0.9. Second by taking statistical features like mean, median and standard deviation the classification is easy and there will be more correctness to recognize the facial expression. Third, dual enhancement of image is there, first at the removal of noise by Adaptive Median Filter and Mutation BFO and second by using histogram equalization.

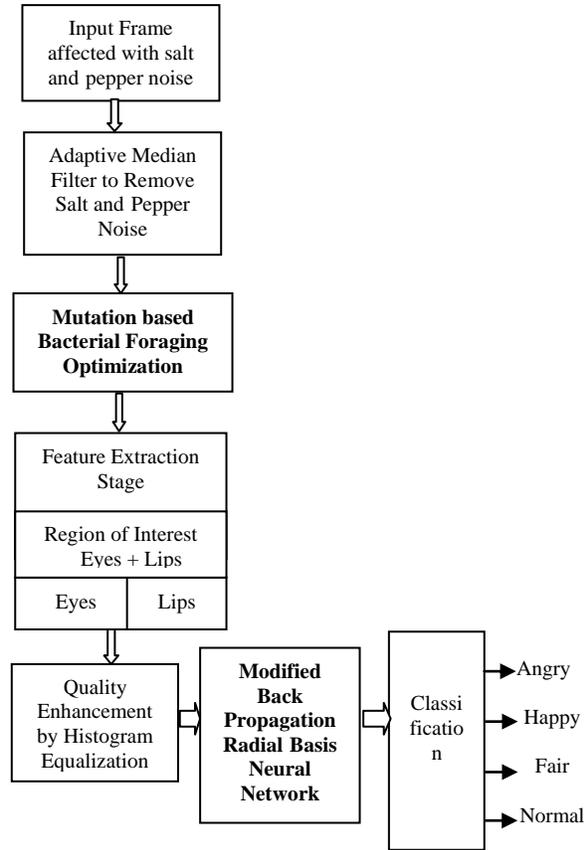• The flowchart and Implementation Overview of Facial Recognition System with proposed method is shown in figure 1 and in figure 2.

```
┌─────────┐
│  Start  │
└─────────┘
     │
     ▼
┌──────────────────┐
│ Noisy Input Frame│
└──────────────────┘
     │
     ▼
┌──────────────────────┐
│ Adaptive Median Filter│
└──────────────────────┘
     │
     ▼
┌──────────────────────────────────────────┐
│ Mutation Based Bacterial Foraging Optimization│
└──────────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────┐
│ Feature Extraction stage Eyes/Lips/Face│
└──────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│ Quality Enhancement by Histogram Processing│
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────┐
│ RBF Neural Network│
└──────────────────┘
     │
     ▼
┌──────────────────┐
│  Classification  │
└──────────────────┘
```

Flowchart of Facial Expression Recognition System

*Block Diagram*

```
┌──────────────────┐
│   Input Frame    │
│ affected with salt│
│ and pepper noise │
└──────────────────┘
        │
        ▼
┌──────────────────┐
│ Adaptive Median  │
│ Filter to Remove │
│ Salt and Pepper  │
│      Noise       │
└──────────────────┘
        │
        ▼
┌──────────────────┐
│ **Mutation based**│
│**Bacterial Foraging**│
│  **Optimization**│
└──────────────────┘
        │
        ▼
┌──────────────────┐
│ Feature Extraction│
│      Stage       │
├──────────────────┤
│ Region of Interest│
│   Eyes + Lips    │
├────────┬─────────┤
│  Eyes  │  Lips   │
└────────┴─────────┘
        │
        ▼
```

Implementation Overview of Facial Expression Recognition System

### A.    Input Noisy Image

The first module shows the input phase. To this module a noisy face image of salt and pepper noise is passed as an input for the system. The input image samples are considered from JAFFE database. The input image is randomly picked up from the database used for training and evaluated for the recognition accuracy.

### B.    Adaptive Median Filter

A median filter is an example of a non-linear filter [20] .It is very good at preserving image detail. To run a median filter:

1.  Consider each pixel in the image
2.  Sort the neighboring pixels into order based upon their intensities
3.  Replace the original value of the pixel with the median value from the list

•  **Algorithm of Adaptive Median Filter**

1. Initialize w = 3 and Zm = 39.

2. Compute Z min, Z max and Z med

3. If Zmin < Zmed < Zmax, then go to step 5. Otherwise w = w+2.

4. If $w \leq Zm$, go to step 2. Otherwise, replace Zxy by Z med

5. If Zmin< Zxy< Zmax, then Zxy is not noisy pixel else replace Zxy by Zmed.

where,

Z = Noisy image.
Zmin = Minimum intensity value.
Zmax = Maximum intensity value.
Zmed = Median of the intensity values.
Zxy = Intensity value at coordinates (x, y).
Zm = Maximum allowed size of the adaptive median filter window.

• A median filter is a rank-selection (RS) filter, a particularly harsh member of the family of rank-conditioned rank-selection (RCRS) filters, a much milder member of that family, for example one that selects the closest of the neighboring values when a pixel's value is external in its neighborhood, and leaves it unchanged otherwise, is sometimes preferred, especially in photographic applications. Median filter is good at removing salt and pepper noise from an image, and also cause relatively little blurring of edges, and hence are often used in computer vision applications.

### C.    Mutation Bacteria Foraging Optimization

During the first stage the input image corrupted with a salt-and-pepper noise of varied densities from 0.05 to 0.9 is applied to the adaptive median filter. In the second stage, both the noisy and adaptive median filter output images are passed as search space variables in the BFO technique [20] to minimize errors due to differences in filtered image and noisy image.

• Bacterial Foraging Optimization with fixed step size suffers from two main problems
    I.   If step size is very small then it requires many generations to reach optimum solution. It may not achieve global optima with less number of iterations.
    II.  If the step size is very high then the bacterium reach to optimum value quickly but accuracy of optimum value gets low.

Similarly, in BFO, chemotaxis step provides a basis for local search, reproduction process speeds up the convergence, elimination and dispersal helps to avoid premature convergence.

To get adaptive step size, increase speed and to avoid premature convergence, the mutation by PSO is used in BFO instead of elimination and dispersal event by equation 1.

$$\theta^i(j+1,k)=\theta^i(j+1,k)+*r_1*C_1(\theta^i(j+1,k)-\theta_{global}) \qquad \ldots\ldots\ldots 1$$

$\theta^i(j,k)$ = Position vector of *i*-th bacterium in *j*-th chemotaxis step and *k*-th reproduction steps.

$\theta_{global}$=Best position in the entire search space

In-preprocessing step of hybrid soft computing technique, adaptive median filter is used to identify pixels that are affected by noise and replacing them with median value to keeps the uncorrupted information as far as possible. The BF-pfPSO follows chemotaxis, swarming, mutation and reproduction steps to obtain global optima. The algorithm of BF-pfPSO is presented below.

• **Step by step algorithm of mutation based BFO**

Initialize Parameters *p, S, Nc, Ns, Nre, Ned, Ped* and *C (i), i=* 1, 2... *S*.
Where,
*p* = Dimension of search space
*S* = Number of bacteria in the population
*Nc* = Number of chemotaxis steps
*Ns* = Number of swimming steps
*Nre* = Number of reproduction Steps
*Pm* = Mutation probability
*C (i)* = Step size taken in the random direction specified by the tumble
$\dot{\theta}$ *(j, k)*= Position vector of the *i*-th bacterium, in *j*-th chemotaxis step, in *k*-th reproduction step and in l-th elimination and dispersal step
*Step 1:* Reproduction loop: *k = k+1*
*Step 2:* Chemotaxis loop: *j = j+1*
  a) For *i*= 1,2…*S*, take a chemotaxis step for bacterium i as follows
*b)* Compute fitness function *J (i, j, k, l)*
c) Let *Jlast =J (i, j, k,)* to save this value since we may find a better cost via a run.
d) Tumble: Generate a random vector $\Delta(i) \in \Re^p$ with each element $\Delta_m(i)$, *m = 1, 2...p*, a random number on [-1 1]
e) Move: Let

$$\theta^i(j+1,k)=\theta^i(j,k)+C(i)\frac{\Delta(i)}{\sqrt{\Delta^T(i)\Delta(i)}}$$

*f)* Compute *J (i, j+1, k, l)*
g) Swim
i) Let *m =0* (counter for swim length)
ii) While *m < Ns* (if have not climbed down too long)
o   Let *m = m+1*
o  If *J (i, j+1, k, l) < Jlast* (if doing better),
Let *Jlast = J (i, j+1, k, l)* and let

$$\theta^i(j+1,k)=\theta^i(j+1,k)+C(i)\frac{\Delta(i)}{\sqrt{\Delta^T(i)\Delta(i)}} \quad \text{And use}$$

this $\theta^i(j+1,k)$ to compute the new *J (j+1, k).*
o Else, let *m = Ns*. This is the end of the while statement.

h) Go to next bacteria *(i+1)* if *i ≠S*

*Step 3:* Update $\theta_{pbest}(j,k)$ *and* $\theta_{global}$. If *j <Nc*, go to step 3.
In this case, continue chemotaxis, since the life of bacteria is not over.
*Step 4:* Reproductions:

a) For the given *k* and *l*, and for each *i = 1, 2…S*, let

$$J^i_{health} = \sum_{j=1}^{Nc+1} J(i,j,k)$$

be the health of bacterium $i$ . Sort bacteria and chemotaxis parameter $C (i)$ in order of ascending cost $J_{health}$ (higher cost means lower health).

b) The $Sr =S/2$ bacteria with the highest $J_{health}$ values die and other $Sr=S/2$ bacteria with the best values split.

**Step 5: (New step): Mutation**

For $i = 1, 2…S$, with probability $Pm$, change the bacteria position by pfPSO.

$$\theta^i ( j+1,k) = \theta^i ( j+1,k) + *r_1 * C_1(\theta^i ( j+1,k) - \theta_{global})$$

*Step 6:* If k < Nre, go to step 2. We have not reached the specified number of reproduction steps. Therefore, we have to start the next generation in the chemotaxis loop

### D.     Pre-Processing and Feature Extraction

The face image passed is transformed to operational compatible format in this phase, where the face image is resized to uniform dimension, the data type of the image sample is transformed to double precision and passed for feature extraction.

As the first step in image processing, the region of interest (ROI) of a lip and an eye or only lips region or only eyes region have been selected independently in the acquired images through the mouse. The ROI image is converted into grayscale image.

### E.     Histogram Equalization

A histogram equalization method has been applied before obtaining the filtered grayscale image. Histogram equalization improves the contrast in the grayscale and its goal is to obtain a uniform histogram. The histogram equalization method also helps the image to reorganize the intensity distributions. New intensities are not introduced into the image. Existing intensity values will be mapped to new values but the actual number of intensity pixels in the resulting image will be equal or less than the original number. In the image sequence, the histogram-equalized image is filtered using average and median filters in order to make the image smoother. Hence, the histogram-equalized image is split into lip ROI and eye ROI regions and then the regions are cropped from the full image. The problem of light intensity variations has been solved.

### F.     Modified Back Propagation and Radial Basis Neural Network

The neurons are trained by hit and trail method, a total of 625 input neurons are taken, hidden neurons are75 and output neurons are 4. Total of 13 pair are trained with the different emotions of happiness, anger, fear and neutral. In this phase epochs and errors are calculated of particular face region. Two parameters are used, total numbers of epochs and errors are calculated in neural network. For more accuracy more computation time is needed. The main accuracy or goal is 0.01 then it takes more computation time. The Classification of Neural Network includes two types of neural networks that were trained based on the input parameters extracted. Back Propagation Neural Network and Radial Basis Neural Network.

Two types of neural networks were trained based on the input parameters extracted that are:

(i) Back Propagation Neural Network
(ii) Radial Basis Neural Network

### (i) Back Propagation Neural Network

The most widely used neural network is the Back Propagation algorithm. This is applied at input to hidden layer, due to its relative simplicity, together with its universal approximation capacity. The learning algorithm is performed in two stages: feed-forward and feed- backward.

In the first phase the inputs are propagated through the layers of processing elements, generating an output pattern in response to the input pattern presented. In the second phase, the errors calculated in the output layer are then back propagated to the hidden layers where the synaptic weights are updated to reduce the error. This learning process is repeated until the output error value, for all patterns in the training set, are below a specified value.

The Back Propagation, however, has two major limitations: a very long training process, with problems such as local minima and network paralysis; and the restriction of learning only static input -output mappings. To overcome these restrictions, RBF Neural Network is used. Table 1 shows the parameters for Back Propagation Neural Network.

TABLE I.          PARAMETERS FOR BACK PROPAGATION nNEURAL NETWORK

| Input Neurons | 625 {25 x25] |
|---|---|
| Hidden Neurons | 75 |
| Output Neurons | 4 |
| Network Size | 625-75-4 |
| Training Pairs | 21 |
| Learning Rate | 0.5 |
| Maximum Epochs | 5000 |
| Activation Functions | Sigmoid, Sigmoid |
| Error Goal | 0.001 |

### (ii)     Radial Basis Functions

Radial Basis Functions (RBF) has attracted a great deal of interest due to their rapid training, generality and simplicity. When compared with traditional multilayer perceptrons, RBF networks present a much faster training, without having to cope with traditional Back Propagation problems, such as network paralysis and the local minima. These improvements have been achieved without compromising the generality of applications. It has been proved that RBF networks, with enough hidden neurons, are also universal approximators.

The RBF is basically composed of three different layers: the input layer, which basically distributes the input data, one hidden layer, with a radially symmetric activation function, hence the network's name and one output layer, with linear

activation function. Table 2 shows the parameters for Radial Basis neural Network

TABLE II.        PARAMETERS FOR RADIAL BASIS NEURAL NETWORK

| Input Neurons | 625 {25 x25] |
|---|---|
| Hidden Neurons | 75 |
| Output Neurons | 4 |
| Network Size | 625-75-4 |
| Training Paris | 21 |
| Learning Rate | 0.5 |
| Maximum Epochs | 5000 |
| Activation Functions | Sigmoid, Radbas |
| Error Goal | 0.001 |

The classification system of Neural Network consisted of three stages.

1. Training of Neural Network
2. Testing of Neural Network
3. Performance Evolution of Neural Network
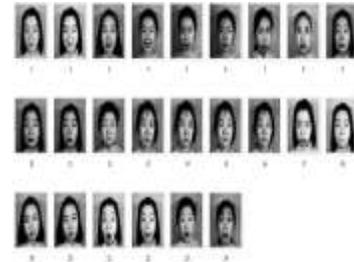
*G.    Classification*

A selected database of 21 images of 4 classes is considered to demonstrate the capability and the accuracy of the recognition stage. The faces presented are the inputs into the training stage where a representative set of facial features are determined. After training, new images are processed and entered into the recognition stage for identification. Then emotions are detected as angry, happy, fear and normal.

IV.    EXPERIMENTAL RESULTS

The proposed algorithm technique is applied on sample images. When the noise level increases, the face images get more affected and sometimes are not visible. Hence in our experiments, we have considered mean and variance varying from 0.05 to 0.9. To start with, applying the salt and peeper noise with mean and variance equal to 0.05 on all the images of the JAFFE face database forms the probing image set. All the images in the JAFFE database without adding any noise are taken as the prototype image set. Hence we get all images in prototype set. Applying our technique on both the sets forms the feature set. In the experimental phase, we take the first image of the first subject from the prototype image set as the query image and the top matching ten images are found from a set of all probe images. If the top matching images lie in the same row (subject) of the prototype query image, then it is treated as a correct recognition. The number of correct recognized images for each query image in the prototype image set is calculated and the results are shown for salt and peeper noise of variance 0.8.

For our experiments, the facial images from the facial image database JAFFE are used. The database contains 213 images of 7 facial expressions (6 basic facial expressions + 1 neutral) posed by 10 Japanese female models. 60 Japanese

subjects have rated each image on 6 emotion adjectives. Figure 3 shows the sample images used in our experiments collected from JAFFE face database. In our experiments, we have taken 1 person images as 21 images for a single person from that 13pairs are trained in the training with different emotions we have used common type of noise namely, salt and pepper impulsive noise that affect the biometric image processing applications. In order to show the robustness of our face recognition method, these noises are introduced in the JAFFE database face images. Figure 3 shows the sample of image database.



Sample images from JAFFE database

The results consist of four sections:

- In **first section (Section 4.1)**, consists of feature extraction stage and enhancement results including preprocessing results of lip feature, eye feature and mouth feature to recognize facial expression of a single person. It consists of noisy image of salt and pepper noise, and then by applying adaptive median filter and mutation bacteria foraging optimization the restored image is appeared. The cropped region is shown in figures and is clearly expressed in the histograms .In the end the enhanced region and enhanced histogram clarifies the results. The statistical features as shown in tables also measure all results.

- **Second section (Section 4.2)**, consists the results for training of neural network. These results are presented in tabular form, and then the report and graph for all of the features is also shown in results.

- **Third section (Section 4.3),** consists the results for testing of frames, which are also shown in tabular form.

- **Fourth section (Section 4.4)** shows the classification results for angry, happy, fear and neutral features.

- **Fifth section (Section 4.5)** shows the comparison of the proposed technique with existing ones **.**

*4.1(Section 1) Feature Extraction Stage and Enhancement Results*

*A.  Preprocessing Results for Lip feature*
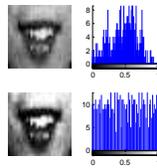


Fig4 (a)          Fig4 (b)          Fig4 (c)

Fig4 (d), (e), (f)

Figure 4 Preprocessing Results for Lip feature
(a)   Noisy Image with noise variance 0.8
(b)   Restored image from Noisy image by using   median filter &
MBFO
(c)   Cropped Lip Region in restored image
(d)   Cropped Lip region and its Histogram
(e)   Enhancement of cropped Lip region by Histogram Equalization
(f)   Histogram after enhancement

TABLE III.        STATISTICAL FEATURES FOR LIPS

| Statistical Features | Noisy Frame | Restored Frame | Cropped Frame | Enhanced Frame |
|---|---|---|---|---|
| Mean | 0.5022 | 0.5016 | 0.5000 | 0.4998 |
| Median | 0.5469 | 0.5429 | 0.4921 | 0.4221 |
| Standard Deviation | 0.4578 | 0.2229 | 0.2134 | 0.2014 |

*B.   Preprocessing Results for Eye feature*


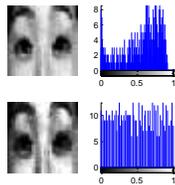
Fig5 (a)              Fig5 (b)              Fig5 (c)



Fig5 (d), (e), (f)

Figure 5  Preprocessing Results for Eye feature
(a)   Noisy Image with noise variance 0.8
(b)   Restored image from Noisy image by using median filter & MBFO
(c )  Cropped Eye Region in restored image
(d)  Cropped Eye region and its Histogram
(e)  Enhancement of cropped Eye region by Histogram Equalization
(f)    Histogram after enhancement

TABLE IV.        STATISTICAL FEATURES FOR EYES

| Statistical Features | Noisy Frame | Restored Frame | Cropped Frame | Enhanced Frame |
|---|---|---|---|---|
| Mean | 0.5018 | 0.5016 | 0.5005 | 0.4996 |
| Median | 0.5510 | 0.5469 | 0.5079 | 0.4921 |
| Standard Deviation | 0.4583 | 0.2229 | 0.2134 | 0.2035 |

*C.  Preprocessing Results for Mouth feature*



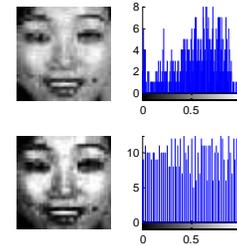Fig6 (a)              Fig6 (b)              Fig6 (c)



Fig6 (d), (e), (f)

Figure 6  Preprocessing Results for (Eye and lip) feature
(a)        Noisy Image with noise variance 0.8
(b)        Restored image from Noisy image by using   median filter &
MBFO
(c )       Cropped Mouth (Eye & Lip) Region in restored image
(d)        Cropped Mouth (Eye & Lip) region and its Histogram
(e)        Enhancement of cropped Mouth (Eye & Lip) region by Histogram
Equalization
(f)        Histogram after enhancement

TABLE V.        STATISTICAL FEATURES FOR MOUTH

| Statistical Features | Noisy Frame | Restored Frame | Cropped Frame | Enhanced Frame |
|---|---|---|---|---|
| Mean | 0.5021 | 0.5016 | 0.5001 | 0.5002 |
| Median | 0.5510 | 0.5469 | 0.5079 | 0.5079 |
| Standard Deviation | 0.4585 | 0.2229 | 0.2937 | 0.2935 |

*4.2 (Section 2) Training of Neural Network Results*

This part shows the results for training of neural network. It consists of a single frame of one subject which includes total of 21 images of different emotions, from which 13 images are trained through neural network, which includes 3 frames for angry feature, 4 for fear feature, 4 for happy feature and 2 for neutral or normal feature. Two parameters epochs and errors are considered for the detection of expression by lip, eye and mouth feature. Epochs specify maximum number of iterations that are to be taken for the detection of feature.

Errors specify total number of errors that are encountered according to the epochs. The more the epochs are, less is the error. The less number of errors gives more accuracy and efficiency in detection of emotions from features. The acceptable range for errors is from 0.1 to 0.9 and the goal is to reduce the errors up to 0.01 for more accuracy. For more clarity reports and graphs are shown which includes both of the two parameters.

## B. *Training of Neural Network Results*

TABLE VI.         TOTAL TRAINING FRAMES

| Frame | Total Frames | Classes | | | |
|---|---|---|---|---|---|
| | | Angry | Fear | Happy | Normal |
| 1 | 21 | 3 | 4 | 4 | 2 |

TABLE VII.         TRAINING OF LIPS

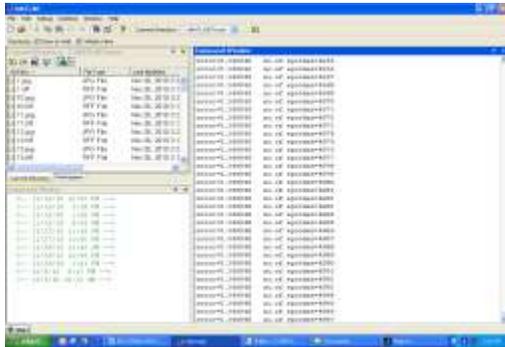| Max Epochs | Maximum Error | Minimum Error |
|---|---|---|
| 5000 | 0.45 | 0.3 |
| 10000 | 0.4 | 0.29 |
| 15000 | 0.5 | 0.2 |



Figure7.  Report of Epochs versus Error for Lips



Figure8. Graph of Epochs versus Error for Lips

TABLE VIII.         TRAINING OF EYES

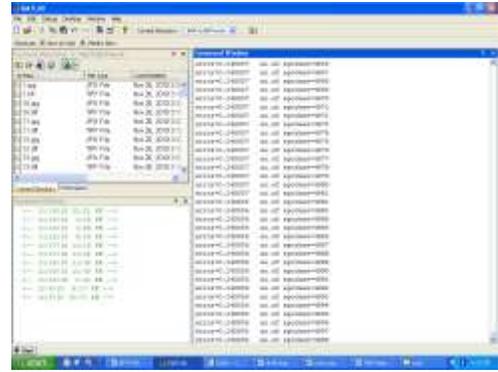| Max Epochs | Maximum Error | Minimum Error |
|---|---|---|
| 5000 | 0.55 | 0.28 |
| 10000 | 0.6 | 0.2 |
| 150000 | 0.7 | 0.1 |



Figure9.  Report of Epochs and Errors for Eyes

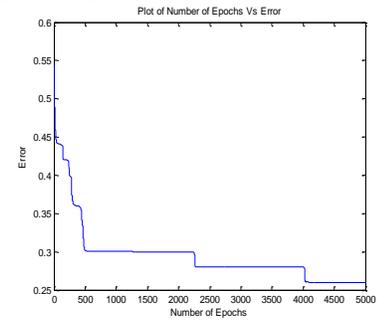

Figure 10 Graph of Epochs versus Error for Eyes

TABLE IX.         TRAINING OF EYES AND LIPS

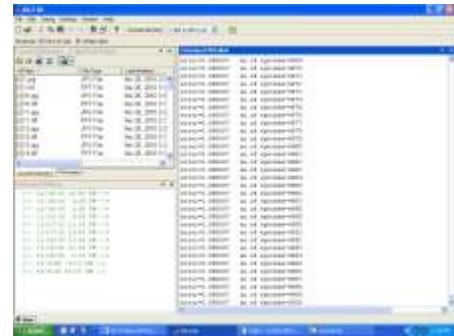| Max Epochs | Maximum Error | Minimum Error |
|---|---|---|
| 5000 | 0.65 | 0.08 |
| 10000 | 0.6 | 0.04 |
| 150000 | 0.5 | 0.02 |



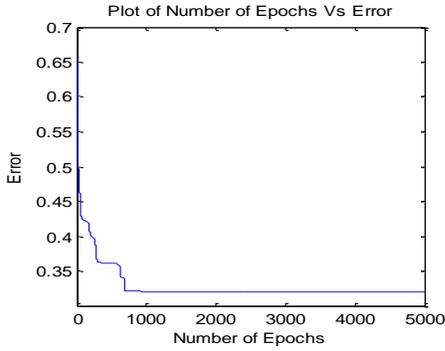Figure11.  Report of Epochs and Errors for Eyes

Figure 12 Graph of Epochs versus Error for Lips and Eyes

### 4.3 (Section 3) Testing of Frames

#### C. Testing of Frames

TABLE X.        TESTING OF LIPS

| Tested Frames | Correct Classification | Wrong Classification | Performance |
|---|---|---|---|
| 13 | 11 | 2 | 84 |

TABLE XI.        TESTING OF EYES

| Tested Frames | Correct Classification | Wrong Classification | Performance |
|---|---|---|---|
| 13 | 12 | 1 | 92 |

TABLE XII.        TESTING OF EYES AND LIP

| Tested Frames | Correct Classification | Wrong Classification | Performance |
|---|---|---|---|
| 13 | 12 | 1 | 92 |

### 4.4 (Section 4) Result of Classification



Angry                          Happy

Fig13 (a)                      Fig13 (b)



Fear                           Neutral

Fig13(C)                       Fig13 (d)

Figure 13        Result of Classification
(a)  Result of Classification of Angry Expression
(b)  Result of Classification of Happy Expression
(c)  Result of Classification of Fear Expression
(d)  Result of Classification of Normal or Neutral Expression

### 4.5 Comparative Analysis

The proposed technique is compared with existing ones [21] and gives more accuracy as compared to other methods. The existing methods include Combined Global and Local preserving features (CGLPF), Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA) and Locality Preserving Projection (LPP).

TABLE XIII.        COMPARISON OF OVERALL PERCENTAGE OF CORRECT RECOGNITION OBTAINED USING MUTATION BFO & AMF,CGLPF,PCA,LDA AND LPP

| Noise Details | Techniques | | | | |
|---|---|---|---|---|---|
| | Mutation BFO&AMF | CGLPF | PCA | LDA | LPP |
| Salt &Pepper Variance=0.05 | 95.9 | 95.8 | 66.9 | 94.725 | 74.525 |
| Salt &Pepper Variance=0.1 | 94 | 94.275 | 64.775 | 90.7 | 71.3 |
| Salt &Pepper Variance=0.15 | 93.32 | 92.825 | 60.075 | 80.125 | 67.275 |
| Salt &Pepper Variance=0.2 | 92.2 | 87.725 | 56.775 | 76.725 | 60.6 |

The number of correct recognized images for each query image in the prototype image set is calculated and results are shown in table 28.The proposed method Mutation BFO & AMF is compared with other existing techniques which is taken from Reference [21], which consists of CGLPF, PCA, LDA, LPP. Table 13 shows that for salt and pepper noise for variance from 0.05 to 0.2.



Figure 14        Comparitive Analysis through graph
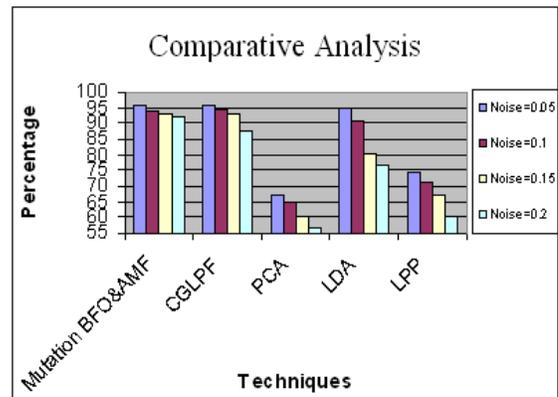
Graph clarifies the efficiency of proposed technique. Our proposed technique removes noise from variance up to 0.9,all other methods removes noise up to the variance of 0.2, means 90% of noise is removed with the accuracy approximately equal to 90% through our proposed method, so it performs better than other conventional techniques and it shows the high robustness of our algorithm.

## V. CONCLUSION

In this work Bacteria Foraging Optimization with mutation is used to remove highly corrupted salt and pepper noise with variance density up to 0.9.

The Bacteria Foraging Optimization with fixed step size requires more computation time with less accuracy, due to mutations the speed of BFO is increased by enhancing the accuracy in terms of quality of images.

This technique can be used as robust face emotion detection algorithm.

In this work a multiple feature options such as face, eyes and lips are used for emotion detection. The global, local features of facial expression recognition images can be independently selected through the mouse for identification for feature extraction.

The Radial Basis Neural Network requires less number of epochs as compares to other neural networks, therefore the proposed method is suitable for identification of emotions in the presence of salt and pepper noise as high as 90%.

Comparative analysis shows that the proposed technique is more efficient in recognizing expressions even under noisier environment.

## FUTURE WORK

Future work includes that the same technique can be used for detection of emotions in the presence of other noise such as speckle noise with adaptive median filter or by wiener filter.

Other neural network, which can be learning through optimization technique, can be used to improve overall significance of the system.

Replacing BFO with other less computational requirement tools improves the computation time. It can also be made as Graphical base system and a number of emotions should be considered more to make algorithm universal.

## REFERENCES

[1] Ekman, P., and Friesen, W., (1978) "Facial *Action Coding System: Investigator's Guide*", Consulting Psychologists Press.

[2] L. Beaurepaire, K.Chehdi, B.Vozel (1997), "Identification of the nature of the noise and estimation of its statistical parameters by analysis of local histograms", *Proceedings of ICASSP-97*, Munich, 1997.

[3] S. Spors, R. Rabenstein and N. Strobel(2001)," Joint Audio Video Object Tracking", IEEE International Conference on Image Processing (ICIP),Thessaloniki, Greece.

[4] Cootes, T., Edwards, G., and Taylor C., (2001) "Active *appearance models*". PAMI, Vol. 23, No. 6, pp. 681–685.

[5] Gukturk S B et.al (2001), A Data-Driven Model for Monocular Face Tracking," Eighth International Conference on Computer Vision (ICCV'01) - Volume 2

[6] Cootes, T. and Kittipanya-ngam, P., (2002) "*Comparing variations on the active appearance model algorithm.*" In BMVC, pp 837– 846.

[7] J. Ahlberg(2002), "An active model for facial feature tracking," *EURASIP Journal on Applied Signal Processing*, vol. 2002, no. 6, pp. 566–571,2002.

[8] Kevin. M.Passino (2002)"Biomimicry of Bacterial Foraging for Distributed Optimization and Control", IEEE Control System Magazine.

[9] G. Votsis, A. Drosopoulos, and S. Kollias(2003), "A modular approach to facial feature segmentation on real sequences," *Signal Processing: Image Communication*, vol. 18, no. 1, pp. 67–89, 2003.

[10] P.Y. Simard, H.S. Malvar (2004), "An efficient binary image activity detector based on connected components", *Proceedings of. IEEE International Conference on Acoustics, Speech, and Signal Processing*, 229–232, 2004.

[11] NoiseModels, *http://homepages.inf.ed.ac.uk/rbf/CVonline/LOCAL_COPIES/VELDHUIZEN/node11.html*

[12] Image Noise, http://en.wikipedia.org/wiki/Image_noise

[13] H.S.M. Al-Khaffaf , A.Z. Talib, R. Abdul Salam, "A Study on the effects of noise level, cleaning method, and vectorization software on the quality of vector data", Lecture Notes Computer Science 299-309.

[14] Kim, D.H., and Cho, C. H. (2005): Bacterial Foraging Based Neural Network Fuzzy Learning. IICAI 2005, 2030-2036.

[15] Mishra, S.: A hybrid least square-fuzzy bacterial foraging strategy for harmonic estimation. IEEE Trans. on Evolutionary Computation, vol. 9(1): 61-73, (2005).

[16] Tripathy, M., Mishra, S., Lai, L.L. and Zhang, Q.P.: Transmission Loss Reduction Based on FACTS and Bacteria Foraging Algorithm. PPSN, 222-231, (2006).

[17] Kim, D.H., Abraham, A., Cho, J.H.(2007)," A hybrid genetic algorithm and bacterial foraging approach for global optimization, Information Sciences, Vol. 177 (18), 3918-3937, (2007).

[18] Biswas, Dasgupta, Das, Abraham(2007)," Synergy of PSO and Bacterial Foraging Optimization –A Comparative Study on Numerical Benchmarks", innovations in Hybrid Intelligent Systems, ASC 44, pp. 255–263, 2007.Springer

[19] K.M.Bakwad, S.S.Pathnaik, B.S.Sohi, S.Devi, M.R.Lohakare (2009)," Hybrid Bacterial Foraging with parameter free PSO", Nature & Biologically Inspired Computing, 2009. NaBIC 2009.

[20] K. M. Bakwad, S.S. Pattnaik, B. S. Sohi, S. Devi1, B. K. Panigrahi, Sastry V. R. S. Gollapudi, (2009)," Bacterial Foraging Optimization Technique Cascaded with Adaptive Filter to Enhance Peak Signal to Noise Ratio from Single Image" IETE Journal of Research, Vol 55, Issue 4.

[21] K.M.Bakwad, S.S.Pathnaik, B.S.Sohi, S.Devi, M.R.Lohakare (2009)," Hybrid Bacterial Foraging with parameter free PSO", Nature & Biologically Inspired Computing, 2009. NaBIC 2009.

[22] Ruba Soundar Kathavarayan and Murugesan Karuppasamy, (2010) "Preserving Global and Local Features for Robust Face Recognition under Various Noisy Environments", International Journal of Image Processing (IJIP) Volume (3), Issue (6)

[23] G.Sofia, M. Mohamed Sathik, (2010)"An Iterative approach For Mouth Extraction In Facial Expression Recognition", Proceedings of the Int. Conf. on Information Science and Applications ICISA 2010 6 February 2010, Chennai, India.

[24] Jagdish Lal Raheja, Umesh Kumar, (2010)"Human Facial Expression Detection from detected in Captured image using Back Propagation Neural Network", International Journal of Computer Science and Information Technology (IJCSIT), Vol 2,No.1, February 2010.

[25] G.Sofia, M. Mohamed Sathik, (2010)"Extraction of Eyes for Facial Expression Identification of students", International Journal of Engineering Science and Technology, Vol.2 (7).

[26] Zhong Zhang, Qun Ding, Mingliang Liu, Hongan Ye (2010),"A Novel Facial Recognition Method", Journal of Communication and Communication and Computer, Vol.7, No.1.

[27] Stelios Krinidis, Ioannis Pitas (2010),"Statistical Analysis of Human Facial Expressions", Journal of Information Hiding and Multimedia Signal Processing, Vol.1, No.3.

AUTHORS PROFILE

**Renu Nagpal**  received diploma in Computer  Engineering  in 2002 from Technical Board,Chandigarh. B. Tech degree in Computer Science and Engineering in 2005 under Punjab Technical University,Jalandhar. Presently she is pursuing her M.Tech from Yadavindra College of Engineering,Guru Khashi Campus,Talwandi Sabo,Bathinda(Punjab) Her interests include, Image Processing, Swarm Intelligence, neural networks,bacteria foraging. She has contributed near about 10 technical papers in various national and international conferences. She is a life member of ISTE.

 **Sumeet Kaur** received her B.Tech in Computer Engineering from Sant Longowal Institute of Engineering & Technology(Deemed University) Punjab in 1999 and her M.Tech from Punjabi University, Patiala in 2007. She has more than 10 research papers in different national and international conferences. Currently, she is working as Assistant Professor in the Department of computer engineering, Yadavindra College of Engineering, Punjabi University Guru Kashi Campus, Talwandi Sabo, Punjab State, India. Her interest areas include encryption, network security, image processing and steganography.

**Pooja Nagpal**  received diploma in Computer  Engineering  in 2004 from Technical Board,Chandigarh. B. Tech degree in Computer Science and Engineering in 2007 under Punjab Technical University,Jalandhar.Since 2008 she is pursuing her M.Tech from Yadavindra College of Engineering,Guru Khashi Campus,Talwandi Sabo,Bathinda(Punjab) Her interests include, Image Processing and Swarm Intelligence .She has contributed  near about 8 technical papers in various national and international conferences. She is a life member of ISTE.

# Design Strategies for AODV Implementation in Linux

Ms. Prinima Gupta

MCA dept., Manav Rachna College of Engineering,
Sector-43, Faridabad. INDIA
prinima_mail@rediffmail.com

Dr. R. K Tuteja

MCA dept., N.C. Institute of Computer Sciences, Israna,
Panipat. INDIA
rk_tuteja2006@yahoo.co.in

*Abstract*—**In a Mobile Ad hoc Network (MANET), mobile nodes constructing a network, nodes may join and leave at any time, and the topology changes dynamically. Routing in a MANET is challenging because of the dynamic topology and the lack of an existing fixed infrastructure. In this paper, we explore the difficulties encountered in implementing MANET routing protocols in real operating systems, and study the common requirements imposed by MANET routing on the underlying operating system services. Also, it explains implementation techniques of the AODV protocol to determine the needed events, such as: Snooping, Kernel Modification, and Netfilter. In addition, this paper presents a discussion of the advantages as well as disadvantages of each implementation of this architecture in Linux.**

*Keywords- Ad-hoc Networking, AODV, MANET.*

## I. Introduction

AODV is an on demand algorithm, meaning that it builds routes between nodes only as desired by source nodes. It maintains these routes as long as they are needed by the sources. Hence, it is considered as a reactive routing protocol. The Ad-hoc On Demand Distance Vector Routing (AODV) protocol is an algorithm used for the implementation of such networks. The connection between nodes is established for the duration of one session, so no need to have a base station in order to establish such a connection between nodes. Nodes discover other target nodes that are out of range by broadcasting the network with Rout Requests (RREQ) that are forwarded by each node. If the destination node get the RREQ, then it sends back Route Reply (RREP) to the source node. After the route has been discovered between source node and destination node, then it's the time to start sending data thru that route.

There are a limited number of such implementations, mostly for the Linux operating system. This paper is an exploration and comparison of several AODV implementations including: National Institute of Science and Technology (NIST), the University of California, Santa Barbara (UCSB), Uppsala University (UU) in Sweden and the University of Illinois, Urbana-Champaign (UIUC). The earliest implementation of AODV is the Mad-hoc implementation.

In this paper we first give an overview of the challenges implementers face when implementing an on-demand ad hoc routing protocol. These challenges emerge, as the on-demand routing model does not easily fit into the standard operating system routing and packet forwarding model. We describe the problems with the routing model of current operating systems and identify the necessary extra events that must be recognized to ensure correct behaviour of on-demand ad hoc routing protocols. Then we describe and discuss the different design strategies that have previously been deployed in implementations of on-demand ad hoc routing protocols, focusing on AODV implementations in Linux. The intention is to give an overview of the developed solutions and point out best practices and experiences learnt.

## II. Background

First, we describe the AODV routing protocol and its basic operation. Then, we describe the available implementations of AODV for the Linux platform. It describes their available functionality and their design philosophy.

### A. AODV Protocol Overview

The AODV routing protocol is a reactive routing protocol; therefore, routes are determined only when needed. Hello messages may be used to detect and monitor links to neighbors. If Hello messages are used, each active node periodically broadcasts a Hello message that all its neighbors receive. Because nodes periodically send Hello messages, if a node fails to receive several Hello messages from a neighbor, a link break is detected. When a source has data to transmit to an unknown destination, it broadcasts a Route Request (RREQ) for that destination. When a RREQ is received by an intermediate node, a route to the source is created. If the receiving node has not received the RREQ before, is not the destination and does not have a current route to the destination, it rebroadcasts the RREQ. If the receiving node is the destination or has a current route to the destination, it generates a Route Reply (RREP). The RREP is unicast in a hop-by-hop fashion to the source. As the RREP propagates, each intermediate node creates a route to the destination. When the source receives the RREP, it records the route to the destination and begins sending data. If multiple RREPs are received by the source, the route with the shortest hop count is chosen.

As data flows from the source to the destination, each node along the route updates the timers associated with the routes to the source and destination, maintaining the routes in the routing table. If a route is not used for some period of time, a node

cannot be sure whether the route is still valid; consequently, the node removes the route from its routing table. If data is flowing and a link break is detected, a Route Error (RERR) message is sent to the source of the data in a hop-by-hop fashion. As the RERR propagates towards the source, each intermediate node invalidates routes to any unreachable destinations. When the source of the data receives the RERR, it invalidates the route and reinitiates route discovery, if necessary.

### B. AODV Implementations

Each implementation was developed and designed independently, but they all perform the same operations. The field studies using the AODV routing protocol have thus far been limited to devices running the Linux operating system, as the current implementations of AODV have all been developed for that platform. Thus real-world testing of the protocol has been limited to homogenous network environments. There are two types of different implementations, user space daemons and kernel modules.

*1) Mad-Hoc Implementation:* was the first available implementation of AODV by Fredrik Lilieblad, Oskar Mattsson, Petra Nylund, Dan Ouchterlony, and Anders Roxenhag running on a Linux 2.2 kernel but does not supports multicast. It uses the method of snooping ARP and data packets, using the libpcap Linux packet capturing facility. It is a user space only solution. It does not comply with an up-to-date version of the AODV specification, and is no longer supported. As such, it does not interoperate properly with the later implementations, and is not recommended for use.

*2) NIST Implementation:* was the only kernel implementation done by the NIST, Department of Commerce's Technology Administration U.S., and Wireless Communications Technologies Group running on a Linux with a 2.4 kernel. It is very fast and efficient reaching the best performance of all implementations. The NIST Implementation of AODV is currently at version 2.1 at time of writing. The latest version has support for multicast AODV, as well as multi-hop Internet gatewaying. The protocol is implemented completely as a Linux kernel module. If uses Netfilter from the 2.4 kernel to capture packets going in and out of the node instead of using the libcap library. It uses a Proc file to update the user about current routes and statistics for that node.

*3) Uppsala University Implementation*: The University of Uppsala also published a user space daemon implementation called AODV-UU which runs on Linux with a 2.4 kernel. Multicast support is available via a patch implemented by a group of researchers from the University of Maryland. The protocol is implemented as a user space daemon, and two loadable Linux kernel modules (kaodv and ip_queue_aodv). It uses the Netfilter library to intercept incoming and outgoing packets, but this is performed in user space.

*4) University California, Santa Barbara Implementation:* was the newest daemon published on 2nd of April 2002 by the University of California, Santa Barbara, running on a Linux with a 2.4 kernel. Similar to the UU implementation, the UCSB

version is implemented as a user space daemon. It similarly uses the Netfilter library for intercepting incoming and outgoing packets from the chosen interface. In fact, the implementation uses directly the UU packet input user space packet queuing module and the kaodv/packet_queue_aodv kernel modules. As such, it suffers from exactly the same problems as the UU implementation in that all packets must pass the boundary between kernel and user space twice.

*5) University of Illinois, Urbana-Champaign Implementation:* The UIUC implementation, is based on their ad-hoc support library (ASL), which is a Linux specific library designed to provide all the services required by ad-hoc routing protocols. As such, the UIUC AODV implementation is a user space daemon compiled against the ASL library. The implementation has not been interoperability tested against the other protocol implementations. Much of the complexity of the user space ad-hoc routing module has been removed to the ASL library, a very desirable feature, as this should allow other, different, ad-hoc routing protocols to be developed using the same library.

## III. IDENTIFYING THE CHALLENGES OF ON-DEMAND AD-HOC PROTOCOLS

When we are discussing the challenges faced when implementing an on-demand ad hoc routing protocol, it is of relevance to recap the routing architecture of current operating systems. In particular, how the functionality is divided and why implementing on-demand protocols is a challenge compared to implementing traditional routing protocols or proactive ad hoc routing protocols.

The routing functionality in modern operating systems is typically divided in two parts:

### A. Packet forwarding function

Consists of the routing function within the kernel, located within the IP layer of the TCP/IP stack, in which packets are directed to the appropriate outgoing network interfaces, or local applications, according to the entries in the kernel routing table. When the IP-layer receives a packet, it inspects a table called the forwarding table. Based on the IP destination address the packet is either directed to a local application listening on the specified port number, dropped, or sent out to the corresponding next-hop neighbor on the specified network interface according to the destination IP address of the packet.

### B. Packet routing function

It typically consists of a user-level program responsible for populating the kernel routing table. The definition of an optimal route is dependent on the routing algorithm; the number of hops to the destination is usually the chosen metric. The program performing the routing is typically implemented in user space as a program running in the background (the routing daemon).

On Linux the route selection process is carried out the following way: when selecting a route for a packet, the kernel

first searches the routing cache for an entry matching the destination IP address of the packet, and if found it forwards the packet to the next hop specified in the routing cache entry. Entries are deleted which have not been used for some time. If no entry for the destination is found in the routing cache, the kernel makes a look-up for the destination in the kernel routing (FIB) table using longest prefix matching. If an entry is found in the table, a new entry for the destination is created and inserted into the routing cache, i.e., the kernel routing table is used to populate the routing cache.

In on-demand routing protocols not all routes are known in advance, they must be discovered as they are needed. In such cases a mechanism is required to notify the on-demand routing protocol that a route discovery cycle must take place for the destination, and any packets already being sent to the destination must be queued while the route discovery cycle completes. For the AODV routing daemon to function, it must be determined when to trigger AODV protocol events. Since the IP stack was designed for static networks where link disconnections are infrequent and packet losses are unreported, most of these triggers are not readily available. Therefore, these events must be extrapolated and communicated to the routing daemon via other means. The events that must be determined are:

*1) When to initiate a RREQ:* Route Requests are needed when the IP layer receives a packet to be transmitted to an unknown destination, i.e., a destination with no matching entry in the route table. The problem of the current network stack architecture is that we only know we need a route after the packet has already crossed the boundary between user space and kernel space.

*2) When and how to buffer packets waiting for a route discovery cycle (or for some other reason) to complete:* When an application attempts to send a packet to a destination for which the routing table has not a valid route, the IP layer should buffer the packet for a period of time while a route discovery cycle takes place. If a route is found, the packets should be reinserted into the IP layer and sent to the destination. If a route is not found, the packets should be discarded and the application program should be notified.

*3) When to update the lifetime of an active route:* On-demand routing protocols typically cache a route that has been discovered for a period of time before deleting it if it is inactive. The IP layer therefore must have the capability to notify the routing protocol when an on-demand route has been used, so that the routing protocol can update its timers for the route.

*4) When to generate a route error message if a valid route does not exist for the next-hop IP address of a received packet:* If a data packet is received from another host and there is no valid route to the destination, the node must send a RERR so that the previous hops and the source stop transmitting data packets along this invalid route. Under normal operation of the IP layer on receiving a packet is to send a destination host unreachable ICMP message to the source of the transmission,

and silently drop the packet. Instead, the IP layer must give notification to the AODV routing protocol such that it knows it should send a route error message to the original source or the packet.

*5) When to generate a RERR during daemon restart:* When a node reboots, the AODV specification requires that it sends Route Error messages to any nodes attempting to communicate with it up until the end of DELETE_PERIOD seconds. This behavior is required in order to ensure no routing loops occur.

These notification and capabilities are not explicitly present in the protocol stacks of modern operating system. The existing implementations have taken a number of different approaches to solving this problem. The next section describes a number of possible approaches that implementers have taken in Linux.

## IV. DESIGN STRATERGIES OF LINUX

Special emphasis is given on implementation techniques for Linux. We highlight the advantages and disadvantages of each technique and discuss how the implementation techniques and the techniques available for programmers have evolved. The alternatives described in this section are:

- *Kernel Modifications:* Modify the source code of an operating system kernel to produce new API for implementation.

- *Snooping:* Use packet capturing facilities.

- *Netfilter:* Use the packet filter and packet mangling architecture.

The Linux netfilter framework can also partly be attributed to the fact that Linux has become the most popular platform for on-demand ad hoc routing protocol implementation development of the thirteen listed implementations on the AODV web page, eight are for Linux only.

### A. Kernel Modification

Another possibility to determine the AODV events is to modify the networking code of an operating system kernel. Royer and Perkins modified the Linux kernel to support their implementation of the AODV protocol. Such an API would require mechanisms as they would require protocols to register interest with the kernel in the relevant routing events (such as the requirement for a new Route Request). The kernel would then inform the ad-hoc routing protocol when a route Request is required; to initiate route discovery, code is added in the kernel at the point where route lookup failures occur. The kernel source was modified so that a route look-up failure would result in a notification to a user space daemon that was a part of the implementation, it would provide a mechanism to buffer packets for which a route request is being performed and to later reinject them; it would maintain timers associated with the route, etc. Figure 1 shows the architecture of the AODV daemon and the required support logic.

Advantages

- The events are explicitly determined and there are no wasted overhead.

- By modifying the kernel data structures and support code directly, there is no overhead of additional protocol accounting, compared to a user space implementation or even a Linux kernel module.
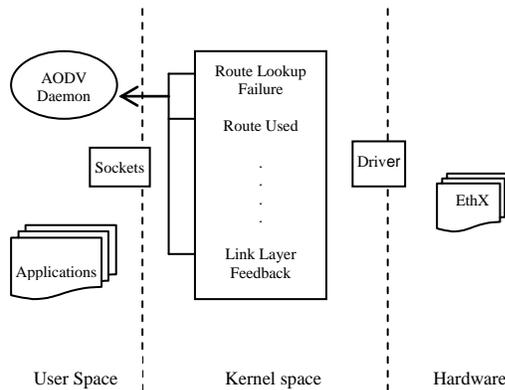


Figure 1.   Kernel Modification Architecture

Disadvantages

- *Difficult Installation procedure*: Installation of the necessary kernel modifications requires a complete kernel recompilation.

- *Less Portable:* A drawback of this approach is that it will require major changes to the operating system kernels, and will not be very portable for existing operating system kernels without requiring users to install a new kernel.

- *Difficult to maintain:* Patches (modifications) might only apply cleanly against a certain version of the Linux kernel. There could even be problems with kernels with the same version number as distributions apply their own set of patches to the Linux kernel source.

The first release of the AODV-UCSB implementation used the kernel modifications approach. Desilva and Das also made an implementation of AODV by modifying the Linux kernel and the in-kernel ARP implementation.

### B.  Snooping

Using code built into the kernel of most operating systems, a user space program can capture all incoming and outgoing packets on a network interface. The process of capturing packets is also known as sniffing or snooping. The code to perform snooping is built into the kernel and is available to user-space programs by using the Packet Capture Library (libpcap). Each packet that is transmitted is passed to the routing daemon using libpcap. When the daemon sees that a packet was transmitted along an active route, the lifetime for that route is updated so that it does not expire, since it is in use. In a similar manner, all the other AODV events may be determined by monitoring the incoming and outgoing packets. By snooping the Address Resolution Protocol (ARP) packets and data packets, AODV can be implemented without any kernel modifications. As such, the routing protocol can be implemented easily in either kernel space or user space. The routing protocol can determine when a route discovery cycle is needed by snooping ARP request packets, as an ARP request is sent to resolve the hardware address for an unknown IP address (if there is an appropriate subnet route entry set up for the correct interface).

The routing protocol can observe incoming and outgoing data packets, and as such can determine when a route is being used, or when a packet is received for which we have no routing information.

Advantage

- *Simple installing and execution:* It does not require any code to run in the kernel-space.

Disadvantages

- *Overhead:* An ARP packet is generated when a node does not know the MAC address of the next hop. Using this inference, if an ARP request packet is seen for an unknown destination and it is originated by the local host, and then a route discovery needs to be initiated. Since route discovery is initiated by outgoing ARP packets, these outgoing packets are unnecessary overhead, and they waste bandwidth.

- *Dependence on ARP:* If the routing table and ARP cache become out of sync, it is possible that the routing protocol may not function properly. For example, if the ARP cache contains an entry for a particular unknown destination, then an ARP packet will not be generated for this destination even though the destination is not known by the routing daemon. Consequently, route discovery will not be initiated. For proper operation the routing protocol must monitor and control the ARP cache in addition to the IP routing table, because disagreement between the two can cause the routing protocol to function incorrectly.

### C.  Netfilter

Netfilter is a packet filtering framework implemented as a set of hooks at well defined places in the Linux TCP/IP networking stack. Netfilter redirects packet flow through user defined code, which can examine, drop, discard, modify or queue the packets for the user-space daemon. Netfilter is similar to the snooping method; however, it does not have the disadvantage of unnecessary overhead or dependence on ARP.

It consists of a number of hooks in the IP layers that are well-defined points in a packet's traversal of the protocol stack. The IPV4 stack has five hooks. The two hooks NF_IP_LOCAL_OUT and NF_IP_LOCAL_IN are for packets incoming to and outgoing from local processes on the current host. Here a routing decision on what to do with the packet is made. If it is an incoming packet, it may be sent to the NF_IP_FORWARD hook before forwarding, sent up to the

NF_IP_LOCAL_IN hook for delivery to a local process, or dropped. If it is an outgoing packet, it is dropped or sent on the NF_IP_POST_ROUTING hook before being released to the appropriate network interface driver for transmission across the network. Incoming packets also traverse the NF_IP_PRE_ROUTING hook as they enter the IP layer, before being subjected to kernel routing. Thus, packets going from and to other hosts can be captured at these two hooks. Routing decisions are made for packets arriving at the network interface of the host after traversing NF_IP_PRE_ROUTING, to see if they are bound for this host or destined to be forwarded. Routing decisions are made for packets sent by local processes after traversing NF_IP_POST_ROUTING.
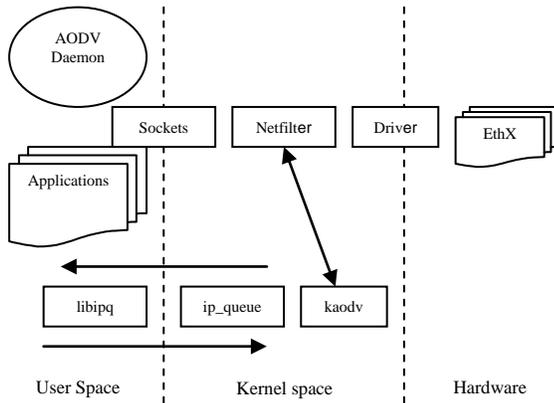


Figure 2.    Netfilter Architecture

These functions can show the direction packets travel through the network stack as they enter from a local process or a network interface.

- NF_ACCEPT: allow the packet to pass to the next registered hook

- NF_DROP: have the choice to either discard the packet

- NF_QUEUE: request that these packets be queued for later reinsertion into the IP layer otherwise this is equal to NF_DROP

- NF_STOLEN: grab the packet. We may reinsert the packet at a later point in time

- NF_REPEAT: call this hook again

These hooks are used by the kaodv kernel module. The kernel module driver, ip_queue module is used to queue these packets for the user-space daemon. There the AODV daemon uses a user space library libipq to make control decisions about each packet. Finally packets that are queued (by returning NF_QUEUE) are buffered by the ip_queue driver, typically (though not necessarily) for user space, see figure 2. These packets are handled asynchronously and thus they can be returned to the IP layer at any later time, or discarded.

The Netfilter architecture can be used for firewall filtering (the Linux ip_tables tool uses Netfilter in version 1.4 of the kernel), all kinds of Network Address Translation (NAT) services, or for other advanced packet processing requirements.

Advantages

- *Highly portable*

- *Easy to install*

Disadvantage

- *Requires a kernel module*: A kernel module is more portable than a kernel modification because it depends only on the Netfilter interface. This interface does not change from one kernel version to the next.

## V. CONCLUSION

AODV is currently one of the most popular ad-hoc routing protocols. These indicate that AODV performs very well both during high mobility and high network traffic load, making it one of the most interesting candidates among today's ad-hoc routing protocols. Implementing a routing protocol is very important to validate its design. Coming up with a clean implementation not only helps better understanding of the protocol nuances, but also allows extensions to explore the protocol design space. In this paper we analysed design possibilities for AODV implementations. We then examined the advantages and disadvantages of three strategies for determining this information. This analysis supported our decision to use small kernel modules with a user-space daemon. We hope that the information in this paper aids researchers in understanding the trade-offs in ad hoc routing protocol implementation development. Further, the description of the design structure and performance of each implementation can assist users in deciding which implementation best fits their needs.

### REFERENCES

[1]   Luke Klein-Berndt, NIST Kernel AODV homepage, http://w3.antd.nist.gov/wctg/aodv_kernel/. September 2003.

[2]   Mad-hoc AODV homepage. http://mad-ho.flyinglinux.net/. September 2003.

[3]   Luke Klein-Berndt, "Kernel AODV". National Institute of Standards and Technology. http://w3.antd.nist.gov/wctg/aodvkernel. 30 Oct 2008.

[4]   Ian Chakeres, UCSB AODV homepage. http://moment.cs.ucsb.edu/AODV/aodv. html, September 2003.

[5]   Erik Nordström, UU AODV homepage. http://user.it.uu.se/~henrikl/aodv/, September 2003.

[6]   Binita Gupta, UIUC AODV homepage. Including ASL library. http://sourceforge.net/projects/aslib/, September 2003.

[7]   E. M. Belding-Royer, "Report on the AODV Interop," University of California Santa Barbara, Tech. Rep. 2002-18, June 2003.

[8]   E. Borgia, "Experimental evaluation of ad hoc routing protocols," in Proc. of IEEE PerCom 2005 Workshops, Kauai Island, Hawaii, March, 8–12 2005.

[9]   "Running AODV-UU in the Network Simulator NS-2.". https://prj.tzi.org/repos/ dmn/aodv-uu-dtn/trunk/ README.ns. 2 Nov 2008

[10]  Erik.Nordstr¨om., AODV-UU. ttp://core.it.uu.se/core/index.php/AODV-UU. Last accessed December 2006.

[11] Ian D. Chakeres and Elizabeth M. Belding-Royer. "AODV routing protocol implementation design". In ICDCSW '04: Proceedings of the 24th International Conference on Distributed Computing Systems Workshops - W7: EC (ICDCSW'04), pages 698–703, Washington, DC, USA, 2004. IEEE.

[12] Douglas E. Comer., "Internetworking with TCP/IP: Principles, Protocols, and Architecture". Prentice-Hall, Inc., Upper Saddle River, NJ, USA, fourth edition, 2000.

[13] Saman Desilva and Samir R. Das., Experimental evaluation of a wireless ad hoc network. In Proceedings of the 9th Int. Conf. on Computer Communications and Networks (IC3N), pages 528–534, Las Vegas, NV, USA, October 2000.

[14] Nova Engineering, "NovaRoam," http://www.novaroam.com/.

[15] C. E. Perkins and E. M. Royer, "The Ad hoc On-Demand Distance Vector Protocol," in *Ad hoc Networking*, C. E. Perkins, Ed. Addison-Wesley, 2000, pp. 173–219.

[16] C. E. Perkins, E. M. Belding-Royer, and S. Das, "Ad hoc On-Demand Distance Vector (AODV) Routing," *RFC 3561*, July 2003.

[17] E.M. Royer & C.E Perkins, "An Implementation Study of the AODV Routing Protocol", Proceedings of the IEEE Wireless Communications and Networking Conference, Chicago, IL, September 2000.

[18] Netfilter homepage. http://www.netfilter.org/. September 2003.

[19] J. Kadlecsik, H. Welte, J. Morris, M. Boucher, and R. Russell, "The netfilter/iptables Project," http://www.netfilter.org/.

[20] IEEE Computer Society, "IEEE 802.11 Standard, IEEE Standard For Information Technology," 1999.

[21] J. Tourrilhes, "Wireless Tools for Linux,". http://www.hpl.hp.com/personal/Jean Tourrilhes/Linux/Tools.html.

[22] H. Lundgren, D. Lundberg, J. Nielsen, E. Nordstrm, and C. F. Tschudin, "A Large-scale Testbed for Reproducible Ad hoc Protocol Evaluations," in IEEE Wireless Communications and Networking Conference 2002 (WCNC), March 2002.

[23] V. Kawadia, Y. Zhang, and B. Gupta, "System Services for Implementing Ad-Hoc Routing: Architecture, Implementation and Experiences," in Proceedings of the International Conference on Mobile Systems, Applications, and Services (MobiSys), San Francisco, CA, June 2003, pp. 99–112.

[24] G. Bianchi, "Performance Analysis of the IEEE 802.11 Distributed Coordination Function," IEEE Journal Selected Areas in Communications, vol. 18, March 2000.

AUTHORS PROFILE

**Prinima Gupta** is presently working as Asst. Professor in MCA Department, Manav Rachna College of Engineering, Faridabad. She has completed Master of Computer Application from Kurukshetra University, Kurukshetra, M.Phil from Vinayaka Mission University, Tamil Nadu and presently doing PhD in Computer Science. She has 5+ years of teaching experience. She published 03 papers in National conferences. Her area of specialization includes Computer Networks and Computer Architecture



**Prof. (Dr.) R. K Tuteja** is presently working as Director (Academics) in NCICS, Israna, Panipat. He has 45 years of teaching experience. He was successfully guided 30 PhD research students and 17 students for M. Phil. Degree. He has published 126 Research papers in National/International Journals. He has worked as Head of Statistics/ Mathematics/ Computers Science & Application Department at M. D. University Rohtak.

# Model Based Test Case Prioritization for Testing Component Dependency in CBSD Using UML Sequence Diagram

Arup Abhinna Acharya
School of Computer  Engineering
KIIT University   Bhubaneswar,
India
aacharyafcs@kiit.ac.in

Durga Prasad Mohapatra
Department of Computer Science
& Engineering
National Institute of Technology
Rourkela, India
durga@nitrkl.ac.in

Namita Panda
School of Computer Engineering
KIIT University
Bhubaneswar, India
npandafcs@kiit.ac.in

*Abstract*—**Software maintenance is an important and costly activity of the software development lifecycle. To ensure proper maintenance the software undergoes regression testing. It is very inefficient to re execute every test case in regression testing for small changes. Hence test case prioritization is a technique to schedule the test case in an order that maximizes some objective function. A variety of objective functions are applicable, one such function involves rate of fault detection - a measure of how quickly faults are detected within the testing process. Early fault detection can provide a faster feedback generating a scope for debuggers to carry out their task at an early stage. In this paper we propose a method to prioritize the test cases for testing component dependency in a Component Based Software Development (CBSD) environment using Greedy Approach. An Object Interaction Graph (OIG) is being generated from the UML sequence diagrams for interdependent components. The OIG is traversed to calculate the total number of inter component object interactions and intra component object interactions. Depending upon the number of interactions the objective function is calculated and the test cases are ordered accordingly. This technique is applied to components developed in Java for a software system and found to be very effective in early fault detection as compared to non-prioritize approach.**

*Keywords- Regression Testing, Object Interaction Graph, Test Cases, CBSD*

## I. INTRODUCTION

Nowadays software development is quality oriented development. Quality can be ensured by very good testing techniques. So optimizing time and cost of testing process is really a challenge for test engineers. Regression testing is a kind of testing which requires maximum effort, time and cost. In fact, it might be hard to run the whole application unattended and to simulate any asynchronous input (e.g., interactive inputs) the application may receive. In such cases, regression testing can last days or weeks and can involve substantial human effort. Hence a technique like

Test case prioritization has to be devised which will lead to early fault detection.

Test case prioritization aims at finding an execution order for the test cases which maximizes a given objective function. Among the others, the most important prioritization objective is probably discovering  faults as early as possible that is, maximizing the rate of fault detection. In fact, early feedback about faults allows anticipating the costly activities of debugging and corrective maintenance, with a related economical return. When the time necessary to execute all test cases is long, prioritizing them so as to discover most faults early might save substantial time, since bug fixing can start earlier.

The major challenges in CBSD are testing component dependency. CBSD uses the reusable components as the building blocks for constructing the complex software system (component based system). Component based system promotes the software quality and productive. This building block approach has been increasingly adopted for software development, especially for large-scale software systems. A component based software often consists of a set of self contained and loosely coupled components allowing plug and- play. The components may be implemented by using different programming languages, executed in various operational platforms distributed across geographic distances; some components may be developed in-house, while others may be the third party off-the-shelf components of which the source code may not be available to the developers. So the cost of maintaining the component based software is comparatively more than the maintenance of conventional software system. So when we want to modify or add a component and apply the regression testing, it incurs more cost and time. So to reduce these two factors we use a test prioritization technique which is based on a criterion like maximum interactions between the components performed due to a test case during component interaction. The test case having maximum interactions given higher priority and executed first so that the debugger

will not sit idle as a result fault will be detected early. In this paper for describing each component we have taken the help of sequence diagrams, then a Object Interaction Graph (OIG) from sequence diagrams is constructed which shows the interrelation among the components. A new test prioritization algorithm is presented which is applied on OIG to count the maximum number of inter component interactions and intra component interactions made by the test cases.

Previous work on test case prioritization [1, 2, 3, 4, 5] is based on the computation of a prioritization index, which determines the ordering of the test cases (e.g., by decreasing values of the index). For example, the coverage level achieved by each test case was used as a prioritization index [3]. Another example is a fault proneness index computed from a set of software metrics for the functions exercised by each test case [1].

P.R. Srivastava [18] suggested prioritizing test cases according to the criterion of increased APFD(Average percentage of Faults detected) value. He proposed a new algorithm which could be able to calculate the average number of faults found per minute by a test case and using this value sorts the test cases in decreasing order. He also determined the effectiveness of prioritized test case(more APFD value) compared to non-prioritized test case(less APFD value). G. Rothermel et. al. [19] have described several techniques for test case prioritization and empirically examined their relative abilities to improve how quickly faults can be detected by those suites. Here more importance is given to coverage based prioritization. The authors applied these techniques to the base version of a program rather than the modified version of a program, hence these techniques are otherwise known as"general prioritization techniques". The objective is to detect faults as early as possible so that the debugger will not sit idle. B. Korel et.al.[9] proposed a new prioritization technique to prioritize the test cases by using several model-based test case prioritization heuristics. Model-based test prioritization methods use the information about the system model and its behaviour to prioritize the test suite for system retesting. An experimental study has been conducted to investigate the effectiveness of those methods with respect to early fault detection. The results from the experiment suggest that system models may improve the effectiveness of test prioritization. The prioritization techniques so proposed are used in traditional software retesting, but in this work we try to use the prioritization techniques in component-based software retesting.

The test case prioritization methods can be categorized in to code-based testing and model based technique. In the code based test prioritization, source code of the system is used to prioritize the test cases. Most of the test prioritization methods [6, 10, 11, 12, 13, 14] are code based. In several test prioritization criteria were presented and their influence on the improvement of the rate of fault detection was investigated.

In model-based test prioritization [7,9] a system's model(s) is used to prioritize tests. System modelling is widely used to model state-based systems, e.g., real time systems. System models are used to capture some aspects of the system behaviour. One type of model-based test prioritization methods [7,8] are appropriate for modifications that involve changes in the model and then in the source code. The second type of model-based test prioritization methods [9] are appropriate for modifications that do not involve any changes in models (changes are only made in the source code). In this paper, we have used UML 2.0 for modelling to concentrate on the second type of model-based prioritization method.

Though several priotization techniques have been proposed previously, but the interdependency issues present in component composition in CBSD, while finding the prioritized test suit, has not been taken care of. Regression testing mainly involves testing the changes occurred in software due to addition of new components. During component composition in CBSD, the inter component dependency leads to lot of errors. So the authors have taken in to consideration the above criteria while generating the prioritized test suit to increase the APFD.

The rest of the paper is organized as follows: Section II describes the problem statement for prioritization along with a brief introduction to CBSD. The proposed model along with a case study and a comparative study are described in Section III and Section IV .The paper concludes in Section V. with the discussion on continuing work in this direction in Section VI. Due to space constraints, this paper does not include descriptions of a system model such as notations and their semantics. Interested readers are referred to any UML book such as [16] or UML manual published by OMG [15].

## II. PROBLEM STATEMENT

In CBSD Component interface is defined as, it is the only way that a component communicates with the external environment. There are two kinds of interface: service providing and service required. When the services are provided by an interface it is called service providing interface and when the interface of a component requiring a service it is called service required interface. All components should be plug-compatible i.e a service required interface can be connected to a service providing interface. We have defined a Component as follows: Component C = ( P, R) , where P=$P_1$, $P_2$, $P_n$ is the set of providing services interface,

R = $R_1$, $R_2$, $R_m$ is the set of required services interface. The providing and required services of a component C is denoted by C.P and C.R respectively and C.P∩C.R=∅. [17]

There are two kinds of special components, one is the component without the required services, the other is the one without the providing services for other components. According to the fact the numbers of two kinds of components can be one or more.

In the Fig.1 the required services of *C1P* C2 are the union of C1.R1 and C2.R2 with the remove of satisfied services in S. With the definition of composition the providing and required services are propagated to the interface of composed component, so the composition could be carried parallel. A Component interaction graph (OIG) is used to describe the interrelation of components. A complete component interaction graph (OIG) makes the testing quite easy. A OIG is a directed graph where OIG = (V, E), V represents a set of nodes. V = VI U VC, VI is the set of interface nodes and VC is the set of component nodes, E represents the set of directed edges.

The interface is denoted by an ellipse and a component with dashed square. The interaction among components can be gained from the OIG directly.

There are two kinds of special components, one is the component without the required services, the other is the one without the providing services for other components. According to the fact the numbers of two kinds of components can be one or more.

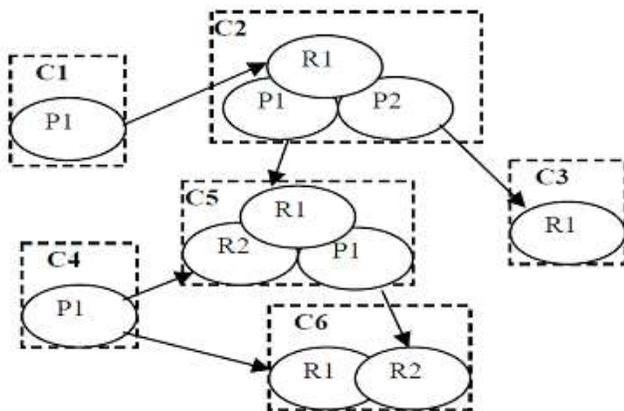The OIG illustration is given in figure 1:



Fig. 1. Object Interaction Graph(OIG)

If there is an existing edge form C1.P1 to C2.R1 in the CIG it means the required service R1 of C2 has been satisfied by the providing service P1 of C1, which is C2.R1= C1.P1.

Practically it is not possible to perform rigorous testing. Tester has to select subset of test cases from the original test suite. This makes test case selection quite challenging. This selected regression suite should cover all the functionality i.e. adequate functional coverage and greater fault exposing potential. Due to squeezed test schedule, testing team may not able to execute all test cases from the selected regression suite. Sequencing of test cases based on some criteria helps testing team to achieve the goals whilst reducing testing cycles. Rothermel at el. [3] defines the test case prioritization problem as follows: **Given:** T, a testsuite; PT, the set of permutations of T; f, a function from PT to the real numbers.

**Problem:** Find T' belongs to PT such that (for all T") (T" belongs to PT) (T"≠ T') [f (T')≥ f(T")].

Here, PT represents the set of all possible prioritizations (orderings) of T and f is a function that, applied to any such ordering, yields an award value for that ordering [3]. The objective of this research is to develop a test case prioritization technique that prioritizes test cases on the basis of detection of fault rate.

### III. PROPOSED MODEL

To facilitate regression testing by optimizing the time and cost, we propose a method to prioritize the test cases by using model based prioritization method by extracting the benefits of Unified Modelling Language(UML). UML provides lifecycle support in software development and is widely used to describe analysis and design specifications of software. It is a big challenge to study the test case generation from UML diagram. In case of a Object Oriented System Design (OOSD), each component is represented by collection of objects. Due to encapsulation, the only way objects can communicate is through message passing. Whenever an event occurs, it is executed through a sequence of occurrence of message passing. We have used sequence diagram from the set of diagrams present in UML 2.0. As Sequence diagram represents various object interactions through message passing, it can act as an input to the proposed model. We are generating an Object Interaction Graph (OIG) from the sequence diagrams present. The methodology we have used for generating the graph has been discussed in Section III(A) Further in Section III(B) we have discussed how to traverse the OIG to calculate the number of inter component object interaction and intra component object interaction. Section III(B) describes about objective function evaluation and the prioritization technique.
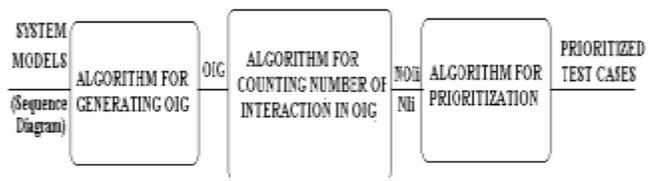


Fig. 2. A Frame Work For Generating Prioritized Test Cases

### A. Generating OIG form System Models

We have used sequence diagram for system modelling. The object interactions can be very well identified using a sequence diagram. During regression testing any modification in the code will have no effect on the sequence diagram. The object interaction can be categorized into two different types. One of them is intra component object interactions and the other one is inter component object interactions. In case of intra component object interaction, the interaction between objects present within a component is considered where as in case of inter component object interaction we consider the object interactions present

between two different components. A sequence diagrams in UML are used to model how an object communicate with other objects in its life time i.e. it is used to capture the dynamic behaviour of a system. The basic elements of a sequence diagram are object s and messages [Booch, Rumbaugh and Jacobson 1998] and it shows a state machine which emphasizes the flow of control from state to state.

A Object Interaction Graph (OIG) is used to describe the interrelation of components. A complete object interaction graph (OIG) makes the testing quite easy. A OIG is a directed graph where OIG = (V, E), V represents a set of nodes. For generating Object Interaction Graph (OIG), each object present in the sequence diagram is represented as a node in the graph. The intra component object interactions form the edges of the graph and represented in BLACK color. The inter component object interactions form the edges of the graph and represented in RED color.

**Algorithm: GENERATE OIG**
**Input:** Sequence Diagrams of various components of the system representing message passing between objects
**Output:** Object Interaction Graph (OIG)// It is a directed graph
1. Initialize OIG to be empty
2. for i=1 to n//n is the total number of objects
3. Add a node $N_i$ to OIG $==N_i$ represents i$^{th}$ node. Object shared by different components treated as a single node.
4. for i=1 to n
5. for j=1 to n
6. for each incoming message from object $O_i$ to $O_j$ ==All guard conditions are ignored
7. if (interaction type==intra)Establish an edge between $O_i$ to $O_j$ (i.e. $N_i$ and $N_j$) and color it as "BLACK" as well as append the pre and post conditions.
8. Else Establish an edge between $O_i$ to $O_j$ (i.e. $N_i$ and $N_j$ )and color it as "RED" as well as append the pre and post conditions.
9. The possible start and end of the scenario sequences are represented with solid arrows.

*B. Traversing OIG*

When the OIG is generated from the system models, it has to be traversed to count the number of inter component and intra component object interactions. NO$I_i$ represents the number of Object Interactions discovered by test case ti with in one component of the software and N$I_i$ represents the number of Object Interactions discovered by test case ti between two different components of the software. We follow the depth first search (DFS) methodology for traversing the graph. The type of interaction is decided depending upon the color of the edge in the graph. If the edge color is found to be "BLACK", it represents an intra component object interaction, where as edges colored as "RED" represents inter component object interaction

**Algorithm**: IN_CALCULATE
**Input:** Test case *ti* & Object Interaction Graph (OIG)
**Output:** $NOI_i$ and $NI_i$
1. Initialize both $NOI_i$ and $NI_i$ to 0.
2. Traverse each interaction in the OIG for *ti* in DFS
3. if (edge color =='BLACK' && current edge is not visited already)
4. $NOIi$ + + ==Increment the value for intra component interaction
5. Else
6. $NI_i$ + + == Increment the value for inter component interaction
7. Return $NOI_i$ and $NI_i$

*C. Generating Prioritized Test Cases*

Once we get the value for NIi and NOIi by using the algorithm described in Section III(B), prioritization process starts. For each test case ti, the value of NIi and NOIi are added. We have considered the total number of intra component interaction where as the total number of inter component object interactions is found out by multiplying it with RP i.e. total number of providing service interface and required service interface . If the faults due to component integration are detected early, it will give a better coverage.

The added result is divided with unit time U to determine value of the objective function i.e. factor criteria FCi. We try to maximize the objective function using a Greedy approach.

**Algorithm: TEST_PRI**
**Input:** Regression Test Suite T
**Output:** Prioritize Test Suite T'
1. Traverse the test suite T, for each test case *ti* present, call **IN_CALCULATE** (*$t_i$*) to calculate $NOI_i$ and $NI_i$
2. Define some unit time U
3. Calculate objective function ( $FC_i$) for test case *ti* as $FC_i= (NOI_i+\textbf{RP}*NI_i)/U.$          *(1)*
   // **RP represents total number of providing service interface**.
4. Generate T' by Sorting the test suit T in ascending order of $FC_i$ for each $t_i$ .
5. Store T' in the test case repository for regression testing.

IV. CASE STUDY- A CELLULAR NETWORK MANAGER

We have taken the case study of a Cellular Network Manager to explain the proposed model. We have taken into consideration two components i.e."Dialing a Phone" and "Cellular Network Connection". From the sequence diagram of both the components given in Fig.3 and Fig.4, corresponding OIG are designed as given in Fig.5.
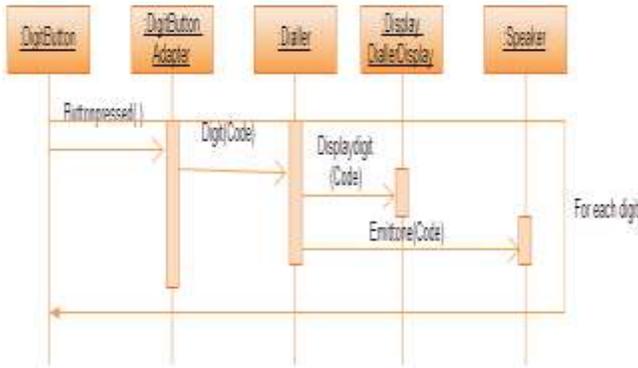
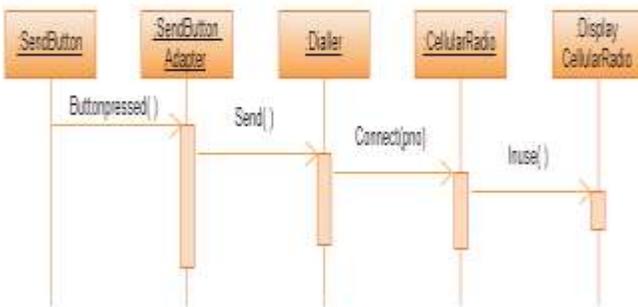Fig. 3. Sequence Diagram for dialing the number



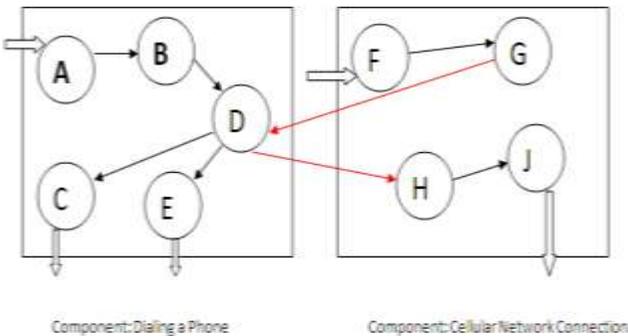Fig. 4. Sequence Diagram for cellular phone connection



Fig. 5. OIG for a Cellular Network Manager

A: Digit Button               F: Send Button
B: Digit Button Adapter   G: Send Button Adapter **(P)**
C: Dialer Display            H: Cellular Radio**(R)**
D: Dialer **(Both P&R)**     J: Cellular Radio Display
E: Speaker

Three test cases are considered to test the prioritization algorithm. The test cases are designed to test the Dialer Display($t_1$), to test the Speaker($t_2$) and to test the Cellular Radio Display($t_3$). The following table contains the value of $NOI_i$, $NI_i$ and $FC_i$. Here the unit time U is considered to be 1 unit.

TABLE I: *OBJECTIVE FUNCTION ($FC_i$) EVALUATION*

| Test Cases | $NOI_i$ | $NI_i$ | $FC_i$ |
|---|---|---|---|
| $t_1$ | 3 | 0 | 3 |
| $t_2$ | 3 | 0 | 3 |
| $t_3$ | 2 | 4 | 6 |

From the table I we conclude that the prioritized test sequence is: **$t_3$, $t_2$, $t_1$ or $t_3$, $t_1$, $t_2$**

The proposed model found to be very effective as it increases the Average Percentage of Fault Detection (APFD) when it is compared with generalized model based method and few code based methods like LOC count and Function count. The comparison made is summarized in Table-II.

**TABLE II**
A COMPARATIVE STUDY

| Name of Prioritized Technique | Approximate Increase in APFD value(%) |
|---|---|
| Code based Approach (LOC count, Function count etc.) | 30 |
| Model based Approach | 35 |
| Model Based Approach using the Dependency Criteria in CBSD | 45 |

## V. CONCLUSION

The cost and time required for regression testing can be minimized by using the prioritization technique discussed in this paper. Here we have proposed a model based prioritization method by considering the number of Object Interactions per unit time as the objective function. Here more importance is given to number of inter component object interactions present because maximum faults are expected to be present when components interact with each other. The proposed model found to be very effective as it increases the Average Percentage of Fault Detection (APFD) when it is applied to few of the projects developed in Java by java 45%-50%. This approach is mainly applicable to test the component composition in case of component based software maintenance.

## VI. CONTINUING WORK

The proposed method can further be extended to prioritize test cases to perform regression testing for real time systems and distributed systems. Here the authors prioritize the test case using a model based approach. The authors are also working on adding new criterion like frequency of data base access number of state changes in UML state chart diagram etc. Two different modelling diagrams can also be integrated to find criterion to generate test cases Requirement specifications can also be used to prioritize the test cases. Test case prioritization for

concurrent systems is also a very challenging area of research due to its dynamic behaviour

## REFERENCES

[1] S. Elbaum, A. Malishevsky, and G. Rothermel.*Test case prioritization: A family of empirical studies.*, IEEE Transactions on Software Engineering, 28(2):159-182, February 2002.

[2] J. M. Kim and A. A. Porter.*A history-based test prioritization technique for regression testing in resource constrained environments.*, In Proceedings of the International Conference on Software Engineering (ICSE), pages 119-129. ACM Press, May 2002.

[3] G. Rothermel, R. Untch, C. Chu, and M. J. Harrold*Test case prioritization.*, IEEE Transactions on Software Engineering, 27(10):929-948, October 2001.

[4] H. Srikanth, L. Williams, and J. Osborne.*System test case prioritization of new and regression test cases.*, In Proceedings of the 4th International Symposium on Empirical Software Engineering (ISESE), pages 62-71. IEEE Computer Society, November 2005.

[5] A. Srivastava and J. Thiagarajan.*Effectively prioritizing tests in development environment*, In Proceedings of the International Symposium on Software Testing and Analysis (ISSTA), pages 97-106. ACM Press, July 2002.

[6] J. Kim, A. Porter, *"A History-Based Test Prioritization Technique for Regression Testing in Resource Constraint Environments,"*, Proc. 24th International Conference on Software Engineering, pp. 119-129, 2002.

[7] B. Korel, L. Tahat, M. Harman, *"Test prioritization Using System Models"*, 21st IEEE International Conference Software Maintenance (ICSM '05), pp. 559-568, 2005.

[8] B. Korel, G. Koutsogiannakis, L. Tahat,*"Model-Based Test Prioritization Heuristic Methods and Their Evaluation"*, 3rd ACM Workshop on Advances in Model Based Testing, A-MOST, 2007.

[9] B. Korel, G. Koutsogiannakis, L. Tahat,*"Application of System Models in Regression Test Suite Prioritization,"*Proc. 24st IEEE International Conference Software Maintenance (ICSM '08), pp. 247-256, 2008.

[10] Z. Li, M. Harman, R. Hierons,*"Search Algorithms for Regression Test Case Prioritization,"*IEEE Transactions on Software Engineering, vol. 33, No. 4, pp. 225-237, 2007.

[11] G. Rothermel, R. Untch, C. Chu, M. Harrold,*"Test Case Prioritization: An Empirical Study,"*Proc. IEEE International Conference on Software Maintenance, pp. 179-188, 1999.

[12] G. Rothermel, R. Untch, M. Harrold,*"Prioritizing Test Cases For Regression Testing,"*IEEE Transactions on Software Engineering, vol. 27, No. 10, pp. 929-948, 2001.

[13] A. Srivastava, J. Thiagarajan,*"Effectively Prioritizing Tests in Development Environment,"*Proc. ACM International Symposium on Software Testing and Analysis, ISSTA-02, pp. 97- 106, 2002.

[14] W. Wong, J. Horgan, S. London, H. Agrawal,*"A Study of Effective Regression Testing in Practice,"*Proc. International Symposium on Software Reliability Eng., pp. 230-238, 1997.

[15] UML 2.0 Reference Manual, Object Management Group, 2003.

[16] Schneider, G., and winters, J.P., Applying Use Cases, Second edition, Addison Wesley,2001.

[17] Arup Abhinna Acharya, Sisir Kumar Jena, *" Component Interaction Graph: A new approach to test component composition"*,Journal of Computer Science and Engineering, Volume 1, Issue 1, May 2010.

[18] P. R. Srivastava, *"Test Case Prioritization"*, Journal of Theoritical And Applied Information Technology 2008 JATIT

[19] G. Rothermel, R. H. Untch, C. Chu ,M. J. Harrold *"Test Case Prioritization:An Emperical Study"*, in Proceedings of the 24th IEEE International Conference Software Maintenance (ICSM '1999) Oxford, U.K,September,1999 .

[20] GB. Korel, G. Koutsogiannakis, *"Experimental Comparsion of Code Based and Model model Based Test prioritization "*, IEEE 2009.

## AUTHORS PROFILE

**Arup Abhinna Acharya** is an Assistant Professor and research scholar in the School of Computer Engineering, KIIT University, Bhubaneswar, Odisha, INDIA. He received his Masters degree from KIIT University Bhubaneswar. His research areas include Object Oriented Software Testing, Software Cost Estimation, and Data mining. Many publications are there to his credit in many International and National level journal and proceedings. He is having eight years of teaching experience. He is a member of ISTE. He can be reached at aacharyafcs@kiit.ac.in.

**Durga Prasad Mohapatra** received his Masters degree from National Institute of Technology, Rourkela, India. He has received his Ph.D. from Indian Institute of Technology, Kharagpur, India. He is currently working as an Associate Professor at National Institute of Technology, Rourkela. His special fields of interest include Software Engineering, Discrete Mathematical Structure, Program Slicing and Distributed Computing. Many publications are there to his credit in many International and National level journal and proceedings. He is a member of IEEE. He can be reached at durga@nitrkl.ac.in.

**Namita Panda** is an Assistant Professor in the School of Computer Engineering, KIIT University, Bhubaneswar, Odisha, INDIA. She received her Master's degree from KIIT University Bhubaneswar. Her research areas include Object Oriented Software Testing, Parallel Processing and Computer Architecture. She has published papers in national and international level proceedings. She is having seven years of teaching experience. She is a member of ISTE. She can be reached at npandafcs@kiit.ac.in.

# EM Wave Transport 2D and 3D Investigations

Rajveer S Yaduvanshi

ECE Deparment

AIT, Govt of Delhi, India-110031

email:yaduvanshirs@yahoo.co.in

Harish Parthasarathy

ECE Department

NSIT, Govt of Delhi, India-110075

email–harishp@nsit.ac in

*Abstract*—**Boltzmann Transport Equation [1-2] has been modelled in close conjunction with Maxwell's equation, investigations for 2D and 3D transport carriers have been proposed. Exact solution of Boltzmann Equation still remains the core field of research. We have worked towards evaluation of 2D and 3D solutions of BTE. Application of our work can be extended to study the electromagnetic wave transport in upper atmosphere i.e. ionosphere. We have given theoretical and numerical analysis of probability density function, collision integral under various initial and final conditions. Modelling of coupled Boltzmann Maxwell's equation taking binary collision and multi species collision terms has been evaluated. Solutions of Electric Field (E) and Magnetic Field (B) under coupled conditions have been obtained. PDF convergences under the absence of electric field have been sketched, with an iterative approach and are shown in figure 1. Also 3D general algorithm for solution of BTE has been suggested.**

*Keywords- Boltzmann Transport Equation, Probability Distribution Function, Coupled BTE-Maxwell's.*

## I. INTRODUCTION

BTE is an integral differential equation used for characterizing carrier transport in semiconductor [5-6] and gases distribution in space f(x,v ,t). The Boltzmann equation applies to a quantity known as the distribution function, which describes this non-equilibrium state mathematically [1] and specifies how quickly and in what manner the state of the gas changes when the disturbing forces are varied. BTE shall compute average behaviour of the system in terms of distribution function of time. Evolution of distribution function is governed by Boltzmann Transport Equation. Boltzmann Transport Equation can be solved by mathematical and numerical techniques. Here distribution function can be a function of seven variables i.e. three physical space, three volume space and one time. It shall provide complete description of the state of gas. The distribution function carries information about the positions and velocities of the particles at any time. There can be two broad methods to solve the BTE. The first method consists in directly discretizing and solving the BTE using standard numerical methods [1-2] for differential equation. The second, called the Monte Carlo (MC) method, solves the BTE as being the stationary solution of a stochastic differential equation [3-4]. BTE structure due to high dimensionality is hard to solve.

One could describe a gas flow in classical physics by giving position and velocity of all molecules at any instant of time. It means information about flow of gas, electron and ions can be known by solving BTE with proper initial conditions and boundary conditions. The solution of BTE is PDF (3-4), which is a function of position and velocity of particles and the time variables. Distribution Function shall be derived by concentration of kinetic energy and moments due to applied force. Distribution function can be also characterized as current density.

A statistical approach [5-7] can be taken instead of defining position and velocity of each molecule. Using the construct of ensemble, a large number of independent systems evolving independently but under same dynamics, can be characterized by density function, which gives the probability that an ensemble member can be found in some elemental volume in phase space, which has been very well explained in reference [8-9].

In a state of equilibrium a gas of particles has uniform composition and constant temperature and density. If the gas is subjected to a temperature difference [8] or disturbed by externally applied electric, magnetic, or mechanical forces, it will be set in motion and the temperature, density, and composition may become functions of position and time, in other words, the gas moves out of equilibrium and specifies how quickly and in what manner state changes when disturbing forces are controlled. Equilibrium can be disturbed by temperature change, external force, magnetic force and mechanical force.

Here we have proposed three different models i.e. when BTE is at equilibrium state, second when force is applied but no field presents and third condition is when both force and field exists simultaneously. We have developed BTE formulations for two dimensional and three dimensional solutions. We have studied transport parameters i.e. charge density, current density, magnetic potential and electric potential, electric field and magnetic field. General purpose algorithms for 3D analysis have also been developed. This developed modeling theory can be very useful for tolerance analysis in chip designing in microelectronics. We have solved coupled equations by finite difference numerical method.

We have organized our paper in five sections. Section 1 gives us the concept of BTE. Section 2 describes BTE formulations [10]. Section 3 presents modeling and simulations of BTE. Section 4 deals with multi species collision Section 5

presents coupled modeling of EM transport. Section 6 concludes the paper.

**Abbreviations used in the text**

v=velocity of particles

f = force on particles

m= mass of particle

x=displacement

BTE=Boltzmann Transport equation

PDF= Probability Density Function

EM= Electromagnetic wave.

## II. BTE FORMULATIONS

Here distribution function can be a function of seven variables for 3D model i.e. three physical space, three volume space and one time. It shall provide complete description of the state of gas. There can be two broad methods to solve the BTE

Standard numerical methods for differential equation.

The Monte Carlo (MC) method solves the BTE as being the stationary solution of a stochastic differential equation.

Here we shall work for numerical solution for all investigations.

**BTE** under Steady State Conditions

$$\frac{\partial f}{\partial t} + v\frac{\partial f}{\partial x} + \frac{F}{m}\frac{\partial f}{\partial v} = 0 \qquad (1)$$

BTE under Equilibrium conditions [2]

$$f(\vec{x}, \vec{v}, t)\left(\frac{\partial}{\partial t} + \vec{v}\,\nabla x + \frac{F}{m}\,\Delta v\right) = 0 \qquad (2)$$

2 D representation of BTE

$$\frac{\delta f}{\delta x}(x, Vx) = -Vx\,\frac{\delta f(x,\,Vx)}{\delta x} \qquad (3)$$

3 D representation of BTE under RTA

$$\frac{\partial f}{\partial t} + v\frac{\partial f}{\partial x} + \frac{F}{m}\frac{\partial f}{\partial v} = -\frac{f-fc}{\tau} \qquad (4)$$

BTE when collision term is accounted for

$$\frac{\partial f}{\partial t} + v\frac{\partial f}{\partial x} + \frac{F}{m}\frac{\partial f}{\partial v} = \left(\frac{\partial f}{\partial t}\right)coll \qquad (5)$$

We have discretized the above 2d and 3 differential equations developed and subjected for numerical solution. The output is the PDF w. r .t position and velocity at different times. The PDF is normalized, where intensity of the image represents its relative value. Solution of BTE under constant EM field and force applied is computed as presented in fig1, it is assumed that there is no collision and initial function is Gaussian in nature. Given below are the 2Dand 3D computed transport models.

$$\frac{df}{dx}(x, Vx, t) + Vx\frac{df}{dx} + \frac{fx}{m}\frac{df}{dVx} = 0$$

$$\frac{df}{dx}(x, Vx, t) + Vx\frac{df}{dx} + \frac{e.Eo}{m}\frac{df}{dVx} = 0$$

Now for six space dimensions seven variables are considered. The distribution function f(x, v, t) for 3D represents x as position vector and v as velocity vector.

$$\frac{df}{dt}(x, y, z, Vx, Vy, Vz, t)$$
$$+ Vx\frac{df}{dx} + Vy\frac{df}{dy} + Vz\frac{df}{dz}$$
$$+ \frac{e}{m}(Ex + Vy\,Bz - Vz\,By)\frac{df}{dVx}$$
$$+ \frac{e}{m}(Ey + Vx\,Bx - Vx\,Bz)\frac{df}{dVy}$$
$$+ \frac{e}{m}(Ez + Vx\,By - Vy\,Bx)\frac{df}{dVz}$$
$$= 0 \qquad (6)$$

Here we have to find $\frac{df}{dt}$, $\frac{df}{dx}$ and $\frac{df}{dv}$. Here of intensity of image represents normalized magnitude of PDF. These plots have been obtained at different time and convergence at equilibrium is presented with initial field as Gaussian.
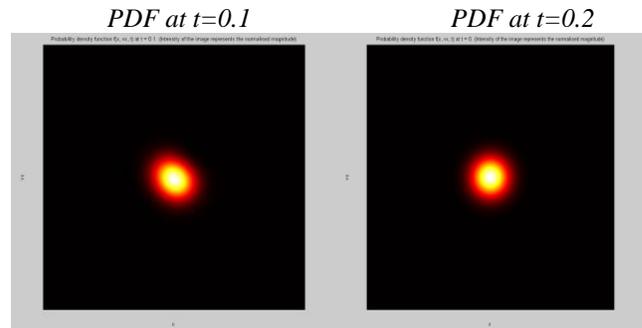
### A. 2D PLOTS

PDF at t=0.1                PDF at t=0.2



Figure. 1(a-b) Intensity of Image

PDF at t=0.3                PDF at t=0.4
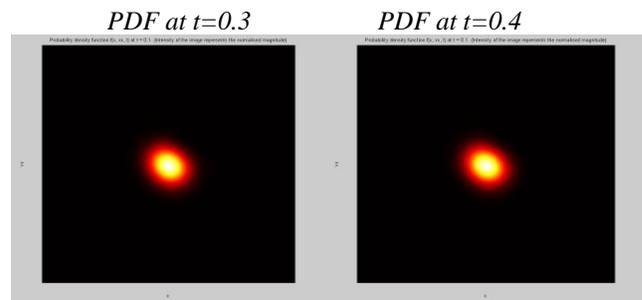


Figure. 1(c-d) intensity of Image

PDF at t=0.7                PDF at t=0.6
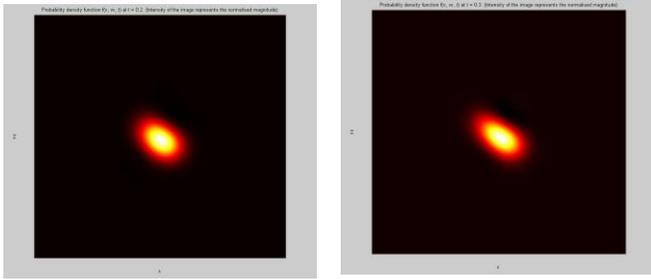
B.   3D   PLOTS



Figure.1 (e-f) Intensity of Image

*PDF at t=0.7*          *PDF at t=0.8*
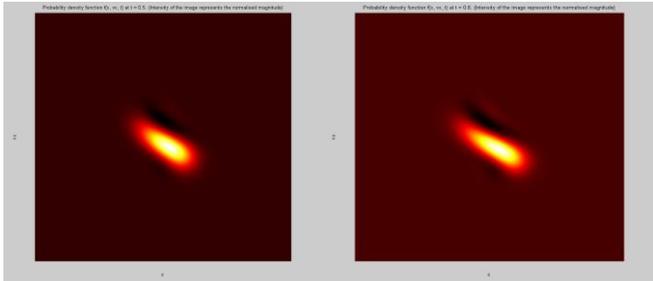


Figure. 1(g-h) Intensity of Image

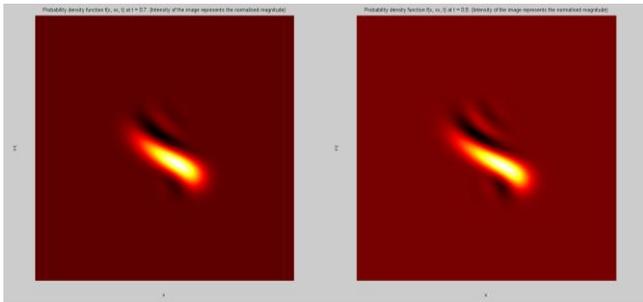*PDF at t=0.9*          *PDF at t=1.0*



Figure.1 (i-J) Intensity of Image
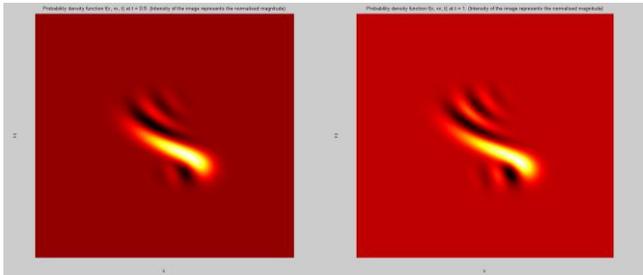
*PDF at t=1.11*        *PDF at t=1.12*
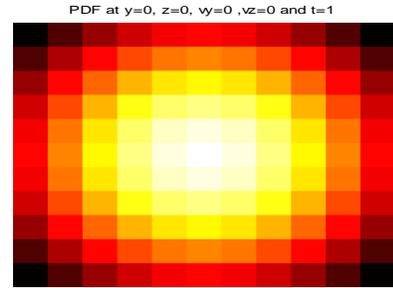


Figure. 1 (l-m) Intensity of Image



Figure. 2(a) Intensity of Image



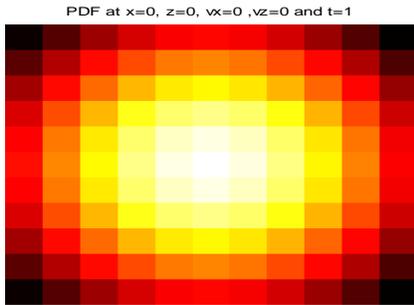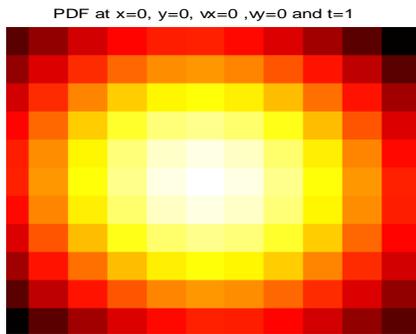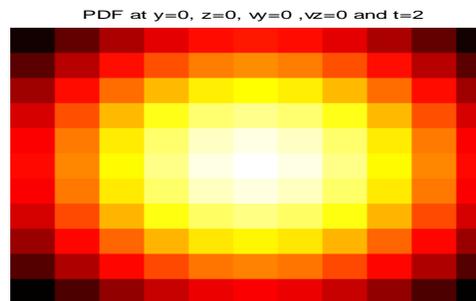Figure 2(b) Intensity of Image



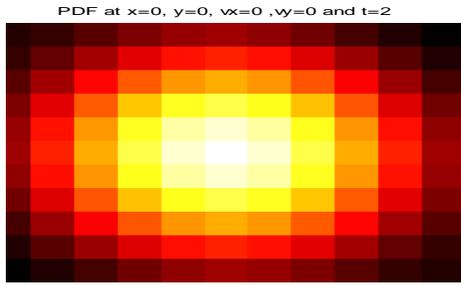Figure 2(c) Intensity of Image



Figure 2(d) Intensity of Image

PDF at x=0, y=0, vx=0 ,vy=0 and t=2

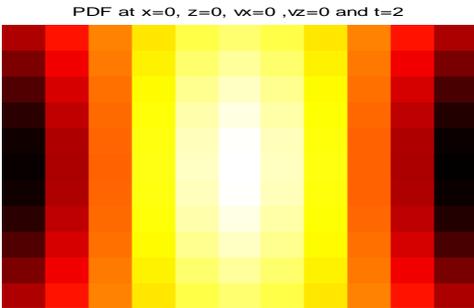Figure 2(e) Intensity of Image

PDF at x=0, z=0, vx=0 ,vz=0 and t=2

Figure 2(f) Intensity of Image

The ensemble is characterized by density function which gives probability that an ensemble member can be found in some elemental volume in phase space.

We have assumed two particles collision term. Equations (7- 8) presents collision terms. Here

$$\frac{df}{dt} + v\frac{df}{dx} + \frac{f}{m}\frac{df}{dv} =$$

$$\int \int v \int \Omega (f'f'1 - ff1) \, v \, \sigma \, d^2 \, \Omega \, d^3 v$$

**(7)**

Where

t= time

v= molecular velocity,

v1= Pre Collision

V'=post collision velocity

V'1=post collision velocity

F= external force

M= mass of the molecule

f1=f(x, v1, t)

$f'$=f (x, v1, t) here prime values represents post collision conditions due to conservation

V=v - v1

V =|V|=|V'|

Ω Solid angle is deflection angle of relative velocity.

σ collision cross section will depend on potential , scattering and velocity of particles. Collision cross section is defined by probability that a collision between two molecules will result in a given pair of post collision velocities. The integration over velocity v1 from -∞ to ∞ in all dimensions while integrating over Ω extends over unit sphere.

$$\left(\frac{\partial f}{\partial t}\right) \text{coll} =$$

$$\int v1 \int \Omega \, (f'\,f'1 - f\,f1) \, v \, \sigma \, d^2 \, \Omega \, d^3 v$$

(8)

Collision integral is five dimensional integral that must be evaluated for every point in physical space, every point in time and every point in velocity space. Collision cross section $\sigma$ will depend on intermolecular potential, pre and post collision velocities and scattering angle Ω [9-11]  , though it is constant for hard sphere molecules. $\sigma$  direction describes the probability density to a certain change of velocity. Collision integral evolution approach can be directly applied to solve high dimensional problems. The term $\int v1 \int \Omega(f'\,f'_1 - f\,f_1)$ v $\sigma$ d² Ω d³v  represents net rate at which molecular enter the point of interest in phase space due to collision and $\int v1 \int \Omega$ $(f'\,f'_1 - f\,f_1)$ v $\sigma$ d²Ω d³v  represents the net rate at which molecules are scattered out . Both terms are integrated over all possible pre collision velocity and all possible collision angles. We have also developed general purpose 3D algorithm for implementation as mentioned below:

**BTE 3DGeneral Purpose Algorithm:**

```
Clear all;
Clc;
Close all;
% for a variable space of 101^6, it might take hours to compute
the result as it may require large amount of memory.
len = 101;
Centre = round ((len+1)/2);
f = zeros (len, len, len, len, len, len);
Stdev = 5;
% initializes the prior PDF as a Gaussian distribution
% that the function can be differentiated properly, and yet not
extend to.
% the ends of variable space. For a variable size of len, even
larger.
% values can be considered.
For x = 1: len
 For y = 1: len
For z = 1: len
For vx = 1: len
For vy = 1: len
For vz = 1: len
 F(x, y, z, vx, vy, vz) =
 exp(-((x-centre)^2+(y-centre)^2+(z-centre)^2+(vx-
centre)^2+(vy-centre)^2+(vz-entre)^2)/(2*stdev^2));
 End
 End
```

```
End
End
End
End
df = zeros (len, len, len,len,len,len);
dx = 0.01;
dy = 0.01;
dz = 0.01;
dvx = 0.01;
dvy = 0.01;
dvz = 0.01;
e_m = 0.01;
ex = 1;
ey = 0;
ez = 0;
bx = 1;
by = 0;
bz = 0;
dt = 0.01;
tic;
for t = 1:2
for x = 1: len
 for y = 1: len
  for z = 1: len
   for vx = 1: len
    for vy = 1: len
     for vz = 1: len
 % df/dx
% First two are at boundaries, so special cases
% must be taken care of separately
 If (x == 1)
 df_dx = f (2, y, z, vx, vy, vz)/dx;
 else if (x == len)
  df_dx =
  -f (len-1, y, z, vx, vy, vz)/dx;
    else
    df_dx =
   (f(x+1,y,z,vx,vy,vz)-f(x-1,y,z,vx,vy,vz))/dx;
                end
                % df/dy
                if (y == 1)
                    df_dy =
f(x, 2,z,vx,vy,vz)/dy;
                else if (y == len)
                    df_dy =
 -f(x, len-1, z, vx, vy, vz)/dy;
                else
                    df_dy =
 (f(x,y+1,z,vx,vy,vz)-f(x,y-1,z,vx,vy,vz))/dy;
                end
                % df/dz
                if (z == 1)
                    df_dz =
f (x, y,2,vx,vy,vz)/dz;
                else if (z == len)
                    df_dz = -f (x,y,len-1,vx,vy,vz)/dz;
                else
```

```
                    df_dz=
(f(x,y,z+1,vx,vy,vz)-f(x,y,z-1,vx,vy,vz))/dz;
                end
                % df/dvx
                if (vx == 1)
                    df_dvx = f(x, y, z, 2, vy, vz)/dvx;
                else if (vx == len)
                    df_dvx =
 -f(x, y, z, len-1, vy, vz)/dvx;
                        else
                        df_dvx =
 (f(x,y,z,vx+1,vy,vz)-f(x,y,z,vx-1,vy,vz))/dvx;
                    end
                    % df/dvy
                    if (vy == 1)
                        df_dvy = f(x, y, z, vx, 2, vz)/dvy;
                    else if (vy == len)
                        df_dvy =
 -f(x, y, z, vx, len-1, vz)/dvy;
                    else
                        df_dvy =
 (f(x,y,z,vx,vy+1,vz)-f(x,y,z,vx,vy-1,vz))/dvy;
                    end

                    % df/dvz
                    if (vz == 1)
                        df_dvz =
   f (x,y,z,vx,vy,2)/dvz;
                    else if (vz == len)
                        df_dvz =
 -f(x, y, z, vx, vy, len-1)/dvz;
                    else
                        df_dvz =
 (f(x,y,z,vx,vy,vz+1)-f(x,y,z,vx,vy,vz-1))/dvz;
                    end
                    % df =
 (vx df/dx + vy df/dy + vz df/dz +
 (ex +   vy bz – vz by) df/dvx + (ey + vx bz – vz bx)
 df/dvy + (ez + vx by – vy bx) df/dvz)*dt
 % since v is cantered around centre; six is    subtracted from
all velocities
 df (x, y,z,vx,vy,vz) =
 (vx-centre)*df_dx+ (vy-centre)*df_dy+
(vz-centre)*df_dz +e_m*((ex+ (vy-centre)*bz-
(vz-centre)*by)*df_dvx+ (ey+ (vx-centre)*bz-
(vz-centre)*bx)*df_dvy+ (ez+ (vx-centre)*by-
(vy-centre)*bx)*df_dvz);
 end
 end
 end
 end
 end
 end
f = f+dt*df;
figure;
 subplot (1, 3,1);
```

```
imshow(mat2gray(f(1:len,1:len,centre,centre,centre,centre)));
 colormap hot;
   xlabel ('x');
   ylabel ('vx');
   title (['PDF at y=0, z=0, vy=0 ,vz=0 and t=' num2str(t)]);
   subplot (1, 3,2);
   g = f (centre, centre, 1: len, 1:len,centre,centre);
   imshow (mat2gray (reshape(g, [len len])));
   colormap hot;
   Xlabel ('y');
   Ylabel ('vy');
   Title (['PDF at x=0, z=0, vx=0, vz=0 and t=' num2str (t)]);
   Subplot (1, 3, 3);
   g = f (centre, centre, centre, centre, 1: len, 1: len);
   Imshow (mat2gray (reshape (g, [len len])));
   Colormap hot;
   Xlabel ('z');
   Ylabel ('vz');
   Title (['PDF at x=0, y=0, vx=0, vy=0 and t=' num2str (t)]);
end
Toc;
```

### III. MULTI SPECIES COLLISIONS

According to hydrodynamic theory of multi species diffusion in a gas is governed by given below equation (10) This problem can be worked for taking collision of similar charges , different charges and opposite charges for computing collision integral of multi species .E/N ratio can be evaluated to predict the outcome of collision for multi species. Here we have modeled for multi species collision integral which will give much better insight in understanding of transport physics of gases, here we assume collision of two species at any time because collision of three species is negligible .also as per authors knowledge no solutions in this regards have been evaluated for simulations of multi species collision term. The collision integral [1-4] which is necessary for computing the transport properties has been worked upon. Collision integral can be written as sum of contributions of particles interactions with long range attractive portion of intermolecular potential defined .Here we have modeled for multi species collision integral which will give much better insight in understanding of transport physics of gases. Here we assume collision of two species at any time because collision of three species is negligible. Also as per authors knowledge no such solution in this regards have been evaluated for simulations of multi species collision term. Flick diffusion is an approximate method for determining multispecies fluid dynamics .we can compare our model, for multispecies with Fick diffusion. Diffusion process SNR can be evaluated by Brownian motion. After few collisions noise starts to creep and contaminates the results. Noise growth is caused when particles frequency tends to infinity as well as to zero. Molecules process is determined by evaluating cross section evaluation. In our approach we have evaluated collision integral taking two different charge particles and only binary collisions are considered .The approach can be extended for computing collision integral

taking different particles into account and summed to evaluate resultant multi species collision integral. For the sake of simplicity no secondary electro emission from the walls is considered. Various parameters i.e. time, ion current, pressure, temperature, number of particles, field, energy distribution function, current, mean free path, E/N ratio, collision cross section of ion and gas molecules can be worked. Solving Boltzmann transport equation governing drift and diffusion of the particles rate of the arrival, time, spectrum etc can be predicted .Ionic flux density is given by Flick's law. Study of collision between electron .Collision kernels can be useful tool for carrying simulations. Relaxation time approximation (RTA) is simplification to BTE Integral .We can determine friction and diffusion coefficients that that are due to collisions. Here friction term represents drag or slowdown of particles due to collision and diffusion term producing a spreading of distribution function.

$$f_i (t, \vec{r}, \vec{v}) \text{ ,where } i = 1, 2, 3, 4\ldots\ldots\ldots N$$

$$m_i = \text{Mass}$$

$$e_i = \text{Charge}$$

Distribution function and derivation of collision integral of multi species taking non-reactive pairs into consideration has been worked and can be written as follows:

$$\frac{\partial f_i}{\partial x} + (v, \nabla r) f_i + \frac{e_i}{m_i} (E + \vec{V}\vec{B}) \nabla_v. f_i =$$

$$\sum_{j=1}^{N} I_{ij}$$

i species colliding with j species and evaluating collision integral as below :

$$I_{ij} =$$

$$\sum_j \int ( - f_i (t, \vec{r}, \vec{v}) . f_j (t, \vec{r}, \vec{v}) + f_i (t, \vec{r}, \vec{v}'). f_j(t, \vec{r}'_i, v_i ). \sigma_{ij} (v, v_i , v', v'_i ) \delta^3 v_i \qquad (9)$$

### 5. Coupled EM Transport

Flow of EM transport through plasma can be modelled as follows:

$$\frac{\partial f_i}{\partial x} + (v, \nabla_r) f + \frac{e_i}{m_i} (E + \vec{V}\vec{B}) \nabla_v. f$$

$$= -\frac{f-fc}{\tau} \qquad\qquad (10)$$

Equation (9) presents multi species domain and equation (10) gives coupled concepts. Electromagnetic field E,B has been coupled with f (x, v ,t) through Maxwell's system. As the upper part of the atmosphere consists shells of electrons and electrically charged atoms and molecules that surround earth from 50 K ms to 1000kms. Plasma contain ionized layer known ionosphere. Positive and negative ions are attracted because of EM force Positive and negative ions are attracted because of EM Force We can derive many other parameters like current density, flux, collision integral etc. Collision integral provides us rate of change of PDF. Coupled system can provide us solution without making approximations. Collision integral can

provide us contribution from particles iterations with long range communication. Power spectrum gives electron density measurement, ion and electron temp, mass, drift velocity. Flux is measurement of intensity of charge. Plasma is a ionized gas . It may responds strongly to EM Field. Aplasma is a collection of particles some positive charged and some negatively charged particles and few neutral. Locally a charge imbalance may exist, which may lead to net electric field in that reason. State of plasma is characterized by distribution function f(r,v,t) . Here we have simulated the above equation (14) by recursive method. We have solved for E & B. For solving these values we have assumed initial function as Gaussian and then evaluated V & A. For solution of A and V values first we need to compute the values of J (current density) and ρ (charge density). After computing all the above values we need to substitute all values obtained thus into the above mentioned master equation or coupled equation for solution by iterative method. For computing coupled solution we can also assume Fermi Dirac/ Gaussian as initial function.

Results of simulations are presented in fig 1,2,3 which gives us total insight of precision solution of Boltzmann Transport Equation. This type of solution approach have been unique of its nature so far as compare to previous one.

$$\vec{E}(\vec{r},t) = \overrightarrow{-\nabla}v(\vec{r},t) - \frac{d}{dt}\vec{A}(\vec{r},t)$$
$$\vec{B}(\vec{r},t) = \overrightarrow{\nabla}' \times \vec{A}'(\vec{r},t)$$

Her E and B are electric and magnetic fields 3D simulation results have been presented in fig 3.

$$V(\vec{r},t) = \frac{1}{4\pi\epsilon 0}\int \frac{e\left(\vec{r}'-t-\frac{|\vec{r},\;\vec{r}'|}{2}\right)}{|\vec{r},\;\vec{r}'|}d^3r'$$

$$\vec{A}(\vec{r},t) = \frac{1}{4\pi}\int \frac{\vec{J}(\vec{r}'-t-\frac{|\vec{r},\;\vec{r}'|}{2})}{|\vec{r},\;\vec{r}'|}d^3r'$$

Here A and V are electric and magnetic potentials Respectively simulation results are presented

$$\rho(\vec{r},t) \quad = \quad -e\int f(\vec{r},\vec{v},t)\,d^3v$$
$$\vec{J}(\vec{r},t) \quad = \quad -e\int \vec{V}(\vec{r},\vec{v},t)\,d^3v$$

Here ρ and J are charge density and current density and simulation results have depicted in fig 16,17 and fig 18,19 respectively. We shall evaluate the value of coupled solution.

$$\frac{\partial f(\vec{r},t)}{\partial t} + (\vec{V},\vec{\nabla})\,ff(\vec{r},\vec{V})$$
$$-\frac{e}{m}(\vec{E}+\vec{V}\times\vec{B},\nabla_{\mathbf{v}})$$
$$f(\vec{r},\vec{v},t) \quad = \quad \frac{f0-f}{\tau}$$
$$\left[\frac{\partial f}{\partial t}\right]_{coll} = \int\sigma(v,\;)\,|V-V'|(-f(\vec{r},\vec{v},t)+$$
$$f(\vec{r},\vec{v}',t)+f(\vec{r},\vec{v},t)f(\vec{r},\vec{v}',t)\,d^3V'\,Sin\theta d\theta$$

Taking initial function as $f_0(\vec{r},\vec{v},t)$ as Gaussian function. We have evaluated value of E, B , A , V, J and ρ. These values have been substituted in the coupled equation and simulated results are obtained for coupled equation.
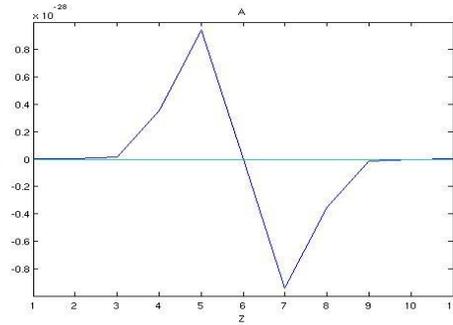


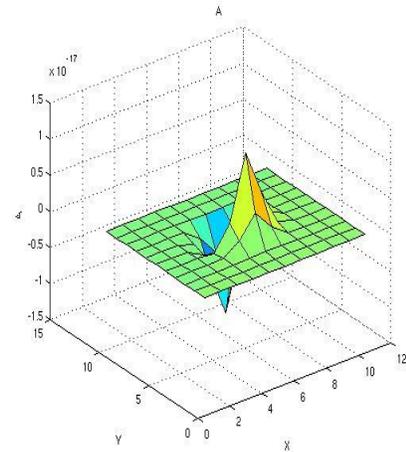Fig.3(a)    A vs Z  at specific value x, y, t mid points



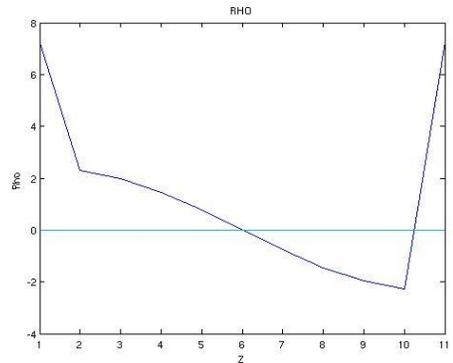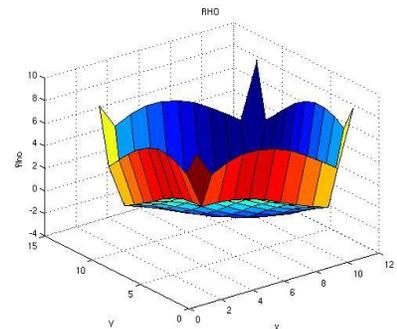Fig3(b)    A vs X,Y at specific value z, t mid points



Fig.3(c )    $\rho(\vec{r},t)$ charge density rho vs. Z at specific mid points x, y, t.

Fig . 3(d)          Rho vs X,Y at specific value z , t  mid points
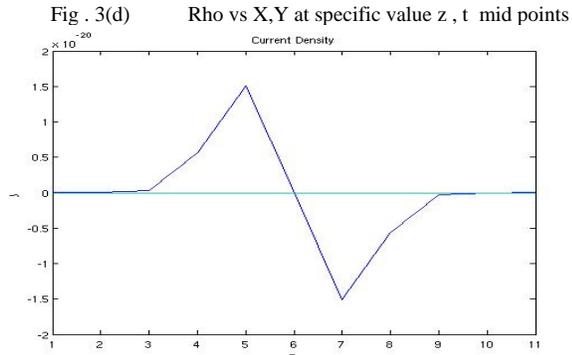


Fig .3( e)  J at specific value x ,y ,t  mid points  w.r.t to Z.
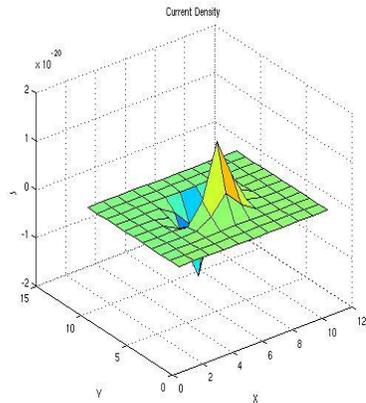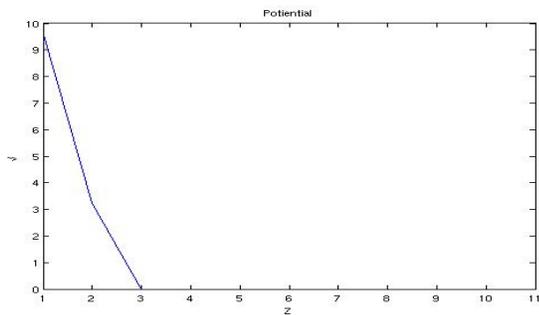


Fig.3(f)  J  vs  X,Y at specific value z, t mid points .



Fig( 3g)  V vs Z  at specific value x , y ,t  mid points
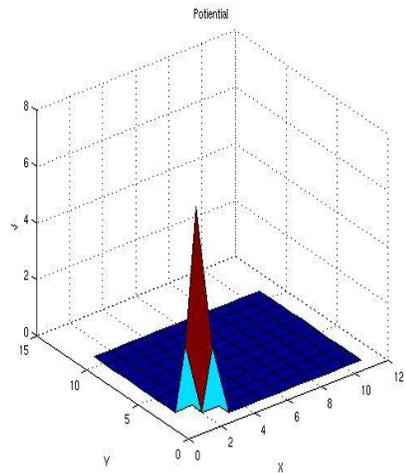


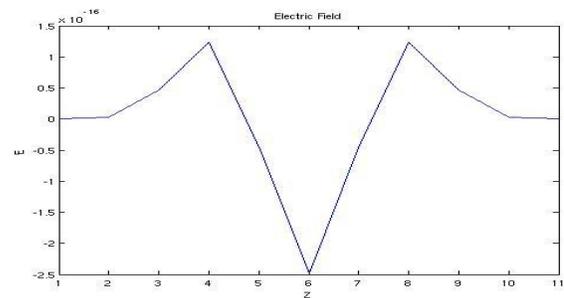Fig .3(h) V  vs  X,Y at specific value z, t mid points



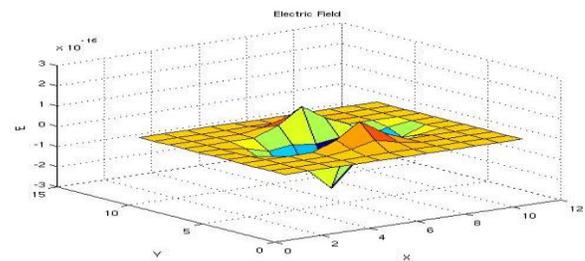Fig .3(i) E vs Z at specific value x, y, z mid points



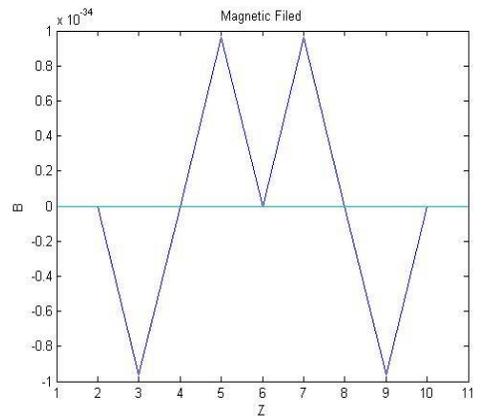Fig.3(j) E vs X,Y  at specific value z,t mid points



Fig 3(k) B vs Z at specific value x , y , t  mid points

Fig3( l) B vs X,Y at specific value z, t mid points



Fig.3(m) PDF vs Vx, Z.



Fig.3(n)  PDF  vs Vy, Vz



Fig.3(o)        PDF vs  Vy, Vz

## IV.    CONCLUSION

Investigations for transport parameters have been worked by means of modeling of Boltzmann transport equation for both linear and nonlinear responses. Collision Integral solution for binary and multi species have been realized and simulated. General algorithms for 3D transport parameters solution have been proposed. 2D transport with and without field have been worked. Multi   species terms have been realized. In addition we have proposed an efficient solution to coupled Boltzmann-Maxwell's equations.

Convergence solution for 2D model have been evaluated. Coupled results seems to feature  pulse nature and presented in results of coupled solutions. This feature with  controlled pulse can be used for fast switching.

Modeling of coupled process is  difficult  because of involvement of two different physics which require time and space analysis. Parameters that control the shape and amplitude of the pulse have been modeled and simulated.

Mismatch of functions having various size of matrix need critical   computations.   Different   dimensions   due interdependency of many parameters need to be resolved during code development for 3D analysis in iterative simulation. Application of our work can be in chip designing in microelectronics, Plasma antenna designs and study of upper atmosphere characteristics.

### REFERENCES

[1]    Statistics Mechanics, Kerson Huang, Wiley 2008

[2]    Bahadir, A.R .and T. Abbasov (2005) "Numerical investigation of the liquid flow velocity over an infinity plate which is taking place in a magnetic field" International journal of applied electromagnetic and mechanics.

[3]         EM Lifshitz and LD Landau Electrodynamics of continuous media"butterworth-Hienemann

[4]   Fermigier, M. (1999), "Hydrodynamique physique" Problèmes résolus avec rappels de cours, collectionsciences sup Physique, edition   Dunod.

[5]    Wait, R. (1979) "The numerical solution of algebraic equations" A Wiley-interscience publication.

[6]    Chorin, A.J., 1986, "Numerical solution of the Navier-Stokes equations", *Math. Comp.*, Vol. 22, pp. 745-762.

[7]    EM Lifshitz and LD Landau   "Fluid Mechanics" Vol 6, Butterworth-Heinemann.

[8]    Holman, J.P., 1990, "Heat Transfer" (7[th] edition), Chapter 12, Special Topics in Heat Transfer, MacGraw-Hill Publishing Company, New York.

[9]    Guermond, J.L. & Shen, J., 2003, "Velocity-correction projection methods for incompressible flows", *SIAM J. Nume. Anal.* Vol. 41, No. 1, pp. 112-134.

[10]  Ramesh   Garg   "Analytical   and   Computational   methods   in electromagnetic" Artech House, London

[11]   JD Jackson, "Classical Electrodynamics" third edition, Wiley.

AUTHOR'S PROFILE

**Rajveer S Yaduvanshi** has been working as Asst Professor in ECE department of AIT, Government of Delhi, Delhi -110031. He has 21 years of teaching experience. He is Fellow of IETE. Author has visited France for radar moderanization program representing India. His research interest has been device physics and satellite communication. He has published six research papers in intenational journals and conferences.

**Professor Harish Parthasarathy** has been working as full professor in ECE deperment of NSIT Dwarka, Delhi-110075. His area of research is DSP. He has published several books in this domain and has guided five PhD students and guiding our more PhD students under Delhi University.

# A Study on Associative Neural Memories

B.D.C.N.Prasad[1]

Dept. of Computer Applications,
P V P Siddhartha Institute of Technology
Vijayawada, India
bdcnprasad@gmail.com

P E S N Krishna Prasad[2]

Dept. of Computer Science & Engineering
Aditya Engineering College
Kakinada, India
surya125@gmail.com

Sagar Yeruva[3]

Department of Computer Applications,
St. Peter's Engineering College, Hyderabad.
sagaryeruva@yahoo.com

P Sita Rama Murty[4]

Dept. of Computer Science & Engineering
Sri Sai Aditya Institute of Science & Technology
Kakinada, India
psramam.1@gmail.com

*Abstract*— **Memory plays a major role in Artificial Neural Networks. Without memory, Neural Network can not be learned itself. One of the primary concepts of memory in neural networks is Associative neural memories. A survey has been made on associative neural memories such as Simple associative memories (SAM), Dynamic associative memories (DAM), Bidirectional Associative memories (BAM), Hopfield memories, Context Sensitive Auto-associative memories (CSAM) and so on. These memories can be applied in various fields to get the effective outcomes. We present a study on these associative memories in artificial neural networks.**

*Keywords-Associative memories; SAM; DAM; Hopfield model; BAM; Holographic Associative Memory (HAM); Context-sensitive Auto-associative Memory (CSAM); Context-sensitive Asynchronous Memory (CSYM)*

## I. INTRODUCTION

Learning is the way we acquire knowledge about the world around us, and it is through this process of knowledge acquisition, that the environment alerts our behavioral responses. Learning allows us to store and retain knowledge; it builds our memories.

Aristotle stated about memory: first, the elementary unit of memory is a sense image and second, association and links between elementary memories serve as the basis for higher level cognition. One author stated, memory stands for the elementary unit and association for recollection between elementary units.

In a neurobiological context, memory refers to the relatively enduring neural alterations induced by the interaction of an organism with its environment. Without such a change, there is no memory. The memory must be useful and accessible to the nerves system that influences the future behavior.

Memory and Learning are intricately connected. When a particular activity pattern is learned, it is stored in the brain where it can be recalled later when required. Learning encodes information. A system learns a pattern if the system encodes the pattern in its structure. The system structure changes as the system learns the information. So, learning involves change. That change can be represented in memory for future behavior.

Over the past century the psychologists have studied learning based on fundamental paradigms*: non-associative* and *associative*. In *non-associative learning* an organism acquires the properties of a single repetitive stimulus. In *associative learning* [Edward Thorndike, B.F. Skinner], an organism acquires knowledge about the relationship of either one stimulus to another, or one stimulus to the organisms own behavioral response to that stimulus.

On the neuronal basis of formation of memories into two distinct categories: STM (short term memory) and LTM (long term memory). Inputs to the brain are processed into STM's which last at the most for a few minutes. Information is downloaded into LTM's for more permanent storage. One of the most important functions of our brain is the laying down and recall of memories. It is difficult to imagine how we could function without both short and long term memory. The absence of short term memory would render most tasks extremely difficult if not impossible - life would be punctuated by a series of one time images with no logical connection between them. Equally, the absence of any means of long term memory would ensure that we could not learn by past experience.

The acquisition of knowledge is an active, on going cognitive process based on our perceptions. An important point about the learning mechanism is that it distributes the memory over different areas, making them robust to damage. Distributed storage permits the brain to work easily from partially corrupted information.

## II. ASSOCIATIVE MEMORIES

The associative memory models[4], an early class of neural models that fit perfectly well with the vision of cognition emergent from today brain neuro-imaging techniques, are inspired on the capacity of human cognition to build calculus makes them a possible link between connectionist models and classical artificial intelligence developments.

Our memories function as an **associative** or **content-addressable**. That is, a memory does not exist in some isolated fashion, located in a particular set of neurons. Thus memories are stored in *association* with one another. These different sensory units lie in completely separate parts of the brain, so it is clear that the memory of the person must be distributed throughout the brain in some fashion.

We access the memory by its contents not by where it is stored in the neural pathways of the brain. This is very powerful; given even a poor photograph of that person we are quite good at reconstructing the persons face quite accurately. This is very different from a traditional computer where specific facts are located in specific places in computer memory. If only partial information is available about this location, the fact or memory cannot be recalled at all.

Traditional measures of associative memory performance are its *memory capacity* and *content-addressability*. Memory capacity refers to the maximum number of associated pattern pairs that can be stored and correctly retrieved while content-addressability is the ability of the network to retrieve the correct stored pattern. Obviously, the two performance measures are related to each other.

It is known that using Hebb's learning rule in building the connection weight matrix of an associative memory yields a significantly low memory capacity. Due to the limitation brought about by using Hebb's learning rule, several modifications and variations are proposed to maximize the memory capacity.

### A. Model

Associative memory maps[4,6] data from an input space to data in an output space. In general, this mapping is from unknown domain points to known range points, where the memory learns an underlying association from a training data set.

For non-learning memory models, which have their origin in additive neuronal dynamics, connection strength's are "programmed" a priori depending upon the association that are to be encoded in the system. Sometimes these memories are referred to as matrix associative memories, because a connection matrix W, encodes associations $(A_i, B_i)_{i-1}^{Q}$, where $A_i \in B^n$ and $B_i \in B^n$. If $(A_i, B_i)$ is one of the programmed memories then $B_i$ is called the association of $A_i$. When $A_i$ and $B_i$ are in different spaces then the model is ***hetero-associative memory.*** *i.*e. it associates two different vectors with one another. if $A_i = B_i$, then the model is ***Auto-associative memory.*** i.e. it associates a vector with itself. Associative memory models enjoy properties such as fault tolerance.

### Associative Neural Memories

Associative neural memories are concerned with associative learning and retrieval of information (vector patterns) in neural networks. These networks represent one of the most extensively analyzed classes of artificial neural networks.

Several associative neural memory models have been proposed over the last two decades. These memory models can be classified into various ways depending on

- Architecture (Static *versus* Dynamic)

- Retrieval Mode (Synchronous v*ersus* Asynchronous)

- Nature of stored association (Auto-associative *versus* Hetero-associative)

- Complexity and capability of memory storage

### 1) Simple Associative Memories

One of the earliest associative memory models is the correlation memory [Anderson, 1972; Kohonen, 1972; Nakano, 1972]. This correlation memory consists of a single layer of **L** non interacting linear units, with the $l^{th}$ unit having a weight vector $w_l \in R^n$. It associates real values input column vectors $x^k \in R^n$ which corresponding real valued output column vectors $y^k \in R^n$ according to the transfer eq.:

$$y^k = Wx^k \tag{1}$$

Where $\{x^k, y^k\}, k = 1,2,....n$ a collection of desired associations and W is an L×N interconnection matrix whose $l^{th}$ row is given by $w_l^T$. This associative memory is characterized by linear matrix vector multiplication retrievals. Hence it is referred to as a ***linear associative memory*** [1](LAM). This LAM is said to be *hetero-associative* because $y^k$ is different (in encoding and/ dimensionality) from $x^k$. If $y^k = x^k$ for all k, then this memory is called *auto-associative*.

The correlation memory is a LAM that employs a simple recording or storage recipe for loading *m* associations $\{x^k, y^k\}, k = 1,2,....n$ into memory. The recording recipe is responsible for synthesizing *W* and is given by

$$W = \sum_{k=1}^{m} y^k (x^k)^T = YX^T \tag{2}$$

Where *W* is the Correlation Matrix of *m* associations.

### 2) Simple Nonlinear Associative Memory Model

The binary-valued associations[1] $\mathbf{x}^k$, {-1, +1} N and $\mathbf{y}^k$ {-1, +1} L and the presence of a clipping nonlinearity *F* operating component wise on the vector Wx, according to

$$y = F[Wx] \tag{3}$$

relaxes some of the constraints imposed by correlation recording of a LAM. Here, **W** needs to be synthesized with the requirement that only the sign of the corresponding components of $\mathbf{y}^k$ and $\mathbf{Wx}^k$ agree. Next, consider the normalized correlation recording recipe given by:

$$W = \frac{1}{n} \sum_{k=1}^{m} y^k \left(x^k\right)^T \qquad (4)$$

This automatically normalizes the $x^k$ vectors. Now, if one of the recorded key patterns $\mathbf{x}^h$ is presented as input, then the following expression for the retrieved memory pattern can be written:

$$\tilde{y}^h = F\left[ y^h + \frac{1}{n} \sum_{k \neq h}^{m} y^k \left(x^k\right)^T x^h \right] = F\left[y^h + \Delta^h\right] \qquad (5)$$

*3) Optimal Linear Associative Memory (OLAM)*

The correlation recording recipe does not make optimal use of the LAM interconnection weights. A more optimal recording technique can be derived which guarantees perfect retrieval of stored memories $\mathbf{y}^k$ from inputs $\mathbf{x}^k$ as long as the set $\{\mathbf{x}^k; k = 1, 2, ..., m\}$ is linearly independent (as opposed to the more restrictive requirement of orthogonality required by the correlation-recorded LAM). This recording technique leads to the optimal linear associative memory[1,7] (OLAM). For perfect storage of $m$ associations $\{\mathbf{x}^k, \mathbf{y}^k\}$, a LAM's interconnection matrix $\mathbf{W}$ must satisfy the matrix equation given by:

$$Y = WX \qquad (6)$$

This equation can always be solved exactly if all $m$ vectors $\mathbf{x}^k$ (columns of $\mathbf{X}$) are linearly independent, which implies that $m$ must be smaller or equal to $n$. For the case $m = n$, the matrix $\mathbf{X}$ is square and a unique solution for $\mathbf{W}$ in Equation (6) exists giving:

$$W^* = YX^{-1} \qquad (7)$$

Which requires that the matrix inverse $\mathbf{X}^{-1}$ exists; i.e., the set $\{\mathbf{x}^k\}$ is linearly independent. Thus, this solution guarantees the perfect recall of any $\mathbf{y}^k$ upon the presentation of its associated key $\mathbf{x}^k$.

*B. Dynamic Associative Memories (DAMs)*

Associative memory performance can be improved by utilizing more powerful architectures than the simple ones considered above. As an example, consider the auto associative version of the single layer associative memory employing units with the sign activation function and whose transfer characteristics are given by Equation (3). Now assume that this memory is capable of associative retrieval of a set of $m$ bipolar binary memories $\{\mathbf{x}^k\}$. Upon the presentation of a key $\hat{\mathbf{x}}^k$ which is a noisy version of one of the stored memory vectors $\mathbf{x}^k$, the associative memory retrieves (in a single pass) an output $\mathbf{y}$ which is closer to stored memory $\mathbf{x}^k$ than $\hat{\mathbf{x}}^k$. In general, only a fraction of the noise (error) in the input vector is corrected in the first pass. Intuitively, we may proceed by taking the output $\mathbf{y}$ and feed it back as an input to the associative memory hoping that a second pass would eliminate more of the input noise. This process could continue with more passes until we eliminate all errors and arrive at a final output $\mathbf{y}$ equal to $\mathbf{x}^k$. The retrieval procedure just described amounts to constructing a recurrent associative memory with the synchronous (parallel) dynamics given by

$$x(t + 1) = F[W\,x(t)] \qquad (8)$$

Where $t = 0, 1, 2, 3, ...$ and $\mathbf{x}(0)$ is the initial state of the dynamical system which is set equal to the noisy key $\hat{\mathbf{x}}^k$. For proper associative retrieval, the set of memories $\{\mathbf{x}^k\}$ must correspond to stable states (attractors) of the dynamical system. In this case, we should synthesize $\mathbf{W}$ (which is the set of all free parameters $w_{ij}$ of the dynamical system in this simple case) so that starting from any initial state $\mathbf{x}(0)$, the dynamical associative memory converges to the "closest" memory state $\mathbf{x}^k$.

DAM has several variations, which are presented below:

*1) Hopfield model*

The Hopfield model[6,8] is a distributed model of an associative memory. The *Hopfield Model* was proposed by John Hopfield of the California Institute of Technology during the early 1980s.

The dynamics of the Hopfield model is different from that of the *Linear Associator Model* in that it computes its output recursively in time until the system becomes stable. We presented a Hopfield model with six units, where each node is connected to every other node in the network is given below.



Figure 1: Hopfield model

Unlike the linear associator model which consists of two layers of processing units, one serving as the input layer while the other as the output layer, the Hopfield model consists of a single layer of processing elements where each unit is connected to every other unit in the network other than itself. The connection weight matrix $\mathbf{W}$ of this type of network is square and symmetric, i.e., $w_{ij} = w_{ji}$ for $i, j = 1, 2, ..., m$. Each unit has an extra external input $I_i$. This extra input leads to a modification in the computation of the net input to the units.

Unlike the linear associator, the units in the Hopfield model act as both input and output units. But just like the linear associator, a single associated pattern pair is stored by computing the weight matrix as follows:

$$W_k = X_k^T Y_k \text{, where } Y_k = X_k$$

$$W = \alpha \sum_{k=1}^{p} W_k \qquad (9)$$

to store *p* different associated pattern pairs. Since the Hopfield model is an auto-associative memory model, patterns, rather than associated pattern pairs, are stored in memory.

After encoding, the network can be used for decoding. Decoding in the Hopfield model is achieved by a collective and recursive relaxation search for a stored pattern given an initial stimulus pattern. Given an input pattern *X*, decoding is accomplished by computing the net input to the units and determining the output of those units using the output function to produce the pattern *X'*. The pattern *X'* is then fed back to the units as an input pattern to produce the pattern *X''*. The pattern *X''* is again fed back to the units to produce the pattern *X'''*. The process is repeated until the network stabilizes on a stored pattern where further computations do not change the output of the units.

If the input pattern *X* is an incomplete pattern or if it contains some distortions, the stored pattern to which the network stabilizes is typically one that is most similar to *X* without the distortions. This feature is called ***Pattern Completion*** and is very useful in many image processing applications.

During decoding, there are several schemes that can be used to update the output of the units. The updating schemes are *Synchronous* (parallel), *Asynchronous (*sequential), or a combination of the two (*hybrid*).

Using the synchronous updating scheme, the output of the units are updated as a group prior to feeding the output back to the network. On the other hand, using the asynchronous updating scheme, the output of the units are updated in some order (e.g. random or sequential) and the output are then fed back to the network after each unit update. Using the hybrid synchronous-asynchronous updating scheme, subgroups of units are updated synchronously while units in each subgroup updated asynchronously. The choice of the updating scheme has an effect on the convergence of the network.

Hopfield (1982) demonstrated that the maximum number of patterns that can be stored in the Hopfield model of *m* nodes before the error in the retrieved pattern becomes severe is around 0.15*m*. The memory capacity of the Hopfield model can be increased as shown by André cut (1972).

Hopfield model is broadly classified into two categories:

- Discrete Hopfield Model

- Continuous Hopfield Model

*2) Brain-state-in-a-Box model*

The "brain-state-in-a-box"[4,6] (BSB) model is one of the earliest DAM models. It is a discrete-time continuous-state parallel updated DAM. The BSB model extends the Linear Associator model and is similar to the Hopfield Model in that it is an Auto-associative model with its connection matrix computed using outer products in the usual way. The operation of both models is also very similar, with differences arising primarily in the way activations are computed in each iteration,

and in the signal function used. The BSB model stands apart from other models in its use of the linear threshold signal function.

*3) Bi-directional Associative Memory (BAM)*

Kosko (1988) extended the Hopfield model by incorporating an additional layer to perform recurrent auto-associations as well as hetero-associations on the stored memories.

The network structure of the Bi-directional Associative Memory model[4,7] is similar to that of the linear associator, but the connections are bidirectional, i.e., wij = wji, for i = 1, 2, ..., m and j = 1, 2, ..., n. Also, the units in both layers serve as both input and output units depending on the direction of propagation. Propagating signals from the X layer to the Y layer makes the units in the X layer act as input units while the units in the Y layer act as output units. The same is true for the other direction, i.e., propagating from the Y layer to the X layer makes the units in the Y layer act as input units while the units in the X layer act as output units. Below is an illustration of the BAM architecture.



Figure 2: BAM model

Just like the linear associator and Hopfield model, encoding in BAM can be carried out by using:

$$W_k = X_k^T Y_k \qquad (10)$$

to store a single associated pattern pair and

$$W = \alpha \sum_{k=1}^{p} W_k \qquad (11)$$

to simultaneously store several associated pattern pairs. After encoding, the network can be used for decoding.

In BAM, decoding involves reverberating distributed information between the two layers until the network becomes stable. In decoding, an input pattern can be applied either on the *X* layer or on the *Y* layer. When given an input pattern, the network will propagate the input pattern to the other layer allowing the units in the other layer to compute their output

values. The pattern that was produced by the other layer is then propagated back to the original layer and let the units in the original layer compute their output values. The new pattern that was produced by the original layer is again propagated to the other layer. This process is repeated until further propagations and computations do not result in a change in the states of the units in both layers where the final pattern pair is one of the stored associated pattern pairs. The final pattern pair that will be produced by the network depends on the initial pattern pair and the connection weight matrix.

Several modes can also be used to update the states of the units in both layers namely *synchronous, asynchronous, and a combination of the two*. In *synchronous* updating scheme, the states of the units in a layer are updated as a group prior to propagating the output to the other layer. In *asynchronous* updating, units in both layers are updated in some order and output is propagated to the other layer after each unit update. Lastly, in synchronous-asynchronous updating, there can be subgroups of units in each layer that are updated synchronously while units in each subgroup are updated asynchronously.

Since the BAM also uses the traditional Hebb's learning rule to build the connection weight matrix to store the associated pattern pairs, it too has a severely low memory capacity. The BAM storage capacity for reliable recall was given by Kosko (1988) to be less than *minimum* (*m*, *n*), i.e., the minimum of the dimensions of the pattern spaces. A more recent study by Tanaka et al (2000) on the relative capacity of the BAM using statistical physics reveals that for a system having *n* units in each of the two layers, the capacity is around 0.1998n.

Mostly, BAM can be classified into two categories:

- Discrete BAM

In a discrete BAM, the network propagates an input pattern $X$ to the $Y$ layer where the units in the $Y$ layer will compute their net input.

- Continuous BAM

In the continuous BAM, the units use the sigmoid or hyperbolic **tangent** output function. The units in the $X$ layer have an extra external input $I_i$ , while the units in the $Y$ layer have an extra external input $J_j$ for $i = 1, 2, ..., m$ and $j = 1, 2, ..., n$. These extra external inputs lead to a modification in the computation of the net input to the units.

*4) Context Sensitive Auto Associative Memory model (CSAM)*

The matrix correlation memories can be very efficiently modulated by contexts in the case in which the key vector and a vectorial context are combined confirming a Kronecker product. The existence of multiplicative contexts enlarges in many directions the cognitive abilities of the correlation memories. One of the abilities of the context-sensitive associative memories is the possibility to implement all the basic logical operations of the propositional calculus. Moreover, these memories are capable of computing some fundamental operations of modal logic. The theory of logic and

the theory of context-dependent associative memories converge to operator formalism.

This model is referred to as Context dependent auto-associative memory neural network [2,3], which is more powerful algorithm that suits to compute the clinical and laboratory factors effectively. Here, we could use the Kronecker product matrix as memory representation in the network structure. The model of this algorithm is presented below.
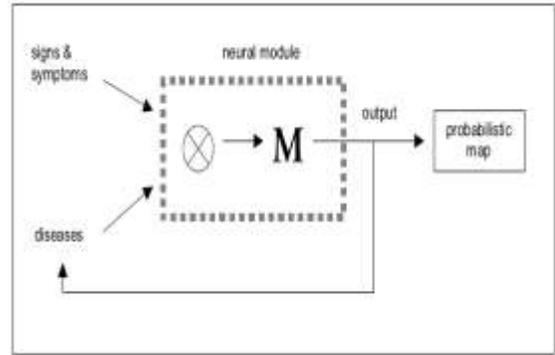


Figure 3: Context-sensitive Auto Associative model

The neural module receives the input of two vectors: one representing the set of possible diseases up to the moment and the other vector corresponding to a new sign, symptom or laboratory result. The action of the neurons that constitute the neural module can be divided into two sequential steps: the Kronecker product of these two entries and the association of this stimulus with an output activity pattern. This output vector is a linear combination of a narrower set of disease vectors that can be reinjected if a new clinical data arrives or can be processed to obtain the probability attribute to each diagnostic decision.

A context-dependent associative memory M acting as a basic expert system is a matrix

$$M = \sum_{i=1}^{k} d_i \left( d_i \otimes \sum_{j(i)} s_j \right)^T \tag{12}$$

Where di are column vectors mapping k different diseases (the set {d} is chosen orthonormal), and sj(i) are column vectors mapping signs or symptoms accompanying the i disease (also an orthonormal set). The sets of symptoms corresponding to each disease can overlap. The Kronecker product between two matrices A and B is another matrix defined by

$$A \otimes B = a(i, j).B \tag{13}$$

Denoting that each scalar coefficient of matrix A, a(i, j), is multiplied by the entire matrix B. Hence, if A is nxm dimensional and B is k x l dimensional, the resultant matrix will have the dimension nk x ml.

### 5) Context Sensitive Asynchronous Memory Model (CSYM)

Context-sensitive asynchronous memory[10,11] is a priming-based approach to memory retrieval. It exploits feedback from the task and environment to guide and constrain memory search by interleaving memory retrieval and problem solving. Solutions based on context-sensitive asynchronous memory provide useful answers to vague questions efficiently, based on information naturally available during the performance of a task.
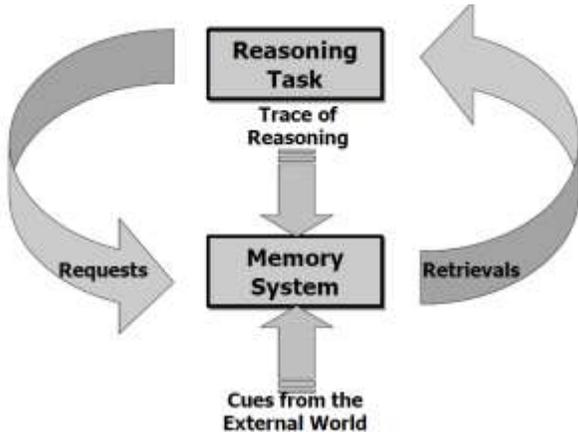


Figure 4: *A Context Sensitive Associative memory*

Reasoning is an important biological activity. Emerged at some point of biological evolution, reasoning is a neural function with a decisive role for the kind of life style that the human beings have developed. The knowledge of the external world, in fact one of the highlights of human culture, is a direct consequence of the capacity of reasoning.

A context-sensitive asynchronous memory differs from these "traditional" approaches by simultaneously providing a means to allow reasoning to proceed in parallel with memory search as well as a means to guide ongoing memory search.

- **Asynchronous retrieval:** Asynchronous retrieval is the autonomous processing of memory retrieval requests. A memory system can perform asynchronous retrieval by using reified retrieval requests and a retrieval monitor which operates in conjunction with the agent's task controller. Asynchrony logically includes spontaneous retrieval (retrieval at the discretion of memory) and anytime retrieval (retrieval on demand of the reasoner).

- **Context-sensitivity:** Context sensitivity is using feedback to guide memory search. Context sensitivity can be achieved through a process, called context-directed spreading activation, which operates hand in hand with the agent's working memory.

### 6) Holographic Associative Memory (HAM)

In 1990 Sutherland in his pioneering work, presented the first truly holographic associative memory[20] with holographic representation and learning algorithm analogous to correlation learning. It is a two-dimensional (2-D) generalized multidimensional phased representation.

Here information is mapped onto the phase orientation of complex numbers operating. It can be considered as a complex valued artificial neural network. Generally speaking, holographic networks are very suitable for those problems where stimuli are long vectors with symmetrically (uniformly) distributed arguments. A longer stimulus vector assures a greater learning capacity, i.e. a greater number of stimulus-response associations that can be learned. Symmetry in arguments assures accuracy in reproducing the learned stimulus-response associations.

Holographic memory is a storage device that is being researched and slated as the storage device that will replace hard drives and DVDs in the future. It has the potential of storing up to 1 terabyte or one thousand gigabytes of data in a crystal the size of a sugar cube.

#### Advantages of Holographic Memory Systems

Aside from having a tremendous amount of storage space for data, holographic memory systems also have the ability to retrieve data very quickly, up to a 1 gigabyte per second transfer rate

The main difference between holographic and conventional neural networks is that a holographic neuron is more powerful than a conventional one, so that it is functionally equivalent to a whole conventional network. Consequently, a holographic network usually requires a very simple topology consisting of only few neurons. Another characteristic of the holographic technology is that it represents information by complex numbers operating within two degrees of freedom (value and confidence). Also an important property is that holographic training is accomplished by direct (almost non-iterative) algorithms, while conventional training is based on relatively slow "back-propagation" (gradient) algorithms. A *holographic neuron* is sketched in *Figure 5* . As we can see, it is equipped with only one input channel and one output channel. However, both channels carry whole vectors of complex numbers. An input vector S is called a stimulus and it has the form

$$S = \left[ \lambda_1 e^{i\theta_1}, \lambda_2 e^{i\theta_2}, \lambda_3 e^{i\theta_3}, \ldots, \lambda_n e^{i\theta_n} \right] \quad (14)$$

An output vector R is called a response and its form is

$$R = \left[ \gamma_1 e^{i\phi_1}, \gamma_2 e^{i\phi_2}, \ldots, \gamma_m e^{i\phi_m} \right] \quad (15)$$
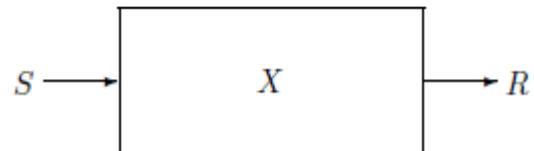


Figure 5: *A Holographic neuron*

All complex numbers above are written in polar notation, so that moduli (magnitudes) are interpreted as confidence levels of data, and arguments (phase angles) serve as actual values of data. The neuron internally holds a complex $n \times m$

matrix $X = \lfloor x_{jk} \rfloor$ which serves as a memory for recording associations.

We present the *basic learning process*. Learning one association between a stimulus *S* and a desired response *R* requires that the correlation between the *j*-th stimulus element and the *k*-th response element is accumulated in the (*j; k*)-th entry of the memory matrix. More precisely:

$$x_{jk} + = \lambda_j \gamma_k e^{i(\phi_k - \theta_j)} \qquad (16)$$

The same formula can be written in the matrix-vector form:

$$X + = \bar{S}^T R \qquad (17)$$

Here $\bar{S}$ denotes the conjugated transpose of the vector *S*.

Mostly Holographic memory is of Different models, but we are considering two of them;

- Dynamically Structured Holographic Associative Memory (DSHAM).
- Composite Holographic Associative Memory (CHAM).

### III. APPLICATIONS

*A. Human Resources Management*

*1) Employee Selection and Hiring - predict on which job an applicant will achieve the best job performance.*
    *Input data*: information about an applicant: personal information, previous jobs, educational levels, previous performance, etc.

*2) Employee Retention - identify potential employees who are likely to stay with the organization for a significant amount of time based on data about an applicant.*
    *Input data*: applicant's hours of availability, previous jobs, educational levels and other routine information.

*3) Staff Scheduling - predict staff requirements for restaurants, retail stores, police stations, banks, etc.*
    *Input data*: time of year, day of week, pay-days, holidays, weather, etc.

*4) Personnel Profiling - forecast successful completion of training program; identify employees most suitable for a certain task.*
    *Input data*: background characteristics of individuals.

*B. Medical*

*1) Medical Diagnosis - Assisting doctors with their diagnosis by analyzing the reported symptoms.*
    *Input data*: patient's personal information, patterns of symptoms, heart rate, blood pressure, temperature, laboratory results, etc.

*2) Detection and Evaluation of Medical Phenomena - detect epileptic attacks, estimate prostate tumor size, detects patient breathing abnormalities when a patient is under anesthesia, etc.*
    *Input data*: patient's personal information, breathing rate, heart rate, patterns of symptoms, blood pressure, temperature, etc.

*3) Patient's Length of Stay Forecasts - forecast which patients remain for a specified number of days.*
    *Input data*: personal information such as age and sex, level of physical activity, heart rate, blood pressure, temperature and laboratory results, treatment procedures, etc.

*4) Treatment Cost Estimation*
    *Input data*: personal information such as age and sex, physiological data, the use of drug or other therapies, treatment procedures, number of recurrences after first treatment, etc.

*C. Financial*

*1) Stock Market Prediction - predict the future movement of the security using the historical data of that security.*
    *Input data*: Open, High, Low, Close, Volume, technical indicators, market indexes and prices of other securities.

*2) Credit Worthiness - decide whether an applicant for a loan is a good or bad credit risk.*
    *Input data*: applicant's personal data, income, expenses, previous credit history, etc.

*3) Credit Rating - assign credit ratings to companies or individuals based on their financial state.*
    *Input data*: current financial state indicators and past financial performance of a company or individual.

*4) Bankruptcy prediction - classify a company as potential bankruptcy.*
    *Input data*: company characteristics and business ratios, such as working capital/total assets, retained earnings/total assets, earnings before interest and taxes/total assets, market value of equity/total debt, and sales/total assets.

*5) Property Appraisal - evaluate real estate, automobiles, machinery and other property.*
    *Input data*: property parameters, environment conditions as well as appropriate demographic, ecological, industrial and other factors.

*6) Fraud Detection - detect and automatically decline fraudulent insurance claims, client transactions, and taxes.*
    *Input data*: transaction parameters, applicant's information and other data of past incidents.

*7) Price Forecasts - forecast prices of raw materials, commodities, and products.*
    *Input data*: previous price movements, economical indicators, market indexes.

*8) Economic Indicator Forecasts - forecast economic indicators for the next week, month, and quarter.*
    *Input data*: social and economical indicators, time-series data of an indicator.

### D. Sales and Marketing

*1) Sales Forecasting - predict future sales based on historical information about previous marketing and sales activities.*
   *Input data*: historical data about marketing budget, number of ads, special offers and other factors affecting sales.

*2) Targeted Marketing - reduce costs by targeting a particular marketing campaign to the group of people which have the highest response rate. Avoid wasting money on unlikely targets.*
   *Input data*: information about customers and their response rate.

*3) Service Usage Forecasting - forecast the number of service calls, customer transactions, customer arrivals, reservations or restaurant covers (patrons) in order to effectively schedule enough staff to handle the workload.*
   *Input data*: season, day-of-week, hour of the day, special events in the city/area, marketing budget, promotional events, weather, etc.

*4) Retail Margins Forecasting - forecast the behavior of margins in the future to determine the effects of price changes at one level on returns at the other.*
   *Input data*: retail prices, expenditures at the retail level, marketing costs, past margin values, price variability and other market characteristics.

### E. Industrial

*1) Process Control - determine the best control settings for a plant. Complex physical and chemical processes that may involve interaction of numerous (possibly unknown) mathematical formulas can be modeled heuristically using a neural network.*

*2) Quality Control - predict the quality of plastics, paper, and other raw materials; machinery defect diagnosis; diesel knock testing, tire testing, beer testing.*
   *Input data*: product/part/machinery characteristics, quality factor.

*3) Temperature and force prediction in mills and factories*
   *Input data*: previous values of temperature, force and other characteristics of mills and factories.

### F. Operational Analysis

*1) Retail Inventories Optimization - forecast optimal stock level that can meet customer needs, reduce waste and lessen storage; predict the demand based on previous buyers' activity.*
   *Input data*: characteristics of previous buyers' activity, operating parameters, season, stock, and budgets.

*2) Scheduling Optimization - predict demand to schedule buses, airplanes, and elevators.*
   Input data: season, day-of-week, hour of the day, special events in the city/area, Weather, etc.

*3) Managerial decision making - select the best decision option using the classification capabilities of neural network.*
   *Input data*: initial problem parameters and final outcome.

*4) Cash flow forecasting - maximize the use of resources with more accurate cash flow forecasts.*
   *Input data*: accounts payable, accounts receivable, sales forecasts, budgets, capital expenditures, stock, season, operating data, etc.

### G. Data Mining

*1) Prediction - use some variables or fields in the database to predict unknown or future values of other variables of interest.*
*2) Classification - map (classify) a data item into one of several predefined classes.*
*3) Change and Deviation Detection - uncover certain data records that are in some way out of the ordinary records; determine which cases/records suspiciously diverge from the pattern of their peers.*
*4) Knowledge Discovery - find new relationships and non-obvious trends in the data.*
*5) Response Modeling - build a neural network based response model.*
*6) Time Series Analysis - forecast future values of a time series.*

### IV. CONCLUSION

We presented a survey on various associative memories which are being used in various applications in recent trends. There are few points that we would like to mention through this article:

Memory is not just a passive store for holding ideas without changing them; it may transform those ideas when they are being retrieved. There are many examples showing that what is retrieved is different from what was initially stored.

Simple Associative memories are static and very low memory so that they cannot be applied in the applications where high memory is required.

Dynamic Associative memories such as Hopfield, BSB, and BAM are Dynamical memories but they are also capable of supporting very low memory, so they cannot be applied in the applications where high memory requirements are there.

A simple model describing context-dependent associative memories generates a good vectorial representation of basic logical calculus. One of the powers of this vectorial representation is the very natural way in which binary matrix operators are capable to compute ambiguous situations. This fact presents a biological interest because of the very natural way in which the human mind is able to take decisions in the presence of uncertainties. Also these memories could be used to develop expert agents to the recent problem domain.

Holographic memories are being used to build the many advanced memory based agents like memory cards, USB Drives, etc.,

Context asynchronous memories are being used to develop experience based agents. Context-sensitive asynchronous memory enables agents to efficiently balance resources between task processing and memory retrieval. By incrementally searching the knowledge base, a context-sensitive asynchronous memory supports anytime retrieval of the best answer found so far and thus enables agents to satisfy.

Finally, we conclude that Context Sensitive Auto Associative memory and Asynchronous Memory and Holographic Associative Memory can be used to solve the real applications which we mentioned.

## V. FUTURE SCOPE

Among these memories we have used Context sensitive auto-associative memory model to implement the expert system for medical diagnosis of Asthma. We are planning to implement the medical based expert systems with the use of Context Sensitive Asynchronous Memory Models (CSYM) and also to implement mining based applications including time series analysis.

## REFERENCES

[1] "Fundamentals of Artificial Neural Networks" – Mohadmad H. Hassoun.

[2] "A Comparative Study of Machine Learning Algorithms as Expert Systems in Medical Diagnosis (Asthma)" – Dr. B D C N Prasad, P E S N Krishna Prasad and Y Sagar, CCSIT 2011, Part I, CCIS 131, pp. 570–576, 2010.

[3] "Context-Sensitive Auto-Associative Memories as Expert Systems in Medical Diagnosis" – Andres Pomi and Fernando Olivera, BioMed Central.

[4] "Neural Networks, A Classroom Approach" – Satish Kumar.

[5] "The Moral Demands of Memory" - Jeffrey Blustein.

[6] "Neural Network Design" – Martin T. Hagan, Howard B. Demuth, Mark Beale.

[7] "Neural Networks and Fuzzy Systems - A dynamical Systems Approach to Machine Intelligence" – Bart Kosko.

[8] Anderson JA, Cooper L, Nass MM, Freiberger W, Grenander U:Some properties of a neural model for memory. AAAS Symposium on Theoretical Biology and Biomathematics 1972 [http://www.physics.brown.edu/physics/researchpages/Ibns/Cooper%20Pub040_SomePropertiesNeural_72.pdf]. Milton, WA. Leon N Cooper Publications.

[9] "An Introduction to Neural Networks" – James A. Anderson.

[10] "Neural Networks Ensembles, Cross Validation and Active Learning" – Anders Krogh and Jesper Vedelsby, MIT Press[1995].

[11] Mizraji E: Context-dependent associations in linear distributed memories. Bulletin Math Biol 1989, 51:195-205.

[12] Valle-Lisboa JC, Reali F, Anastasía H, Mizraji E: Elman topology with sigma-pi units: An application to the modelling of verbal hallucinations in schozophrenia. Neural Networks 2005,18:863-877.

[13] Mizraji E, Pomi A, Alvarez F: Multiplicative contexts in associative memories. BioSystems 1994, 32:145-161.

[14] Pomi-Brea A, Mizraji E: Memories in context. BioSystems 1999,50:173-188.

[15] Mizraji E, Lin J: A dynamical approach to logical decisions. Complexity 1997, 2:56-63.

[16] Mizraji E, Lin J: Fuzzy decisions in modular neural networks. Int J Bifurcation and Chaos 2001, 11:155-167.

[17] Pomi A, Mizraji E: A cognitive architecture that solves a problem stated by Minsky. IEEE on Systems, Man and Cybernetics B (Cybernetics) 2001, 31:729-734.

[18] "Context-Sensitive Asynchronous Memory"- Anthony G. Francis, Jr.

[19] "Associative Memories in Medical Diagnostic"- Jorge M. Barreto, Fernando de Azevedo, Carlos I. Zanchin, Lycia R. Epprecht.

[20] "Context-Sensitive Associative Memory: Residual Excitation" in Neural Networks as the Mechanism of STM and Mental Set - Victor Eliashberg

[21] "Retrieving Semantically Distinct Analogies" - Michael Wolverton.

[22] "Improving Array Search Algorithm Using Associative Memory Neural Network" - Emad Issa Abdul Kareem and Aman Jantan.

[23] "Hippocampal Auto-Associative Memory" - Nicolas P. Rougier.

[24] "Phrase Detection and the Associative Memory Neural Network" - Richard C. Murphy.

[25] "Human Recognition in Passive Environment using Bidirection Associative Memory" - Vivek Srivastava and Vinay kumar Pathak.

[26] "Generalized Asymmetrical Bidirectional Associative Memory" - Tae-Dok Eom, Changkyu Choi, and Ju-Jang Lee.

[27] UCI Machine Learning Repository:

[28] "Characteristics of MultiDimensional Holographic Associative Memory in Retrieval with Dynamically Localizable Attention" – Javed I . Khan.

[29] "Bidirectional Associative Memory Neural Network method in the Character Recognition" - Yash Pal Singh, V.S.Yadav, Amit Gupta, Abhilash Khare.

[30] Cross SS, Harrison RF, Lee Kennedy R: Introduction to neural networks. The Lancet 1995, 346:1075-1079.

[31] E. Oja, "A simplified neuron model as a principal component analyzer," J. Math. Biol., vol. 15, pp. 267–273, 1982.

[32] D. Psaltis and F. Mok, "Holographic Memories," Sci. Amer., Nov. 1995, pp. 70–76.

[33] H.-M. Tai and T. L. Jong, "Information storage in high-order neural networks with unequal neural activity," J. Franklin Inst., vol. 327, no.1, pp. 16–32, 1990.

[34] M. Wenyon, Understanding Holography. New York: Arco, 1978.

[35] B.Widrow and M. E. Hoff, "Adaptive switching circuits," IRE WESCON Convention Rec., pt. 4, pp. 96–104, 1960.

[36] D. Willshaw, "Holography, associative memory and inductive generalization," in Parallel Models of Associative Memory, G. E. Hinton and J. E. Anderson, Eds. Hillsdale, NJ: Lawrence Erlbaum, 1985.

[37] McKerrow, P.J. (1991). Introduction to Robotics. Sydney, Australia: Addison-Wesley.

[38] McNamara, T.P. (1992). Priming and constraints it places on theories of memory and retrieval. Psychological Review, 99, 650-662.

[39] McNamara, T. P., & Diwadkar, V. A. (1996). The context of memory retrieval. Journal of Memory and Language, 35, 877-892

[40] Moore, A.W., & Atkeson, C.G. (1995) Memory based neural networks for robot learning. Neurocomputing, 9(3) pp. 243-269.

[41] http://archive.ics.uci.edu/ml/machine-learning-databases/

## AUTHORS PROFILE

Dr. B D C N Prasad, currently is a Professor & Head of Department of Computer Applications at Prasad V. Potluri Siddardha Institute of Science and Technology, Vijayawada, Andhra Pradesh, India. He received Ph.D. in Applied Mathematics from Andhra University, Visakhapatnam,India in 1984. His research interests includes Machine Intelligence, Data Mining, Rough Sets and Information Security in Computer Science and Boundary value problems and Fluid Dynamics in Mathematics. He has several publications in mathematics and computer science in reputed national and international journals. He is a member of ISTAM , ISTE and also he is a national executive member of Indian Society for Rough Sets.

Mr. P E S N Krishna Prasad, currently is a Research Scholor under the guidance of Dr. BDCN Prasad in the area of Machine Intelligence and Neural Networks. He is working as Associate Professor in the Department of CSE, Aditya Engineering College, Kakinada, Andhra pradesh, India. He is a member of ISTE. He has presented and published papers in several national and International conferences and journals. His areas of interest are Artificial Intelligence, Neural Networks and Machine Intelligence.

Mr. P Sita Rama Murty, currently is a Research Scholor, in the area of ATM networks and Information Secuirty. He is working as Assistant Professor in the department of CSE, Sri Sai Aditya Institute of Science and Technology, Kakinada, Andhra Pradesh, India

# Adaptive Channel Estimation Techniques for MIMO OFDM Systems

Md. Masud Rana

Dept. of Electronics and Radio Engineering
Kyung Hee University
South Korea
mrana928@yahoo.com

Md. Kamal Hosain

School of Engineering
Deakin University
Victoria-3217, Australia

*Abstract*— **In this paper, normalized least mean (NLMS) square and recursive least squares (RLS) adaptive channel estimator are described for multiple input multiple output (MIMO) orthogonal frequency division multiplexing (OFDM) systems. These CE methods uses adaptive estimator which are able to update parameters of the estimator continuously, so that the knowledge of channel and noise statistics are not necessary. This NLMS/RLS CE algorithm requires knowledge of the received signal only. Simulation results demonstrated that the RLS CE method has better performances compared NLMS CE method for MIMO OFDM systems. In addition, the utilizing of more multiple antennas at the transmitter and/or receiver provides a much higher performance compared with fewer antennas. Furthermore, the RLS CE algorithm provides faster convergence rate compared to NLMS CE method. Therefore, in order to combat the more channel dynamics, the RLS CE algorithm is better to use for MIMO OFDM systems.**

*Keywords- MIMO; NLMS; OFDM; RLS*

## I. INTRODUCTION

Recently, multiple input multiple output (MIMO) channels have been introduced to achieve high data speed requisite by the next-generation communication systems [1]. The use of MIMO channels provides higher spectral efficiency versus single input single output (SISO), single input multiple output (SIMO), and multiple input single-output (MISO) channels, when the available bandwidth is inadequate. Furthermore, the diversity gain of the MIMO channels is nearly of second order when channel matrix has full rank. Consequently, by employing MIMO channels, not only the mobility of the wireless communications can be increased but also the algorithm can be more robust against fading, which makes it efficient for the requirements of the next-generation wireless services such as wireless local area networks (WLANs), worldwide interoperability for microwave access (WiMAX), wireless fidelity (WiFi), cognitive radio, and 3rd generation partnership project (3GPP) long term evolution (LTE) [2].

In SISO flat channels, channel estimation (CE) and its precision do not have a drastic impact on the performance of the receiver. Whereas in outdoor MIMO channels, the precision and speed of convergence of the channel estimator can drastically affect the performance of the receiver [3]. In SISO communications, the channel estimators may or may not use the training sequence or not. Although the distribution of the training symbols in a block of data affects the performance of systems [4], but due to simplicity, it is conventional to use the training symbols in the first part of each block. If the training sequence is not used, the estimator is called the blind channel estimator. A blind channel estimator uses information latent in statistical properties of the transmitting data [5]. In full-rank MIMO channels, the use of an initial training data is mandatory, and without it, the channel estimator does not converge [2], [5].

Orthogonal frequency division multiplexing (OFDM) systems have attracted much attention as a promising technology in wireless communication systems. In OFDM systems, the whole spectrum is divided into several sub-carriers, and before each OFDM block the cyclic prefix (CP) is inserted. So, OFDM systems can mitigate the effects of multipath and have high spectrum efficiency. Therefore, OFDM is the important technique for next-generation communication. However, it is strict about CE to exploit the coherent demodulation, detection and decoding [7].

Several CE techniques have been proposed to mitigate interchannel interference (ICI) in OFDM systems. In [8], the least square (LS) CE has been proposed to minimize the squared differences between the received and estimated signal. The LS algorithm, which is independent of the channel model, is commonly used in equalization and filtering applications. But the statistics of channels in real world change over time and inversion of the large dimensional square matrix turns out to be ill-conditioned. To further improve the accuracy of the estimator, Wiener filtering based iterative CE has been investigated [9], [10]. However, this scheme also requires high complexity and knowledge of channel correlations.

The most important research topic in the wireless communications is the adaptive CE where the channel is rapidly time-varying. The time-varying multipath channel can be represented by a tap-delayed line with time varying coefficients and fixed tap spacing. An adaptive algorithm is a process that changes its parameters as it gain more information of its possibly changing environment. Among numerous iterative techniques that exist in the open literature, the popular category of approaches which are obtain from the minimization

of the mean square error (MSE) between the output of the filter and desired signal to perform CE [10-15].

In this paper, normalized least mean (NLMS) square and recursive least squares (RLS) adaptive channel estimator are described for MIMO OFDM systems. These CE methods uses adaptive estimator which are able to update parameters of the estimator continuously, so that knowledge of channel and noise statistics are not required. This NLMS/RLS CE algorithm requires knowledge of the received signal only. This can be done in a digital communication system by periodically transmitting a training sequence that is known to the receiver. Simulation results show that the RLS CE method has better performances compared NLMS CE method for MIMO OFDM systems. In addition, the utilizing of more multiple antennas at the transmitter and/or receiver provides a much higher performance compared with fewer antennas. Furthermore, the RLS CE algorithm provides faster convergence rate compared to NLMS CE method. Therefore, in order to combat the channel dynamics, the RLS CE algorithm is better to use for MIMO OFDM systems.

We use the following notations throughout this paper: bold face and upper lower letter are used to represent matrix and vector. Superscripts $x^*$ and $x^T$ denote the conjugate and conjugate transpose of the complex vector x respectively, and the symbol E(.) denotes expectation.

The remainder of the paper is organized as follows. The NLMS/RLS CE scheme is presented in section II, and its performance is analyzed in section III. Finally, some concluding remarks are given in section IV.

## II. CE METHODS

### A. LMS CE Method

An adaptive algorithm is a process that changes its parameters as it gain more information of its possibly changing environment. Among numerous iterative techniques that exist in the open literature, the popular category of approaches which are obtain from the minimization of the MSE between the output of the filter and desired signal to perform CE as shown in Fig. 1.
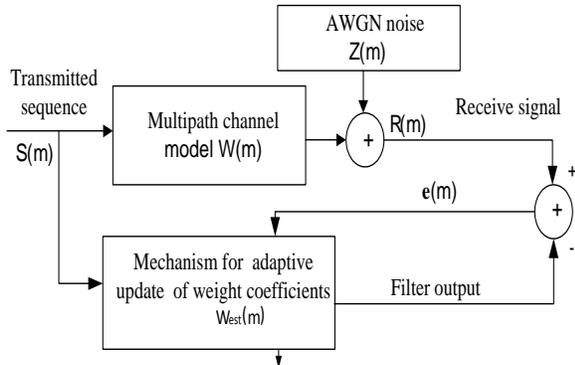


Figure 1. Scheme for adaptive CE.

The signal S(m) is transmitted via a time-varying channel W(m), and corrupted by an additive noise estimated by using

any kind of CE method. The main aim of most channel estimation algorithms is to minimize the mean squared error (MMSE) i.e., between the received signal and its estimate [16-20]. In the Fig 1, we have unknown multipath fading channel, that has to be estimated with an adaptive filter whose weight are updated based on some criterion so that coefficients of adaptive filter should be as close as possible to the unknown channel. The output from the channel can be expressed as:

$$R(m) = \sum_{l=0}^{L-1} W(m,l) \, S(m-l) + Z(m), \qquad (1)$$

where S(m-l) is the complex symbol drawn from a constellation s of the lth paths at time m-l, L is the channel length, Z(m) is the AWGN with zero mean and variance $\sigma^2$. The above equation can be rewritten as vector notation [1]:

$$R(m) = W(m)S(m) + Z(m), \qquad (2)$$

The output of the adaptive filter is

$$Y(m) = W_{est}(m)S(m), \qquad (3)$$

where $W_{est}(m)$ is the estimated channel coefficients at time m. The priori estimated error signal needed to update the weights of the adaptive filter is

$$e(m) = R(m) - Y(m)$$
$$= W(m)S(m) + Z(m) - W_{est}(m)S(m) \qquad (4)$$

This error signal is used by the CE to adaptively adjust the weight vector so that the MSE is minimized. Now the cost function $j(m) = E[e(m)e^*(m)]$ for the adaptive filter structure is

$$j(m) = E[R(m)R(m)] - E[S]R^*W_{est}(m) - R(m)W_{est}(m)$$
$$E[S] - W_{est}(m)W_{est}(m)E[S(m)S(m)]S$$
$$= \sigma_r^2 - C(m)W_{est}(m) - W_{est}(m)C(m)$$
$$+ D(m)W^T_{est}(m)W_{est}(m), \qquad (5)$$

where $\sigma_r^2$ is the variance of the received signal, $C(m) = E[S(m)R(m)]$ is the cross-correlation vector between the tap input vector S(m) and the received signal r(m), and $D(m) = E[S(m)S^T(m)]$ is the correlation matrix of the tap input signal S(m). Now taking the gradient vector with respect to $W_{est}(m)$:

$$\Delta j(m) = -2C(m) + 2D(m) W_{est}(m)$$
$$= -2S(m)R^*(m) + 2S(m)S(m)W_{est}(m). \qquad (6)$$

According to the method of steepest descent, if $W_{est}(m)$ is the tap-weight vector at the mth iteration then the following recursive equation may be used to update $W_{est}(m)$:

$$W_{est}(m+1) = W_{est}(m) - 1/2 \ \eta\Delta j(m)$$

$$= W_{est}(m) + \eta S(m)[R^*(m) - W_{est}(m)S(m)]$$

$$= W_{est}(m) + \eta S(m)e^*(m), \qquad (7)$$

where $W_{est}(m+1)$ denotes the weight vector to be computed at iteration $(m + 1)$ and $\eta$ is the LMS step size which is related to the rate of convergence. The smaller step size means that a longer reference or training sequence is needed, which would reduce the payload and hence, the bandwidth available for transmitting data. The term [ $\eta S(m)e^*(m)$ ] represents the correction factor or adjustment that is applied to the current estimate of the tap-weight vector. In order to improved system performance, taking into account the variation in the signal level at the filter input and selects a normalized step size parameter i.e.,

$$W_{est}(m+1) = W_{est}(m) + \frac{\eta S(m)}{S(m)S(m)}e^*(m), \qquad (8)$$

The iterative procedure is started with an initial guess $W_{est}(0)$. Therefore, the NLMS based CE is least sensitive to the scaling of its input signal variation. Therefore, this algorithm is able to sense the best possible channel coefficients are changing.

*B.  RLS CE Method*

The RLS CE requires all the past samples of the input and the desired output is available at each iteration. The objective function of a RLS CE algorithm is defined as an exponential weighted sum of errors squares:

$$c(m) = \sum_{m=1}^{n} \lambda^{n-m}\mathbf{e}^H(m)\mathbf{e}(m) + \delta\lambda^n \ W^H(m)W(m), \quad (9)$$

where $\delta$ is a positive real number called regularization parameter, $e(m)$ is the prior estimation error, and $\lambda$ is the exponential forgetting factor with $0 < \lambda < 1$. The prior estimation error is the difference between the desired response and estimation signal:

$$\mathbf{e}(m) = H(m) - W^H(m) \ S(m) \qquad (10)$$

The objective function is minimized by taking the partial derivatives with respect to W(n) and setting the results equal to zero.

$$\frac{\delta C(m)}{\delta W(m)} = 0 = -2\sum_{m=1}^{n} \lambda^{n-m}S(m)e^H(m) + 2\delta\lambda^n W(m)$$

$$= -2\sum_{m=1}^{n} \lambda^{n-m}S(m)[H(m) - W^H(m)S(m)]^H + 2\delta\lambda^n W(m)$$

$$W(m)[\sum_{m=1}^{n} \lambda^{n-m}S(m)S^H(m) + \delta\lambda^n I] = \sum_{m=1}^{n} \lambda^{n-m}S(m)H^H(m)$$

$$R_s(m)W(m) = R_{sh}(m)$$

$$W(m) = R_s^{-1}(m) \ R_{sh}(m) \qquad (11)$$

where $\mathbf{R}_s(m)$ is the transmitted auto-correlation matrix

$$R_s(m) = \sum_{m=1}^{n} \lambda^{n-m}S(m)S^H(m) + \delta\lambda^n I = \lambda R_s(m\text{-}1) + S(m)S^H(m)$$

and $\mathbf{R}_{sh}(m)$ is the cross correlation matrix i.e.,

$$R_{sh}(m) = \sum_{m=1}^{n} \lambda^{n-m}S(m)H^H(m) = \lambda R_{sh}(m\text{-}1) + S(m)H^H(m) .$$

According to the Woodbury identity, the above $\mathbf{R}_{sh}(m)$ can be written as

$$R_{sh}^{-1}(m) = \lambda^{-1}R_{sh}^{-1}(m\text{-}1) - \frac{\lambda^{-2}R_{sh}^{-1}(m\text{-}1)S(m)S^H(m)R_{sh}^{-1}(m\text{-}1)}{1+ \lambda^{-1}S^H(m)R_{sh}^{-1}(m\text{-}1)S(m)} \quad (12)$$

For

convenience of computing, let $D(m) = R_{sh}(m)$ and

$$K(m) = \frac{\lambda^{-1}D(m\text{-}1)S(m)}{1+\lambda^{-1}S^H(m)D(m\text{-}1)S(m)} \qquad (13)$$

The K(m) is referred as a gain matrix. We may rewrite (9) as:

$$D(m) = \lambda^{-1}D(m\text{-}1) - \lambda^{-1}K(m)s^H(m)D(m\text{-}1) \qquad (14)$$

So simply (13) to

$$K(m) = D(m)S(m) = R_{sh}^{-1}(m)S(m) \qquad (15)$$

Substituting (14), (15) into (11), we obtain the following RLS CE formula

$$W(m) = W(m\text{-}1) + K(m)[H(m) - W^H(m\text{-}1)s(m)]^H$$

$$= W(m\text{-}1) + K(m)\varepsilon^H(m), \qquad (16)$$

where $\varepsilon(m)$ is a prior estimation error as

$$\varepsilon(m) = H(m) - W^H(m\text{-}1)S(m) \qquad (17)$$

Therefore, equation (17) is the recursive RLS CE algorithm to update channel coefficient.

### III.  ANALYTICAL RESULTS

The error performance of the aforementioned iterative estimation algorithm is explored by performing extensive computer simulations. In these simulations, we consider 2 by 2, 4 by 4, 6 by 6, and 8 by 8 MIMO OFDM systems. The data symbol is based on Q-PSK modulation. The forgetting factor is 0.9 and learning rate is 0.4, and signal to noise ratio is 15 dB. From the simulation results, one can observed that the RLS CE method has better performances compared NLMS CE method. In addition, the utilizing of more multiple antennas at the transmitter and/or receiver provides a much higher BER

performance compared with fewer antennas. Furthermore, the RLS CE algorithm provides faster convergence rate compared to NLS CE method. Therefore, in order to combat the channel dynamics, the RLS CE algorithm is better to use for OFDM systems. But the RLS CE algorithm is suffered from a computational complexity point of view. In addition, the RLS algorithm has the recursive inversion of an estimate of the autocorrelation matrix of the input signal as its cornerstone; problems arise, if the autocorrelation matrix is rank deficient.
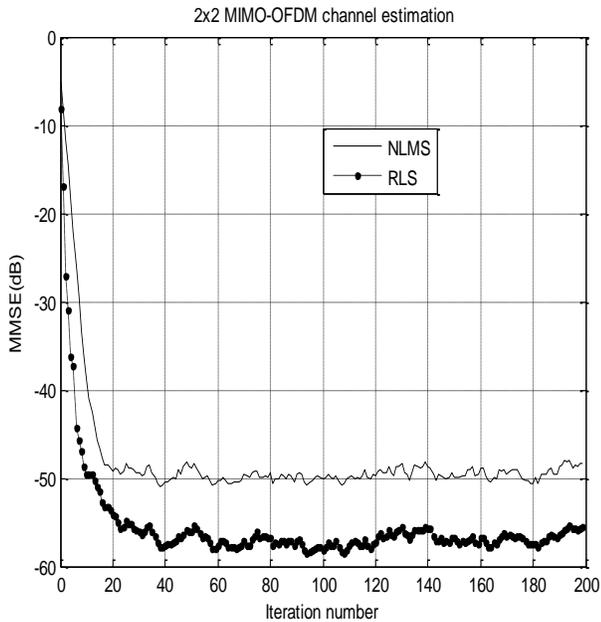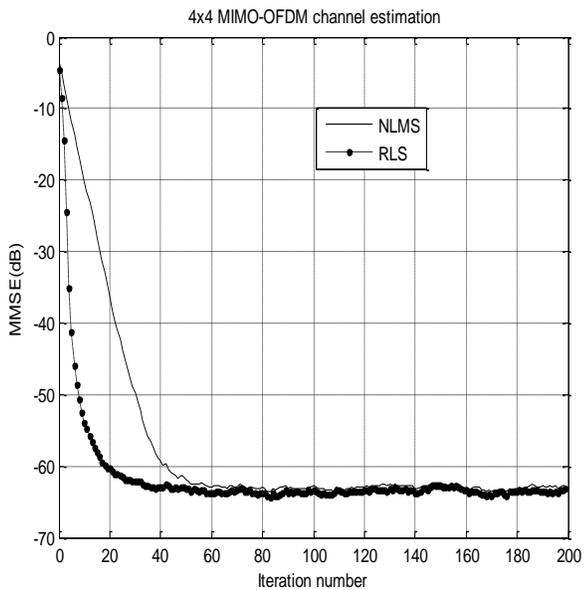


Figure 2. 2 by 2 MIMO systems
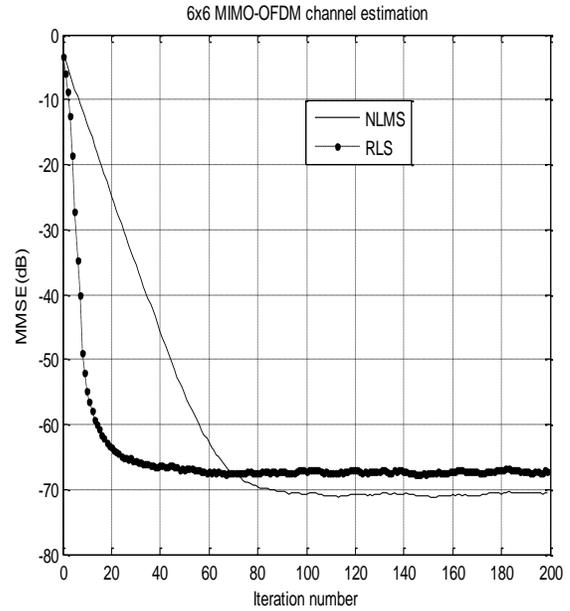


Figure 3. 4 by 4 MIMO systems



Figure 4. 6 by 6 MIMO systems



Figure 5. 8 by 8 MIMO systems

## IV. CONCLUSION

Recently, multiple input multiple output (MIMO) transmission has been well known as one of the most important practical technique to combat fading as well as increase the channel capacity of wireless communication systems. In this paper, NLMS and RLS adaptive channel estimator are described for MIMO OFDM systems. Simulation results demonstrated that the RLS CE method has better performances compared NLMS CE method for MIMO OFDM systems. In addition, the utilizing of more multiple antennas at the transmitter and/or receiver provides a much higher

performance compared with fewer antennas. Therefore, in order to combat the channel dynamics, the RLS CE algorithm is better to used for MIMO OFDM systems.

## REFERENCES

[1] G. J. Foschini and M. J. Gans, "On limits of wireless communications in a fading environment when using multiple antennas," *Wireless Pers. Commun.*, vol. 6, no. 3, pp. 311–335, Mar. 1998.

[2] E. Karami, "Tracking performance of least squares MIMO channel estimation algorithm," *IEEE Trans. On Wireless Comm.* vol. 55, no.11, pp. 2201-2209, Nov. 2007.

[3] V. Pohl, P. H. Nguyen, V. Jungnickel, and C. Von Helmolt, "How often channel estimation is needed in MIMO systems," *in Proc. IEEE Global Telecommun. Conf.*, no. 1, pp. 814–818. Dec. 2003,

[4] S. Addireddy, L. Tong, and H. Viswanathan, "Optimal placement of training for frequency-selective block-fading channels," *IEEE Trans. Inf. Theory*, vol. 48, no. 8, pp. 2338–2353, Aug. 2002.

[5] L. Tong, "Blind sequence estimation," *IEEE Trans. On Commun.,* vol. 43, no. 12, pp. 2986–2994, Dec. 1995.

[6] A. Vosoughi and A. Scaglione, "Channel estimation for precoded MIMO systems," *in Proc. IEEE Workshop Statist. Signal Process.*, pp. 442–445, Sep. 28–Oct. 1, 2003.

[7] J. Li, J. Ma, S. Liu "RLS channel estimation with superimposed training," *in Proc. Int. Con. on Comm. Technology*, 2008.

[8] A. Ancora, C. B. Meili, and D. T. Slock, "Down-sampled impulse response least-squares channel estimation for LTE OFDMA," *in Proc. Int. Con. on Acoustics, Speech, and Signal Processing*, pp. III–293–III–296, Apr. 2007.

[9] L. A. M. R. D. Temino, C. N. I Manchon, C. Rom, T. B. Sorensen, and P. Mogensen, "Iterative channel estimation with robust wiener filtering in LTE downlink," *in Proc. Int. Con. on VTC*, pp. 1-5, Sept. 2008.

[10] J. J. V. D. Beek, O. E. M. Sandell, S. K. Wilsony, and P. O. Baorjesson, "On channel estimation in OFDM systems," *in Proc. Int. Con. on VTC*, vol. 2, pp. 815-819, July 1995.

[11] S. Haykin, "Adaptive Filter Theory," *Prentice-Hall International Inc*, 1996.

[12] F. Adachi, H. Tomeba, and K. Takeda, "Frequency-domain equalization for broadband single-carrier multiple access," IEICE Trans. on Commun., vol. E92-B, no. 5, pp. 1441–1456, May 2009.

[13] S. Yameogo, J. Palicot, and L. Cariou, "Blind time domain equalization of SC-FDMA signal," in Proc. Vehicular Technology Conference, Sept. 2009, pp. 14.

[14] S. Y. Park, Y. G. Kim, and C. G. Kang, "Iterative receiver for joint detection and channel estimation in OFDM systems under mobile radio channels," IEEE Trans. on Vehicular Technology, vol. 53, no. 2, pp. 450–460, Mar. 2004.

[15] J. Berkmann, C. Carbonelli, F. Dietrich, C. Drewes, and W. Xu, "On 3G LTE terminal implementation-standard, algorithms, complexities and challenges," in Proc. Wireless Commun. and Mobile Computing Conference, Aug. 2008, pp. 970–975.

[16] W. C. Jakes, Ed., Microwave mobile communications. New York: Wiley-IEEE Press, Jan. 1994.

[17] B. Karakaya, H. Arslan, and H. A. Curpan, "Channel estimation for LTE uplink in high Doppler spread," in Proc. Wireless Commun. And Networking Conference, Apr. 2008, pp. 1126–1130.

[18] R. C. Alvarez, R. Parra-Michel, A. G. O. Lugo, and J. K. Tugnait, "Enhanced channel estimation using superimposed training based on universal basis expansion," IEEE Trans. on Signal Process., vol. 57, no. 3, pp. 1217–1222, Mar. 2009.

[19] A. Kalayciogle and H. G. Ilk, "A robust threshold for iterative channel estimation in OFDM systems," Radio Engineering journal, vol. 19, no. 1, pp. 32–38, Apr. 2010. 10

[20] Q. Li, G. Li, W. Lee, M. il Lee, D. Mazzarese, B. Clerckx, and Z. Li, "MIMO techniques in WiMAX and LTE: a feature overview," IEEE Commun. Magazine, vol. 48, no. 5, pp. 86–92, May. 2010.

# The Impact of Social Networking Websites to Facilitate the Effectiveness of Viral Marketing

Abed Abedniya

Faculty of Management (FOM)
Multi Media University (MMU)
Cyberjaya, Malaysia
Abed.abedniya@gmail.com

Sahar Sabbaghi Mahmouei

Institute of Advanced Technology (ITMA)
University Putra Malaysia (UPM)
Serdang, Malaysia
Sabbaghi.sahar@gmail.com

*Abstract*—**The Internet and the World Wide Web have become two key components in today's technology based organizations and businesses. As the Internet is becoming more and more popular, it is starting to make a big impact on people's day-to-day life. As a result of this revolutionary transformation towards the modern technology, social networking on the World Wide Web has become an integral part of a large number of people's lives. Social networks are websites which allow users to communicate, share knowledge about similar interests, discuss favorite topics, review and rate products/services, etc. These websites have become a powerful source in shaping public opinion on virtually every aspect of commerce. Marketers are challenged with identifying influential individuals in social networks and connecting with them in ways that encourage viral marketing content movement and there has been little empirical research study about of this website to diffuse of viral marketing content. In this article, we explore the role of social network websites which has influence on viral marketing, and the characteristics of the most influential users to spread share viral content. Structural equation modeling is used to examine the patterns of inter-correlations among the constructions and to empirically test the hypotheses.**

*Keywords-Social networks website, viral marketing, structural equation modeling*

## I. INTRODUCTION

Given the increasing popularity of the internet nowadays, businessmen and entrepreneurs have started to explore the concept of marketing on the Web. The World Wide Web has an exorbitance of ways to promote a business and most importantly, the internet caters to a wide range of audience that is perhaps interested to the business' products and services. The most powerful and influential form of advertising is passing the information from one person to another. This form of advertising called "Word-Of-Mouth". Many marketers and researchers believe that word-of-mouth Communication has become a hot subject in marketing. Word-of-mouth is becoming a main base for interactive marketing communication among offline communication strategists. Many reasons have been presented to explain this growth in marketing strategists. Word-of-mouth marketing is such a successful marketing strategy because it fostered "familiarity, personal connection, care and trust" between the consumer and

the translator of the information [10]. Another reason for growing this strategy is that many people like to talk about their purchase with products and services for a variety of reasons. Psychologists believe these customer behaviors may arise through ownership or a need to share their purchase experiences in order to help others. These conversations are then passed to family, friends and other people in social networks [2]. While the underlying principle of Word-Of-Mouth marketing is well-established and acknowledged [39]; [51], the Internet fosters new marketing strategies [3], one of which is viral marketing. Viral marketing is said to be the electronic version of traditional word-of-mouth advertising and product communication [5]. Viral marketing can be described as a marketing technique that uses e-mail messages containing powerful advertising messages and promotional offers that are specifically designed for its recipients to forward to their family, friends, or others on their e-mail contact list [5]. The reason behind providing viral marketing with its specific name is because like human and computer viruses it also "multiplies rapidly in a cell, commandeering the cells resources to do the virus' bidding" [15], [41]. At the base of viral marketing is the transmission of viral message through internet users by peers. This is an opportunity for Marketers that can transfer information between internet users without the involvement during these transmissions. Viral marketing works because friends are better at target marketing than any database [7]. Due to advancements in computer technology and internet people all over the world can now interact and communicate with virtually anyone else who has access to a computer and the internet. These advancements in communication and technology have opened up huge opportunities for businesses to appeal to much larger markets than ever before. The Internet, with the help of instant messaging and social networking sites, has hyper-accelerated the rate at which people talk to each other and has greatly expanded the range of topics they discuss and how they discuss them. According to Alexa.com four social network websites (MySpace.com, Facebook.com, Twitter.com and Hi5.com) belong to the top ten of the global traffic ranking [53]. The nature of these social networking websites assist people to convey a message through a potentially self-replicated, growing campaign where, ideally, one person tells two people who tell two more people each. Among the global websites Facebook.com, is the leading social

network website that currently has more than 55 million active users, with an average of 50,000 new registrations per day [14]. StudiVZ.net, the most popular social network website in Germany, witnessed an impressive increase in user subscriptions since its launch in 2005, currently having more than 3 million members [46]. Social networking is rapidly expanding; Williamson in 2008 estimates that there was an 11 percent increase of people visiting social networking sites between 2007 and 2008, with "79.5 million people—41% of the U.S. Internet user population" visiting the sites in 2008[53]. Furthermore, the trend will continue, and, by 2013, the number will increase to 52 percent. According to an eMarketer study in 2008, nearly six out of ten United States users now communicate with businesses and believe that the businesses must "interact with their consumers" and "deepen the brand relationship" via social networking website. Its big instrument for viral marketers that develop their process because they can be driven by content integrated into consumer/user profiles and such content is increased through new user acquisition and retention.

## II. LITERATURE REVIEW

In this section, we review previous studies on the social network websites and viral marketing, and provide a set of hypotheses to examine the motivational characteristics of these websites on facilitated the viral marketing.

### A. Viral Marketing

In 1997, Juvertson and Draper was developed the viral marketing by describe free email service for Hotmail, they explain term "viral marketing" simply as "network-enhanced word-of-mouth" [25]. However after them, many researchers used different terminology to explain what viral marketing is. According to [49], some of the terminology used to describe electronic WOM includes "Interactive Marketing" [4], viral marketing[25], Internet communication [8], Internet word-of-mouth and word-of-mouse [17], online feedback mechanisms [8], stealth marketing [27], buzz marketing [43], electronic word-of-mouth communication [20], interactive or electronic word-of-mouth advertising [37], and electronic referral marketing [9],defines viral marketing as "any strategy that encourages individuals to pass on a marketing message to others, creating the potential for exponential growth in the message's exposure and influence"[52].

However, viral Marketing using informal communication among consumers in social networks to promote and grow products, services and brands. Viral marketing as any strategy that encourages individuals to pass on a marketing message to others, creating the potential for exponential growth in the message's exposure and influence [52]. Base on social theory suggests that people tend to connect with others who share common interests [19]. Purchasing decisions are often strongly influenced by people who the consumer knows and trusts by his or her social network and their community. Viral messages can reach and potentially influence many receivers, and are usually perceived by consumers to be more reliable and credible than firm-initiated ones, since the senders of viral are mostly independent of the market [55]. After the emergence of the Internet, marketers tend to look at this phenomenon of how

they can increase awareness about their services and products with immediate and low cost in many countries with different cultures. [36] researcher suggest viral contents such as joke, picture, game and video through the Internet, often distributed through independent third-part sites, are usually personal, more credible than traditional advertising and we must employ these material for viral in executions by website. Despite the increasing shift of advertising spending to viral marketing [27], the factors critical to viral marketing effectiveness remain largely unknown to both marketing academics and practitioners [16].

Now with increased social networking in the internet, viral marketing could encourage consumers to spread their message and entertaining media to their friends and also in turn encourage their friends to diffuse forward these message in a largely chain reaction of consumer awareness. Furthermore, viral marketing is often also stealth marketing, encouraging customers to feel they just happened to hear about the product or service rather than to feel directly marketed to [54]. Considering the viral marketing capabilities to increase customer awareness of our products and services can be concluded that there can be significant benefits to be gained from viral marketing. Before review the influential factors on viral marketing we should determine the main viral marketing characteristic that can be more effective on success of its campaign. In order to implementing this purpose we done research and after having studied we find main following characteristic as important indicator:

- Rapid diffusion to audience reaches:

  Kaikati, Helm, Welker believe viral marketing can reach audience and spread exponentially within a short period of time [26], [18]. This rapid diffusion can significantly boost the speed of the adoption of the marketed product or service [11]. Viral marketing can help achieving substantial audience, reaching as marketers get access to diverse audiences through social contacts [18].

### B. Marketing through Social Networks sites

Nowadays, social sites have become most popular place for internet user on the Internet. These sites are instrument for building virtual communities with same age, education, lifestyles, idea and interests or to create a variety of activities for individual. Today, these sites have been able to break geographical boundaries and set many people with different cultures and nationalities together. These sites are helped to be broken down into sub networks, based on demographic or geographical preferences.

Marketers believe that members of social sites who share information with other members and friends are best target for participation in viral marketing. Their involvement in these social sites allows marketer to spread more the viral content because they naturally want to share information to other members and send interesting content to friends.

More than half of social network site users already tell members of their social network about products they have used [29]. Online social network members are also more interested in viewing the profile pages of companies [29]. The online

social networks are favorable places for executing a purpose to reach groups of consumers who share common interests and comments same viral messages that can be spread quickly by consumers who truly share common interests and preferences.

Many factors will make social network website an interesting instrument for marketing strategy. Many users can join easily to these websites without pay money where they could share their opinion about anything and also make recommendation on these sites. Some research showed 78% of global consumers say they trust and believe other people's recommendations for products and services more than any other medium [10], actually it has been shown that many consumer attempt to know another persons' opinion in the social network websites when considering the purchase of products and services because they belief and trust on their friends opinion. Many members believe that their friends in these sites are better resources than companies advertising for buy products or services. Members of social networks serve two roles; they both supply and consume content. The creators of content are typically highly engaged consumers and, as a result, influential [10]. When one person who affects or influences others perceives message as valuable, he or she could be turned into a viral. This event make these website as powerful weapon for marketers, because users don't feel that the information is being pushed at them, but referred to them by a trusted friend in a trusted network.

Nevertheless, in spite of the fact that we have been studied about word-of-mouth and offline social network but there is little research that studies the phenomenon of viral marketing through social network websites and the factors affecting the effectiveness of this kind of marketing. Among the limited number of studies with this line of research, we found social network websites influence facilitated viral marketing. Therefore we want to explain how the motivational characteristics of these websites have impact on main character of viral marketing.

## III. CONCEPTUAL FRAMEWORK AND HYPOTHESIS

Based on our review of the related research, we found very few studies that investigate the impact of motivational characteristic of social network website on marketing special viral marketing. The proposed research model is depicted in *"Fig. 1"*. As we know the primary users' purpose of these websites are used for entertainment to pleasure purposes rather than utilitarian purposes. Our model which we call the Social Network Website Influence model (SNWI) focuses on motivational factor of these websites. It is expected that natural motivators and social influences of user behavior are dominant predictors of usage [44]. In this model, the networks' characteristics include perceive playfulness, critical mass, community driven, peer pressure, perceive ease of use and perceive usefulness involvement. The dependent variable is the Rapid diffusion to audience reaches. To build our model, we examine these factors on dependent variable have paid. We suppose that these factors will have a direct effect on rapid diffusion to the audience reaches of viral marketing content. In order to prove that we prepare model that includes six motivational components. The model used for the

effective components of these websites to end-users' intentions in order to publish viral content among their friends.

### A. Playfulness

Current or potential users of social network website believe that perceived playfulness by social network website will bring him/her a sense of enjoyment and pleasure. In the study of global network, [38] describe that the playful actions expressed by social network websites such as Facebook Applications could be viewed as representations of characteristics of online playfulness. Researchers are believed that perceived playfulness exhibits an important role in the usage or continuance use of Social network website. By base on [32] and [30], "Perceived playfulness" is defined as intrinsically enjoyable or interesting". They found that perceived playfulness has a direct effect on extent viral marketing content. Therefore, we propose that:

H1: Playfulness website has a significant positive effect on rapid diffusion content in viral marketing.

### B. Critical mass

Critical mass is a subjective measure of the point where enough of one's friends participate in a social network to make it valuable. In the context of Social network websites, this subject refer to user perceives this website to have a significant number of users that they can associate with them due to friendship, common interests, and share content for example. A critical mass of users who actively occupied content transfer, information exchange and knowledge sharing activities is crucial to keep users online community up and running over time. If a social network website may claim they have many users and member and current or potential users perceive enough active members that they can associate with, consequently critical mass achieved or sustained for those members. But if these members of social network website perceive there are not enough active members for associate with, therefore this subject has not been achieved or sustained for member. Researchers believe that perceived critical mass has an initiation to intention to use other communication technologies, such as groupware [31], and instant messaging [45]; hence, we expect perceived critical mass to impact intent to use social network website, as hypothesized below:

H2: Perceived critical mass has a significant positive effect on rapid diffusion content in viral marketing.

### C. Community-driven

One of the strong motivational characters of social network website is community driven. This character is a new way to find, share and transfer information with internet users. Social network website with this character allows their member to ask any question and to receive answer from other members.

Some researchers on information system and information technology pay attention to community driven as a good example of harnessing user generated content [32],[28].They believe that community driven provides the best opportunity for social network website Internet users with non-stop access to any kind of information from multiple domains. Nowadays we can find many sub communities of people who share

commonalities between social network website, such as fan of movie, music, alumni of particular university and an economic welfare group. Many users during these communities find old friends that they lost connection with them many years and reconnect with them or discover new friends. We argue that community driven in social network website can influence between members to transfer and share content in internet. This makes the next hypothesis:

H3: community driven has a significant positive effect on rapid diffusion content in viral marketing.

### D. Peer pressure

Push factors, such as peer pressure, were identified as being a strong influence upon decisions to join a social network website. Many members special teenager typically join the social network website such as Facebook because a friend invites them to join [6]. Some researchers believe external influence such as peer pressure is important external determinants that should be accounted on participation in social network [40]. Social network websites grab more members and expand their network influence that makes people connect together and participate in community. This leads to the following hypothesis:

H4: the degree of peer pressure to participate has significant positive effect on rapid diffusion content in viral marketing.

### E. Perceived Ease and  F.  Use and Usefulness

Perceived ease of use is related to user extent to which people find using a new technology will be comfortable and perceived usefulness is identified with a person who believes that using of a new technology will increase their productivity or performance. Accordingly, both perceived usefulness and perceived ease of use are likely to affect users' self-disclosure intentions [35]. Some researchers believe that both perceived ease of use and usefulness are strong determinants of user acceptance, adoption, and usage behavior [42], [50]. If users find a social networking website very useful and done without waste of time to communicate with others members, they are attracting to regularly update their profiles and disclose more information to others. In similar way, if members believe very few efforts to use the site, they may also frequently update their status and share more personal information to the public. Perceived ease of use is a belief that it would be easy to acquire the knowledge for using the information technology or system [35]. Perceived usefulness is a belief that the target information technology or system will help the user in performing his or her task [35]. This leads to the following hypothesizes:

H5: Perceived ease of use has a significant positive effect on rapid diffusion content in viral marketing.

H6: Perceived usefulness has a significant positive effect on rapid diffusion content in viral marketing.



Figure 1.   Research Model for Social Network Site influence on Viral Marketing

## IV.   METHODOLOGY

### A.  Data Collection

In order to perform this research, we used an online survey to collect data. The data was collected through online survey questionnaires distributed to students enrolled at major Malaysian universities. In the questionnaire survey we define social network website and its motivational characteristics. We introduce some popular social network websites for respondents such as (Facebook, Friendster, Tagged or Myspace) as a frame of reference for their responses.  In our study around 150 students were involved during three weeks, of these participants, 41% were females and 59% were males. We focus on adult person because they more potentials to use of these website therefore about 39% of the sample between the age 18-20 and 61% between 21-30 ages

### B.  Development of Measurement Scales

All theoretical constructs in the study involved and operationalized by multiple item scales. In order to maintain content validity of the adopted scales in the field of social network websites were verified the scales. Some of the scale items were slightly rephrased to reflect the current research context. Additionally, we removed low correlation coefficient item. Table 1 summarizes sources used to operationalize model constructs. All needs-related items were attached firmly on a seven-point Likert scale.

TABLE I.   CONSTRUCT OPERATIONALIZATION.

| Construct Name | Construct Type | Sources |
|---|---|---|
| Rapid diffusion to audience reaches | endogenous | Helm (2000) and Dobele (2005) |
| Playfulness | exogenous | Moon and Kim (2001) and Rao (2008) |
| Critical mass | exogenous | Van Slyke (2007) |
| Community-driven | exogenous | Moon(2007) |
| Peer pressure | exogenous | Ajzen (2002) and roger (2003) |
| Perceived Ease of Use | exogenous | Venkatesh and Davis (2000) and Ohbyung Kwon (2010) |
| Usefulness | exogenous | Venkatesh and Davis (2000) and Ohbyung Kwon (2010) |

TABLE II.   RESULTS FROM THE CFA OF STUDY CONSTRUCTS

| Construct | Composite Reliability | Average Variance Extracted (AVE) | Cranach's Alpha |
|---|---|---|---|
| Playfulness | 0.885 | 0.529 | 0.835 |
| Critical mass | 0.847 | 0.649 | 0.706 |
| Community-driven | 0.855 | 0.545 | 0.792 |
| Peer pressure | 0.844 | 0.524 | 0.748 |
| Perceived Ease of Use | 0.901 | 0.694 | 0.843 |
| Usefulness | 0.891 | 0.673 | 0.825 |
| Rapid diffusion to audience reaches | 0.918 | 0.693 | 0.882 |

Notes: All factor loadings are significant at $p = .05$ (i.e., $t > 2.0$); b: Only remaining items
After the purification process are shown.

## C. Analyses and Results

The purpose of the statistical analysis is to explain the relationship between the dependent variable and independent variables. In our study dependent variable or endogenous variable is Rapid Diffusion to Audience Reach in social network website and independent variables or exogenous variables are playfulness, critical mass, community driven, peer pressure, perceived ease of use and perceive usefulness. We tested our model with structural equation modeling (SEM) technique. Structural Equation Modeling (SEM) is a statistical technique for testing and estimating causal relations using a combination of statistical data and qualitative causal assumptions. In this research we using 2 step SEM approach using Structural Equation Modeling Software EQS. At first to determine the composite reliabilities, convergent and discriminate validity of the multi-item measure, we purified by confirmatory factor analysis (CFA). Second, we used the structural model for evaluate the proposed hypotheses.

### 1) Evaluation of the Measurement Model

At first we evaluated Convergent Validity and Discriminant Validity in order to construct validity. Convergent validity is the degree to which an operation is similar to other operations that it theoretically should also be similar to. High correlations between the test scores would be evidence of a convergent validity. Discriminant validity describes the degree to which the operationalization is not similar to other operationalizations that it theoretically should not be similar to. For convergent validity we must evaluated three criteria: indicator reliability, composite reliability and average variance. We tested these criteria for each indicator and latent variables. As can be seen in table 2, Average Variance Extracted (AVE) for our construct between 0.524 and 0.694 and Cronbach's Alpha is between 0.706 and 0.882. When Cronbach's Alpha is above 0.7, showing that Internal Consistency is assured [34]. Putting together the results from the different criteria, Convergent Validity can be assumed. To ensure Discriminant Validity, require that the AVE for any latent variable has to be bigger than the squared correlation between this variable and all other latent variables in the model [13]. As can be inferred from Tables 2 and 3, this requirement is indeed ensured for all latent variables.

TABLE III.   CORRELATION BETWEEN LATENT VARIABLES

| Construct | diffusion to audience reaches | Playfulness | Critical mass | Community-driven | Peer pressure | Perceived Ease of Use | Usefulness |
|---|---|---|---|---|---|---|---|
| Rapid diffusion to audience reaches | 1.000 | | | | | | |
| Playfulness | 0.498 | 1.000 | | | | | |
| Critical mass | 0.514 | 0.662 | 1.000 | | | | |
| Community-driven | 0.519 | 0.404 | 0.627 | 1.000 | | | |
| Peer pressure | 0.429 | 0.463 | 0.530 | 0.513 | 1.000 | | |
| Perceived Ease of Use | 0.515 | 0.501 | 0.682 | 0.650 | 0.553 | 1.000 | |
| Usefulness | 0.448 | 0.475 | 0.512 | 0.587 | 0.555 | 0.620 | 1.000 |

### 2) Evaluation of the Structural Model

TABLE IV.   SUMMARY OF HYPOTHESES TESTS

| Hypothesis | Supported/Not Supported |
|---|---|
| H1: Playfulness has a significant positive effect on rapid diffusion context in viral marketing. | Supported |
| H2: Perceived critical mass has a significant positive effect on rapid diffusion context in viral marketing. | Supported |
| H3: community driven has a significant positive effect on rapid diffusion context in viral marketing. | Supported |
| H4: the degree of peer pressure to participate has significant positive effect on rapid diffusion content in viral marketing. | Supported |
| H5: Perceived ease of use has a significant positive effect on rapid diffusion context in viral marketing. | Supported |
| H6: Perceived usefulness has a significant positive effect on rapid diffusion context in viral marketing. | Supported |

We used EQS software to test the hypotheses for structural model analysis. Perceived community driven has the strongest effect on rapid diffusion in social network website (M = .47, p < .01), followed in order perceived playfulness by (M = .36, p < .01), critical mass (M = .33, p < .01), perceived usefulness (M = .29, p < .01), perceive ease of use (M = .11, p < .01 and peer pressure (M=.11, p<.01). Therefore all hypotheses are supported. As anticipated, the usage of social network websites is founded to be positively effect on rapid diffusion of content to audience reach in viral marketing.

## V. DISCUSSION

Our study examined the motivational characters of social network websites that contributing the intention to use and diffusion content of viral marketing. The validity of our model (SNWI) and the relationship among its constructs were tested using structural equation modeling. Our model demonstrated that social network websites are significantly influenced by their motivational characters on viral marketing. As the result of our model, Social Network Website Influence (SNWI), community driven is an important driver of rapid diffusion content in social network website with a significant positive path coefficient of 0.47. Social network websites with high level of community driven are predicted to be more likely to share and diffusion viral content. The satisfaction of the perceive playfulness through social network websites is another important driver of rapid diffusion of viral contents with a significant path coefficient of 0.36. Our result shows that critical mass with a significant positive path coefficient of 0.33 is one of influenced character in these websites. Social network website with high level of critical mass has more influence on potential users to believe and participate in viral activity. Perceive usefulness with a significant positive path coefficient of 0.29 showed when users find these websites as usefulness technology will increase their productivity or performance, they will do more activity in these websites and lead to more share and diffusion viral contents. The perceived ease of use and peer pressure by a social network websites was found to have no impact on user participation, as the path coefficient, though positive, was insignificant. But we believe these characters are influenced but Intensity of less than community driven and perceive playfulness. Consequently, social network websites plays an important role in viral marketing influence. Therefore, it can be assumed that social network websites' characters can potentially have a stronger effect on user to share viral content. Social network website is based on network effects, which increases the possibility that a message will reach the right people.

## VI. CONCLUSIONS, RECOMMENDATIONS AND FUTURE STUDY

People are increasingly using the Internet to communicate with others, bring out information, find recommendations, increase knowledge and interact with family friends. This study demonstrates that social networking websites by motivational characters are gaining rapidly in popularity and can fill the role of a way to reach and interact between members. Viral marketing is a sufficient marketing strategy and an important tool for all businesses with limited marketing resources. We

can assimilate social network websites and marketing strategies in order to changing and developing consumer behavior expectations to reach company goal. Consumers are and will continue to share their opinions on brands and products with or without company interaction. Hence, it is best opportunity for company to be corporate or engaged in these communications sharing and influencing on it with positively influencing the message therefore facilitate action and brand awareness through integrated viral marketing strategies. Anything is available to embrace and profit from the incorporation of viral marketing and social network website into an integrated marketing and communications strategy. This strategy provides an opportunity to increase brand awareness and exponentially employ the most influential marketing strategy. Entry into social network websites and viral marketing has low barriers to entry and consequently any company can do it. Based on this research and other studies, we recommend that business owners focus on these websites for brand awareness and introduce new products by spreading viral content. This strategy has low cost and more influence between customers for marketers because customers more trust to their friends rather than company advertisement.

Future study is needed to better isolate this effect, perhaps in a more controlled on member characters. It may also prove fruitful to consider alternate models with member behavior, and to draw parallels with member character influence on viral marketing across social network websites. It can be more effective rather than websites character because viral campaign can predict member interaction and make attractive contents for more influence for rapid diffusion.

## REFERENCES

[1] Ajzen, I. (2002). Constructing a TpB Questionnaire: Conceptual and MethodologicalConsiderations.http://www.people.umass.edu/aizen/pdf/tpb.measurement.pdf, last access on 10.11.2007.

[2] Allsop, D.T., Bassett, B.R. and Hoskins, J.A. 2007, "Word of mouth research: principles and applications", Journal of Advertising Research, Vol. 47 No. 4, pp. 398-411.

[3] Arnott,D.C.&Bridgewater,S.(2002).Internet, interaction and implications for marketing. Marketing Intelligence and Planning, 20(2), 86-95.

[4] Blattberg, R.C. and Deighton, J. (1991), "Interactive marketing: exploiting the age of addressability", Sloan Management Review, Vol. 33 No. 1, pp. 5-14.

[5] Bidgoli, Hossein (2004), The Internet Encyclopedia. Hoboken, NJ: John Wiley and Sons. p. 568

[6] Boyd, D. M. and Ellison, N. B.,"Social Network Sites: Definition, History, and Scholarship",Journal of Computer-Mediated Communication, 13(1), 2007, pp. 210-230.

[7] Bulkeley, W. (2002). " E-commerce: Advertisers find a friend in viral marketing [Electronic version]. ", The Wall Street Journal Europe, p. 25.

[8] D. Cruz, Honda, Slough, C. Fill,(2008), "Evaluating viral marketing: isolating the key criteria", Emerald Group Publishing Limited 0263-4503-Marketing Intelligence & Planning Vol. 26 No. 7, 2008-pp.743-758

[9] De Bruyn, A. and Lilien, G.L. (2004), "A multi-stage model of word-of-mouth through electronic referrals", eBusiness Research Centre Working Paper, The Pennsylvania State University, State College, PA, February.

[10] Datta, P., Chowdhury, D., & Chakraborty, B. (2005). Viral marketing: New form of wordof- mouth through internet [Electronic version]. The Business Review, Cambridge, 3(2), 69-75.

[11] Dobele, A., Lindgreen, A., Beverland, M., Vanhamme, J. & Van Wijk, R. (2007), "Why pass on viral messages? Because they connect emotionally", Business Horizons, Vol. 50 No. 4, pp. 291-304

[12] D.Sledgianowski,S.kulviwat,(2009), "Using social network sites: The effects of playfulness, critical mass and trust in a hedonic context", Journal of Computer Information Systems,pp.74-83

[13] Fornell, C. and Larcker, D. F. (1981). Evaluating Structural Equation Models with Unobservable Variables and Measurement Error. Journal of Marketing Research, 18, 39-50.

[14] Facebook.com Official Statistics Page. Retrieved Sept. 12, 2009. http://facebook.com/press/info.php?statistics.

[15] Gattiker, Urs (2004), The Information Security Dictionary. New York: Springer. p. 348

[16] Godes D, Mayzlin D, Chen Y, Das S, Dellarocas C, Pfeiffer B, Libai B, Sen S, Shi M, Verlegh P,(2005). " The firm's management of social interactions", Mark Lett 2005;16(3/4):415–28.

[17] Goldenberg, J., Libai, B. and Muller, E. (2001), "Talk the network: a complex systems look at the underlying process of word-of-mouth", Marketing Letters, Vol. 12 No. 3, pp. 211-23.

[18] Helm, S. (2000), "Viral marketing – Establishing customer relationship by „word-of-mouse'"", Electronic Markets, Vol. 10 No. 3, pp. 158-161

[19] Hill, S., Provost, F. and Volinsky, C. Network-Based Marketing: Identifying Likely Adopters via Consumer Networks. Statistical Science, 21, 2, (2006), 256-276.

[20] Hennig-Thurau, T., Gwinner, K.P., Walsh, G. and Gremler, D.D. (2004), "Electronic word-of-mouth via consumer-opinion platforms: what motivates consumers to articulate themselves on the internet?", Journal of Interactive Marketing, Vol. 18 No. 1, pp. 38-52.

[21] Ilie, V., Van Slyke, C., Green, G., and Lou, H., "Gender Differences in Perceptions and Use of Communication Technologies: A diffusion of innovation approach", Information Resources Management Journal, 18(3), 2005, pp. 13-31.

[22] J.Cao,T.Knotts,J. Xu.M.Chau,(2009), "Word of Mouth Marketing through Online Social Networks", Proceedings of the Fifteenth Americas Conference on Information Systems, San Francisco, California August 6th-9th 2009

[23] Jason Y.C. Ho, Melanie Dempsey,(2010)."Viral marketing: Motivations to forward online content". Journal of Business Research 63 (2010) pp.1000-1006

[24] J.Cao, t. Knotts,J.Xuz ,M.Chau,(2009), "Word of Mouth Marketing through Online Social Networks", Americas Conference on Information Systems (AMCIS)

[25] Juvertson, S. and Draper, T. (1997), "Viral marketing", Draper Fisher Juvertson website, available at: www.dfj.com/cgi-in/artman/publish/printer _steve_tim_may97.html (accessed 12 March 2006).

[26] Kaikati, A.M. & Kaikati, J.G. (2004), "Stealth marketing: How to reach consumers surreptitiously", California Management Review, Vol. 46 No. 4, pp. 6-22

[27] Knight K. Forecast: WOM to exceed $1 billion in 2007. BizReport 2007. http://www.bizreport.com/2007/11/forecast_wom_to_exceed_1_billion_in_2007.html.

[28] L.Huang, Z.Xia,(2009), "Measuring User Prestige and Interaction Preference on Social Network Site", 2009 Eigth IEEE/ACIS International Conference on Computer and Information Science

[29] Li, C. (2007), "How consumers use social networks". Cambridge, MA: Forrester Research,

[30] Lin, C. and Yu, S. "Consumer Adoption of the Internet as a Channel: The Influence of Driving and Inhibiting Factors, Journal of Academy of Business, 9(2), 2006, pp. 112-117.

[31] Lou, H., Luo, W., and Strong, D. "Perceived Critical Mass Effect on Groupware Acceptance. European Journal of Information Systems, 9(2), 2001, pp. 91-103.

[32] Moon, J. and Kim, Y., "Extending the TAM for a World-Wide-Web Context", Information & Management, 38, 2001, pp. 217-230.

[33] Moon, Jeremy, and David Vogel. 2008. "Corporate Social Responsibility, Government, and Civil Society." In The Handbook of Corporate Social Responsibility, ed. A. Crane, A. McWilliams, D. Matten, J. Moon and D. Siegel. Oxford: Oxford University Press.

[34] Nunnally, J. C. (1978). Psychometric Theory. 2nd Edition. McGraw-Hill, New York, NY. Oleson, M. (2004). Exploring the Relationship between Money Attitudes and Maslow's Hierarchy of Needs. International Journal of Consumer Studies, 28 (1), 83-92.

[35] Ohbyung Kwon, Yixing Wen, (2010), "An empirical study of the factors affecting social network service use", Computers in Human Behavior 26 (2010) 254–263

[36] Porter, L., & Golan, G. J. (2006). " From subservient chickens to brawny men: A comparison of viral advertising to television advertising"

[37] Phelps, J.E., Lewis, R., Mobilio, L., Perry, D. and Raman, N. (2004), "Viral marketing or electronic word-of-mouth advertising: examining consumer responses and motivations to pass along e-mail", Journal of Advertising Research, December, pp. 333-48.

[38] Rao, V.(2008), " Facebook Applications and playful mood: the construction of Facebook as a "third place"", Proceedings of the 12th international conference on Entertainment and media in the ubiquitous era, Tampere, Finland, 2008.

[39] Richins,M.L. (1983). Negative word-of-mouth by dissatisfied customers: A pilot study. Journal of Marketing, 47(1), 68-78.

[40] Rogers, E. (2003). Diffusion of Innovations. Fifth Edition. Free Press, New York, NY. Schneider, B. and Alderfer, C. (1973). Three Studies of Measures of Need Satisfaction in Organizations. Administrative Science Quarterly, 18 (4), 489-505.

[41] Robert H. Luke, Jarrod Freeman.(2008), " viral marketing: a significant learning opportunity marketing strategy and business practive enviroment", MMA Fall Educators' Conference – 2008

[42] Taylor, S., P. A. Todd.(1995), " Understanding information technology usage: A test of competing models", Inform. Systems Res. 6(2) 144–176.

[43] Thomas, G.M. (2004), "Building the buzz in the hive mind", Journal of Consumer Behaviour, Vol. 4 No. 1, pp. 64-72.

[44] Van der Heijden, H.,"Factors Influencing the Usage of Websites: The Case of a Generic Portal in The Netherlands", Information & Management, 40, 2003, pp. 541-549.

[45] Van Slyke, C., Illie, V., Lou, H., and Stafford, T., "Perceived Critical Mass and the Adoption of a Communication Technology", European Journal of Information Systems, 2007, 16, pp. 270-283.

[46] V.Bolotaeva, T.Cata,(2010), "Marketing Opportunities with Social Networks", Journal of Internet Social Networking and Virtual Communities ,Vol. 2010 (2010), Article ID 109111, 8 pages

[47] Venkatesh, V. and Davis F.D., "A Model of the Antecedents of Perceived Ease of Use: Development and Test", Decision Sciences, 27, 1996, pp. 451-481.

[48] Venkatesh, V., Morris, M. G., Davis, G. B., and Davis, F. D., "User Acceptance of Information Technology: Toward a Unified View", MIS Quarterly, 27(3), 2003, pp. 425-478.

[49] Vilpponen, A., Winter, S. and Sundqvist, S. (2006), "Electronic word-of-mouth in online environments: exploring referral network structure and adoption behaviour", Journal of Interactive Advertising, Vol. 6 No. 2, pp. 71-86.

[50] Viswanath Venkatesh,(2000), "Determinants of Perceived Ease of Use: Integrating Control, Intrinsic Motivation, and Emotion into the Technology Acceptance Model", Information Systems Research, _ 2000 INFORMS Vol. 11, No. 4, December 2000, pp. 342–365

[51] Wilson, J.R. (1991). Word of Mouth Marketing, J. Wiley, New York.

[52] Wilson, R.F. (2000), "The six simple principles of viral marketing", Web Marketing Today, Vol. 70, pp. 1-3.

[53] Williamson, D. (2008). Social Network Marketing: Slow Growth Ahead for Ad Spending. eMarketer Research. http://www.emarketer.com/Report.aspx?code=emarketer_2000541

[54] Xu, Y., Zhang, C., Xue, L., & Yeo, L. L. (2008). " Product adoption in online social network. ", Proceedings of the International Conference on Information Systems 2008, Paris, France.

[55] Brown, J., Broderick, A.J., and Lee, N. (2007). Word Of Mouth Communication within Online Communities: Conceptualizing the Online Social Network. Journal of Interactive Marketing, 21(3), 2-20.

AUTHORS PROFILE

Abed Abedniya is currently doing his Master degree in faculty of Management (FOM), Multi Media University (MMU) in Malaysia. Abed has received his B.Sc in industrial engineering field in 2005 from Iran Azad University. His research interest includes control project, CRM, consumer behavior, quality management, Marketing, e-commerce and image processing

Sahar Sabbaghi Mahmouei is currently doing her Master degree in Institute of Advanced Technology and Research (ITMA), Universiti Putra Malaysia, UPM. Sahar has received her B.Sc in software computer engineering field in 2006 from Iran Azad University. Her research interest includes image processing, machine vision, artificial Intelligence and e-commerce.

# A Face Replacement System Based on Face Pose Estimation

Kuo-Yu Chiu
Department of Electrical
Engineering
National Chiao-Tung University,
Hsinchu, Taiwan(R.O.C)
E-mail:Alvin_cgr@hotmail.com*

Shih-Che Chien
Department of Electrical
Engineering
National Chiao-Tung University,
Hsinchu, Taiwan(R.O.C)

Sheng-Fuu Lin
Department of Electrical
Engineering
National Chiao-Tung University,
Hsinchu, Taiwan(R.O.C)
E-mail:sflin@mail.nctu.edu.tw*

*Abstract*—**Face replacement system plays an important role in the entertainment industries. However, most of these systems nowadays are assisted by hand and specific tools. In this paper, a new face replacement system for automatically replacing a face with image processing technique is described. The system is divided into two main parts: facial feature extraction and face pose estimation. In the first part, the face region is determined and the facial features are extracted and located. Eyes, mouth, and chin curve are extracted by their statistical and geometrical properties. These facial features are used as the information for the second part. A neural network is adopted here to classify the face pose according to the feature vectors which are obtained from the different ratio of facial features. From the experiments and some comparisons, they show that this system works better while dealing with different pose, especially for non-frontal face pose.**

*Keywords- Facial feature• Face replacement• Neural network• Support vector machine (SVM)*

## I. INTRODUCTION

For entertainment and special effects industries, the ability of automatically replacing a face in a video sequence with that of another person has huge implications. For example, consider a stunt double in full-view of the camera performs a dangerous routine, and the stunt double's face could be automatically replaced latter with that of the desired actor for each instance by a post-processing. While few of the recent films have achieved good results when performing face replacement on the stunt doubles, there are still some limits, such like the illumination conditions in the environment should be controlled and the stunt double has to wear a special custom-fit mask with reflective markers for tracking [1].

In order to accurately replace a face in a photograph or a frame of video, we separate the system into two main parts. The first part is facial feature extraction and the second part is face pose estimation. Generally, the common approach of face region detection is to detect the face region by using the characteristic of the skin color. After locating the face region, the facial features can be obtained and determined by the geometric relation and statistical information. For example, the most common pre-processing method is to detect skin regions by a built skin tone model. R.L. Hsu et al. [2] proposed a face detection method based on a novel light compensation

technique and a nonlinear color trans-formation. Besides, there are still many color models used for the human skin-color [3]-[5]. For example, H. K. Jee et al. [6] used the color, edge, and binary information to detect eye pair from input image with support vector machine. Classifier boost methods are used to detect face region in paper [7]-[8]. However, neural network-based approaches required a large number of face and non-face training examples [9]-[11]. C. Garcia et al. [12] presented a novel face detection based on a convolutional neural architecture, which synthesized simple problem-specific feature extractors. There are also several algorithms for facial feature extraction. C. H. Lin [13] located facial feature points based on deformable templates algorithm. C. Lin [14] used the geometric triangle relation of the eyes and the mouth to locate the face position. Yokoyama [15] synthesized the color and edge information to locate facial feature.

The second part of face replacement system is face pose estimation. It is assumed that the viewpoint is on a fixed location and the face has an unknown pose that needs to be determined by one or more images of the human head. Previous face pose estimation algorithms can be roughly classified into two main categories: window-based approaches [16]-[19] and feature-based approaches [20]-[23]. Window-based approaches extract face block from the input image and analyze the whole block by statistical algorithms. Among window-based approaches, multi-class classification method divides the whole head pose parameter space into several intervals and determine head pose [16]-[17]. For example, Y. M. Li et al. [18] used the technique of support vector regression to estimate the head pose, which could provide crucial information and improve the accuracy of face recognition. L. Zhao et al. [19] trained two neural networks to approximate the functions that map a head from an image to its orientation. Windowed-based approaches have the advantage that they can simplify the face pose estimation problem. However, the face pose is generally coupled with many factors, such as the difference of illumination, skin color, and so on. Therefore, the learning methods listed above require large number of training samples.

On the other hand, the feature-based approaches extract facial features from human face by making use of the 3D structure of human face. These approaches are used to build 3D models for human faces and to match the facial features, such

---

*Corresponding author

as face contour and the facial components of the 3D face model with their projection on the 2D image. Y. Hu et al. [20] combined facial appearance asymmetry and 3D geometry to estimate face poses. Besides, some sensors are used to improve feature location. For instance, D. Colbry et al. [21] detected key anchor points with 3D face scanner data. These anchor points are used to estimate the pose and then to match the test image to 3D face model. Depth and brightness constraints can be used to locate features and to determine the face pose in some researches [22]-[23].

This paper is organized as follows. The face region detection and facial feature extraction system are introduced in Section 2. Section 3 describes the face pose estimation system. The face replacement system will be exhibited in Section 4. Section 5 shows the experimental results and comparisons. Finally, the conclusions and the future works are drawn in Section 6.

## II. FACIAL FEATURE EXTRACTION

Facial feature extraction plays an important role in face recognition, facial expression recognition, and face pose estimation. A facial feature extraction system contains two major parts: face region detection and facial feature extraction. According to the skin color model, the candidate face regions can be detected first. Then, the facial features can be extracted by their geometric and statistic properties from the face region. In this section, face region detection and facial feature extraction will be described.

### A. Face Region Detection

The first step of the proposed face replacement system is to detect and to track the target face in an image. A skin color model is used here to extract the skin color region which may be a candidate face region. The skin color model is built in *YCbCr* color space [24]. This color space is attractive for skin color modeling because it can separate chrominance from luminance. Hence, an input image is first transformed from RGB color space to *YCbCr* color space. Then the skin-color pixels are obtained by applying threshold values which are obtained from training data. After the skin color region is extracted, the morphological operator and 4-connectivity are then adopted to enhance the possible face region. The larger connected region of skin-color pixels are considered as the face region candidate and the real face region is determined by eye detection. Skin color region with eyes is defined as the face region. SVM classifier [25] is used here to detect eyes. Three sets of eye data are used for training. Eye images with frontal pose (set A) or profile pose (set B) are trained as the positive patterns. For negative patterns, non-eye images (set C) such as nose, lips, and ears are included for eye detection. All the training sets for eye detection are shown in Fig.1.
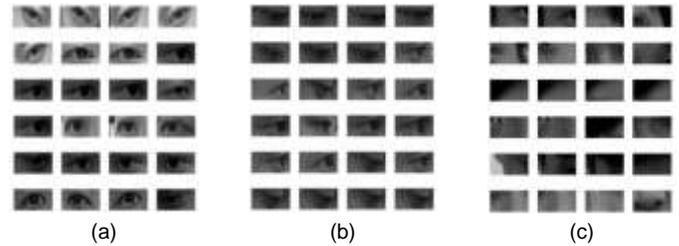


Figure 1. Training data of SVM. (a) Eye images of frontal pose. (b) Eye images of half-profile or profile pose. (c) Non-eye images.

Hence, for an input image as Fig.2a, the skin color region, which may be a face candidate, can be extracted after applying skin color model, as Fig.2b. Morphological operator and 4-connectivity is used then to eliminate noise and enhance the region shape and boundary, as Fig.2c. The skin color region is defined as a face when the eyes can be found by using SVM-based eye detection, as Fig.2d.
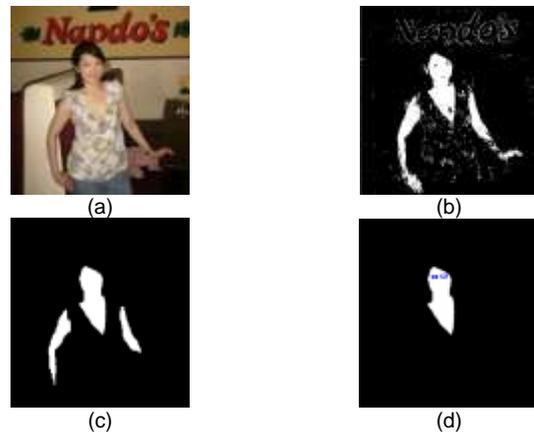


Figure 2. Face region detection. (a) Original image. (b) Binary image after applying the skin color model. (c) Possible face candidate regions after applying morphological operator and 4-connectivity. (d) The remaining face region after applying eye detection.

### B. Facial Feature Extraction

After the eyes are located and the face region is determined, the other facial features, such like lips region and the chin curve, can be easily extracted according to their geometrical relationship. In this section, the locating of right tip and left tip of lips region, the construction of chin curve, and the hair region segmentation are described.

To extract the lips region, the property that the lips region is on the lower part of the face and the color of lips region is different from the skin color is considered. Since the lips region is redder than the skin color region, a red to green function $RG(x,y)$ is employed to enhance the difference of lips color and skin color [26]. From the experimental results, the function $RG(x,y)$ is defined as follows:

$$RG(x, y) = \begin{cases} \dfrac{R(x, y)}{G(x, y)+1}, & \text{if } R(x, y) + G(x, y) + B(x, y) > 50, \\ 0, & \text{if } R(x, y) + G(x, y) + B(x, y) \le 50. \end{cases} \quad (1)$$

The $RG(x,y)$ has higher value when the value of red channel is larger then the value of green channel, which is probably a pixel of lips region. The possible lips region with higher red value is shown in binary image as Fig.3b. Besides, the edge information is also taken into account here to improve the lips region locating. In the *YCbCr* color space, the Sobel operator is employed to find the horizontal edge in luminance (*Y*) channel. The edge information is shown in Fig.3c. Using the union of redder region and edge information, the left and right tip points of lips region can be determined. The results of left and right tip points of lips region locating is shown in Fig.3d.



(a)  (b)  (c)  (d)

Figure 3. Lips region locating. (a) Original image. (b) Binary image of function $RG(x,y)$. (c) Horizontal edge by using Sobel operator. (d) The left and right tip points of lips region.

The next facial feature which is going to be extracted is the chin curve. There are two advantages to extract the chin curve: one is to separate the head from the neck and the other is to estimate the face pose. Since the chin curve holds strong edge information, a face block image is transformed into gray value image, as Fig.4a, and then the entropy function is applied to measure the edge information. Large entropy value contains more edge information, as shown in Fig.4b. The equation of entropy is defined as follows:

$$E = -\sum_{M,N} P(I_{m,n}) \log P(I_{m,n}) \qquad (2)$$

where $I_{m,n}$ represents the gray level value of point $(m,n)$ and $P(I_{m,n})$ is the probability density function of $I_{m,n}$. Using the lips position found before and the face block information, the searching region for the chin curve can be pre-defined. Five feature points, $x_1$, $x_2$, $x_3$, $x_4$, and $x_5$, are used to represent the chin curve. These five feature points are the intersections of chin curve and horizontal or vertical extended line from lips. The feature points, $x_1$ and $x_5$, are the intersections of chin curve and horizontal extended line from left tip of lips and right tip of lips respectively. The feature points, $x_2$, $x_3$, and $x_4$, are the intersections of chin curve and vertical extended line from left tip, middle, and right tip of lips. These feature points are shown in Fig.4c. Since the chin curve may not be symmetric, two quadratic functions, defined as: $y = ax^2 + bx + c$, are adopted here to construct the chin curve. The features $x_1$, $x_2$, and $x_3$ are used to find out the left quadratic function $f_{Lc}$ and the features $x_3$, $x_4$, and $x_5$ are used to find out the right quadratic function $f_{Rc}$ by using lease square method. The result of chin curve fitting is shown in Fig.4d.
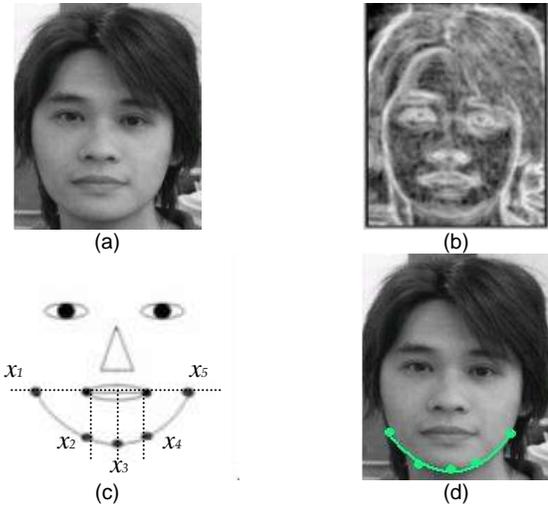


(a)  (b)

(c)  (d)

Figure 4. Chin curve construction. (a) Input gray level image. (b) The entropy of input gray level image. (c) Five feature points, $x_1$, $x_2$, $x_3$, $x_4$, and $x_5$, represent the most left point to the most right point respectively. (d) The function of chin curve fitting.

Hence, for an input image as Fig.5a, the skin color region can be found first as Fig.5b. Using the information of curve fitting function, the face region can be separated from the neck, as shown in Fig.5c.
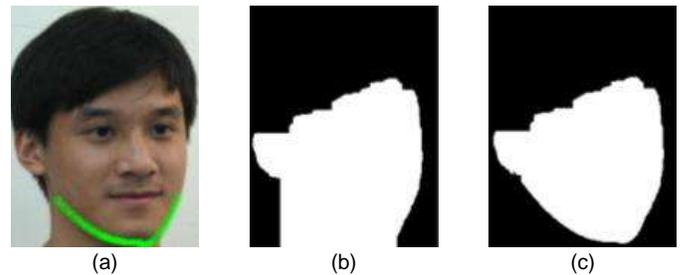


(a)  (b)  (c)

Figure 5. Face region segmentation. (a) Input image with curve fitting function. (b) Skin color region. (c) Face region only by using curve fitting information.

After the face region is found, the hair region can be defined easily. It is known that the hair region is above the face region. Hence, if an appropriate block above the face region is chosen and the skin color region is neglected, the remaining pixels, as Fig.6a, can be used as the seeds for seed region growing (SRG) algorithm. The hair region then can be extracted. The hair region extraction result is shown in Fig.6b.



(a)  (b)

Figure 6. Hair region extraction. (a) The remaining pixels are used as the seeds for SRG after the skin color region is neglected. (b) The result of hair region extraction.

### III. POSE ESTIMATION

In this section, how to estimate the face pose is detailed. All the different face poses are described with three angle parameters, namely the yaw angle $\alpha$, the tilt angle $\beta$, and the roll angle $\gamma$, as shown in Fig.7.
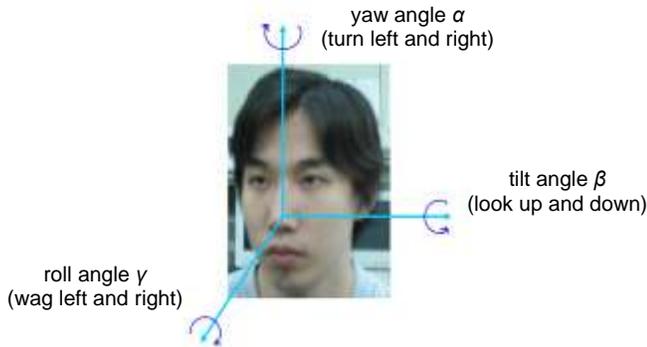


Figure 7. Parameters of face pose estimation.

Since the profile pose is much different from the frontal pose, two different methods are proposed here. All the different poses are roughly divided into two classes according to the number of eyes extracted in SVM-based eye detection system. When there are two eyes extracted, the face pose belongs to class A which is more frontal. Otherwise, if only one eye is extracted, then the face pose belongs to class B which is more profile. The examples of class A and class B are shown in Fig.8a and Fig.8b respectively.



Figure 8. Two kinds of face pose. (a) Frontal face pose with two eyes extracted. (b) Profile face pose with only one eye extracted.

#### A. Pose Angle Estimation of Class A

For an input image of class A, it will be normalized and rotated first so that the line crossing two eyes is horizontal. In other words, the roll angle $\gamma$ of the input face should be found out first. The roll angle $\gamma$ is defined as the elevation or depression angle from left eye. Using the relative vertical and horizontal distance of the two eyes, the roll angle $\gamma$ can be obtained. Set $x_6$ and $x_7$ as the center of left eye and right eye respectively as shown in Fig. 9a, the roll angle $\gamma$ is defined by:

$$\gamma = \tan^{-1}\left(\frac{y_{x_7} - y_{x_6}}{x_{x_7} - x_{x_6}}\right) \tag{3}$$

where $x$ and $y$ represent the x-coordinate and y-coordinate respectively. Using the information of roll angle $\gamma$, the image can be rotated to horizontal as Fig. 9b. For an input image Fig. 9c, the normalization result is shown in Fig. 9d.
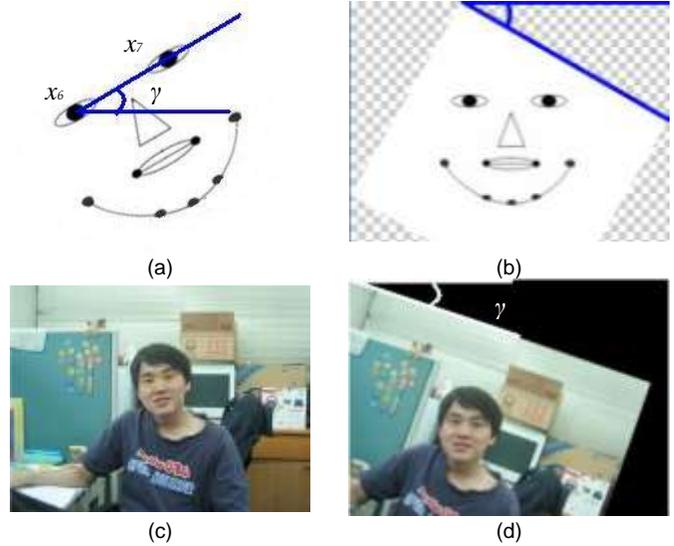


Figure 9. The roll angle $\gamma$. (a) The definition of $x_6$ and $x_7$. (b) The rotated image with horizontal eyes. (c) Input image. (d) Normalization result of input image (c).

After retrieving the roll angle information, the face can be normalized to horizontal. Five scalars, $v_1$, $v_2$, $v_3$, $v_4$, and $v_5$, are used as the input of neural network to estimate the face pose in class A. The first scalar $v_1$ is defined as:

$$v_1 = \frac{L_1}{L_2} \tag{4}$$

where $L_1$ is the horizontal distance between the left tip of lips and the constructed chin curve $f_c$ and $L_2$ is the distance between the right tip of lips and $f_c$. The scalar $v_1$ is relative to the yaw angle $\alpha$. It is close to 1 when the yaw angle $\alpha \approx 90°$, as Fig. 10a. When the face turns to right as Fig. 10b, $L_1$ is smaller than $L_2$ and the scalar $v_1$ is smaller than 1. Contrarily, when the face turns to left as Fig. 10c, $L_1$ is larger than $L_2$ and the scalar $v_1$ is larger than 1.
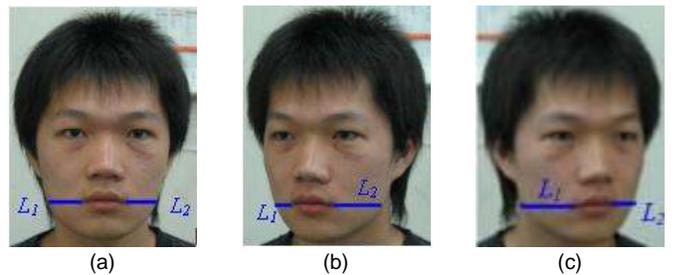


Figure 10. The relationship between scalar $v_1$ and the yaw angle $\alpha$. (a) Scalar $v_1$ is close to 1 when the face is frontal. (b) Scalar $v_1$ is smaller than 1 when the face turns to right. (c) Scalar $v_1$ is larger than 1 when the face turns to left.

The second scalar $v_2$ is defined as the ratio of $L_3$ and $L_4$:

$$v_2 = \frac{L_3}{L_4} \tag{5}$$

where $L_3$ is the vertical distance between the middle point of two eyes, defined as $x_8$, and the constructed chin curve, and $L_4$

is the vertical distance between the center of the lips and $x_8$, as Fig. 11a. The scalar $v_2$ is relative to the tilt angle $\beta$ as Fig. 11b. The third scalar $v_3$ is defined as:

$$v_3 = \frac{L_3}{L_5} \qquad (6)$$

where $L_5$ represents the distance of $x_6$ and $x_7$. The scalar $v_3$ is relative to the tilt angle $\beta$, as Fig. 11c, and the yaw angle $\alpha$, as Fig. 11d, simultaneously.
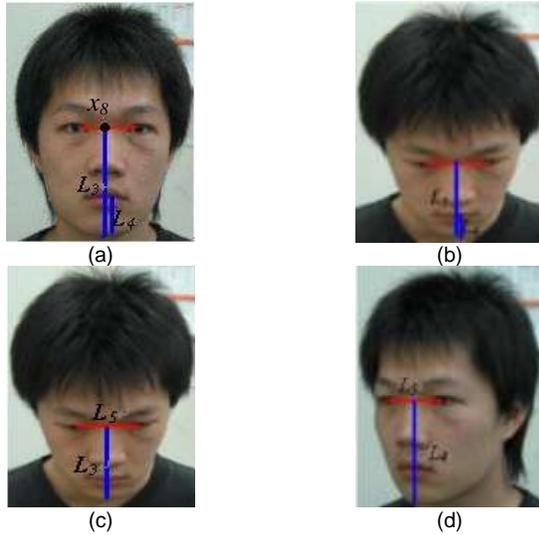


(a)            (b)

(c)            (d)

Figure 11. The relationship between scalars, $v_2$ and $v_3$, and pose parameter, tilt angle $\beta$ and yaw angle $\alpha$. (a) The definitions of $x_8$, $L_3$ and $L_4$. (b) The relationship between scalar $v_2$ and tilt angle $\beta$. (c) The relationship between scalar $v_3$ and tilt angle $\beta$. (d) The relationship between scalar $v_3$ and yaw angle $\alpha$.

Before defining the last two scalars, another two parameters, $L_6$ and $L_7$, are defined first. Connecting the feature point $x_3$ of the chin curve and two tip points of the lips, the extended lines will intersect the extended line crossing $x_6$ and $x_7$ with two intersections. These two intersections are defined as $x_9$ and $x_{10}$ from left to right respectively as Fig. 12a. Parameter $L_6$ is then defined as the distance between $x_6$ and $x_9$, and $L_7$ is the distance between $x_7$ and $x_{10}$. The definitions of parameters $L_6$ and $L_7$ are shown in Fig. 12b. Then, the forth scalars $v_4$ is defined as:

$$v_4 = \frac{L_6 \cdot L_7}{L_5} \qquad (7)$$

and the last scalars $v_5$ is defined as:

$$v_5 = \frac{L_6}{L_7}. \qquad (8)$$

Scalar $v_4$ is relative to tilt angle $\beta$, as shown in Fig. 12c, and the scalar $v_5$ is relative to yaw angle $\alpha$, as shown in Fig. 12d.
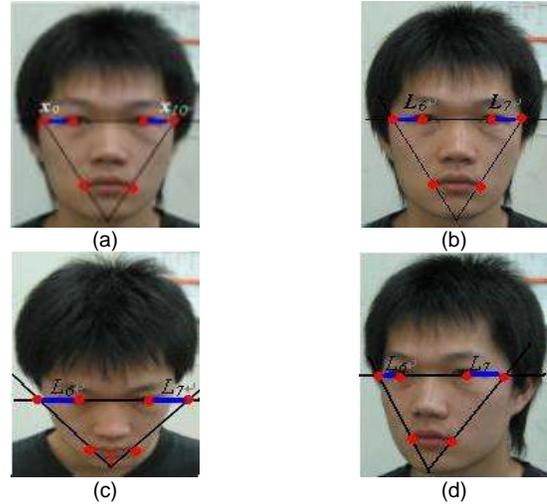


(a)            (b)

(c)            (d)

Figure 12. The relationship between scalars, $v_4$ and $v_5$, and pose parameter, tilt angle $\beta$ and yaw angle $\alpha$. (a) The definitions of $x_9$ and $x_{10}$. (b) The definitions of $L_6$ and $L_7$. (c) The relationship between scalar $v_4$ and tilt angle $\beta$. (d) The relationship between scalar $v_5$ and yaw angle $\alpha$.

### B. Pose Angle Estimation of Class B

The face is classified to class B if only one eye can be found when applying eye detection. For the face in class B, there are also five scalars used as the input of neural network to estimate the face pose. Feature points $x_{11}$ and $x_{12}$ represent the intersection points of face edge and the horizontal extended line crossing the eye and the lips respectively. Feature point $x_{13}$ the tip point of chin curve which is found with the largest curvature and feature point $x_{14}$ is the only extracted eye center. With these four feature points which are shown in Fig. 13a, the first scalar $v'_1$ is defined as:

$$v'_1 = \frac{L_9}{L_8} \qquad (9)$$

where $L_8$ is the distance between $x_{14}$ and face edge $x_{11}$ and $L_9$ is the distance between $x_{14}$ and the middle point of the lips. These two parameters are shown in Fig. 13b. The first scalar $v'_1$ is relative to the yaw angle $\alpha$ as Fig. 13c.
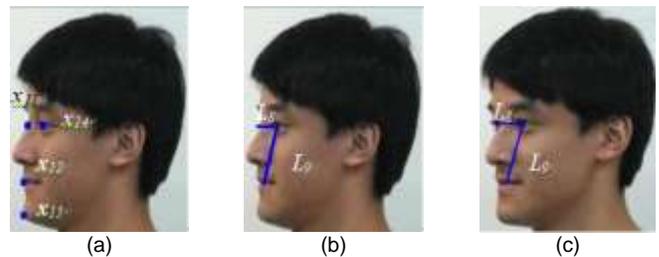


(a)        (b)        (c)

Figure 13. The relationship between scalar $v'_1$ and yaw angle $\alpha$. (a) The definition of feature points $x_{11}$, $x_{12}$, $x_{13}$, and $x_{14}$. (b) The definition of parameters $L_8$ and $L_9$. (c) The scalar $v'_1$ is relative to the yaw angle $\alpha$.

The scalar $v'_2$ is the slope of the line crossing $x_{12}$ and $x_{13}$ as shown in Fig. 14a and it is defined by:

$$v'_2 = m_{\overline{x_{12} \cdot x_{13}}} \qquad (10)$$

where $m$ represents slop. The scalar $v'_2$ is relative to the tilt angle $\beta$ as Fig. 14b. The scalar $v'_3$ is the angle $\theta$ which is defined by:

$$\theta = \angle x_{11} x_{14} x_{12}, \quad 0° < \theta < 90° \qquad (11)$$

and it is shown in Fig. 14c. The scalar $v'_3$ is relative to the tilt angle $\beta$ as Fig. 14d and the yaw angle $\alpha$ as Fig. 14e, simultaneously.
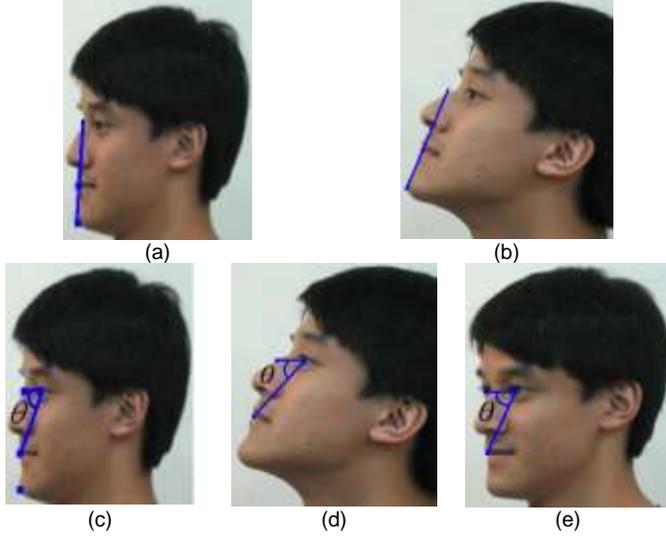


Figure 14. The relationship between scalars, $v'_2$ and $v'_3$, and pose parameter, tilt angle $\beta$ and yaw angle $\alpha$. (a) The line crossing $x_{12}$ and $x_{13}$. (b) The scalar $v'_2$ is relative to the tilt angle $\beta$. (c) The definition of angle $\theta$. (d) The scalar $v'_3$ is relative to the tilt angle $\beta$. (e) The scalar $v'_3$ is relative to the yaw angle $\alpha$.

Connecting $x_{14}$ with middle point and right tip point of lips, the extended line will intersect the horizontal line passing $x_8$ with two intersections. $L_{10}$ is defined as the distance between these two intersections as Fig. 15a. Then the scalar $v'_4$ is defined as:

$$v'_4 = L_{10} \cdot L_8 \qquad (12)$$

and the scalar $v'_5$ is defined as:

$$v'_5 = \frac{L_{10}}{L_9}. \qquad (13)$$

The scalar $v'_4$ and $v'_5$ are relative to the tilt angle $\beta$, Fig. 15b, and the yaw angle $\alpha$, Fig. 15c, simultaneously.
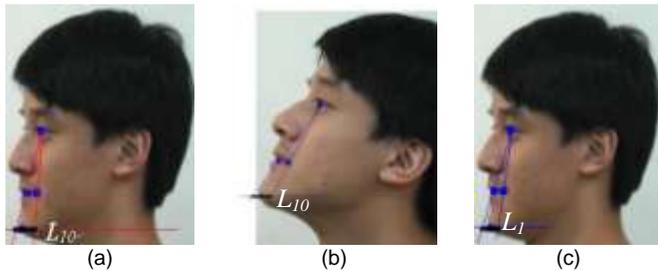


Figure 15. The relationship between scalars and pose parameter (a) The definition of $L_{10}$. (b) The scalar $v'_4$ is relative to the tilt angle $\beta$. (c) The scalar $v'_5$ is relative to the yaw angle $\alpha$.

## IV. FACE REPLACEMENT

In this section, the procedure of face replacement is detailed. For an input target face, the face pose is estimated first and then the face with similar face pose is chosen from the database as the source face to replace the target face. However, there are some problems when replacing the face, such like the mismatch of face size, face position, face pose angle, and skin color. Hence, image warping and shifting are adopted first to adjust the source face so that it is much similar as the target face. Color consistency and image blending are used later to reduce the discontinuousness due to the replacement. All the details are described below.

### A. Image Warping and Shifting

After the face pose angle of target face is determined and the face region of source face is segmented, the target face is going to be replaced by the source face. However, the resolution, face size, and face pose angle may not be exactly the same. Hence, image warping is adopted here to deal with this problem.

Image warping is applied according to features matching. It is a spatial transformation that includes shifting, scaling, and rotating. In this paper, an affine matrix with bilinear interpolation is used to achieve image warping. The affine transformation matrix is defined by:

$$\begin{bmatrix} X' \\ Y' \\ 1 \end{bmatrix} = \begin{bmatrix} m_1 & m_2 & m_3 \\ m_4 & m_5 & m_6 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} \qquad (14)$$

where $(X',Y')$ is the feature point coordinate of target face, $(X,Y)$ is the feature point coordinate of source face, and $m_1,\ldots,m_6$ are parameters. For faces in class A, six feature points, two eyes ($x_6$ and $x_7$), the center of lips, and feature points $x_1$, $x_3$, and $x_5$ of chin curve, are used to solve the matrix by the least square method as Fig. 16a, while four feature points, $x_{11}$, $x_{12}$, $x_{13}$, and $x_{14}$, are used for faces in class B as Fig. 16b.
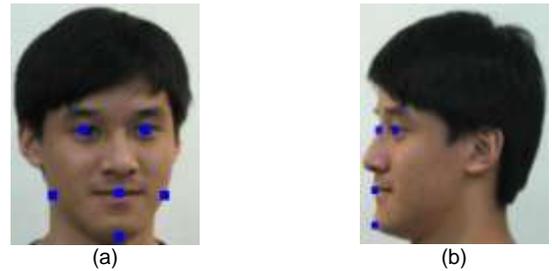


Figure 16. Feature points for image warping. (a) Six feature points are used for class A. (b) Four feature points are used for class B.

After the source face is warped, a suitable position is going to be found to replace the target face. A better face replacement is achieved when more face and hair regions are matched in both source face and target face. The source face is first pasted on so that the coordinates of pasted feature point are the same for source face and target face. The pasted feature point is chosen as the middle point of chin curve, $x_3$, for class A and tip points of chin curve, $x_{13}$, for class B. Later, the pasted source

face is shifted around the pasted feature point to a best position with most matching points. A matching degree function $M(x,y)$ for a pasting point $(x,y)$ is used to evaluate the degree of matching, which is defined as:

$$M(x,y) = \sum_{(i,j)\in I} [h(F_s(i,j), F_t(i,j)) + h(H_s(i,j), H_t(i,j))] \qquad (15)$$

where $F_s(i,j)$ and $F_t(i,j)$ are binary face images which have value 1 only for face region pixel in source and target images respectively, $H_s(i,j)$ and $H_t(i,j)$ are binary hair images which have value 1 only for hair region pixel in source and target images, and $I$ is the region of interest which is larger than the pasted region. The function $h(a,b)$ in equation (15) is defined by:

$$h(a,b) = \begin{cases} +1, & \text{if } a = b = 1, \\ 0, & \text{if } a = b = 0, \\ -1, & \text{if } a \neq b. \end{cases} \qquad (16)$$

For each point near the pasted feature point, the matching degree can be calculated. The point with highest matching degree will be chosen as the best position to paste the source face on. For example, the face region (white) and hair region (red) for source face image and target face image are shown in Fig. 17a and 17b respectively. When the target face is randomly pasted by the source face as Fig. 17c, there are more "-1", denoted as the red region, and less "+1", denoted as the white region. This means that the matching degree is low. After calculating all the matching degree of nearby points, the best pasting point with most "+1" and least "-1" can be found, as Fig. 17d.
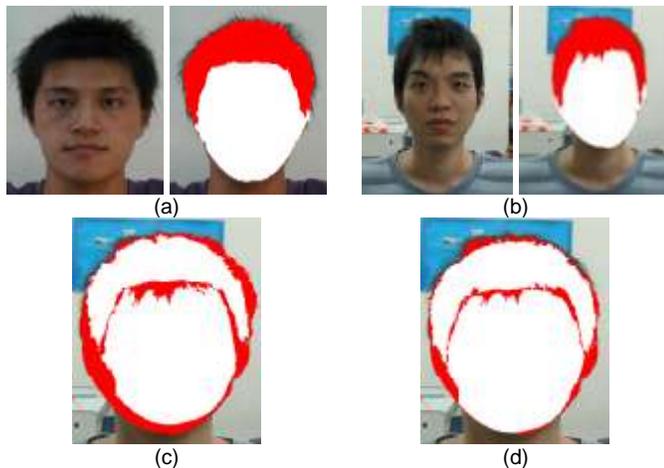


Figure 17. Image shifting according to matching degree. (a) Source face image (b) Target face image (c) Face replacement with lower matching degree. (d) Face replacement with highest matching degree.

## B. Color Consistency and Image Blending

Because of the difference of luminance and human races, the skin color of target face may not be similar to the source face. To solve the problem, skin color consistency is adopted here. The histogram of both source face and target face are analyzed first and the mean of skin color of target face is shifted to the same value as the mean of source face. For example, the source face as Fig. 18a is darker than the target face as Fig. 18b. If the face replacement is applied without adopting skin color consistency, the skin color of face region and necks region of the result is different, as shown in Fig. 18c. To avoid this situation, the mean of histogram of the target face is shifted to the same value as the source face, as Fig. 18d. Then, the skin color of the face region and necks region will be similar after replacement. The result of face replacement with skin color consistency is shown in Fig. 18e.
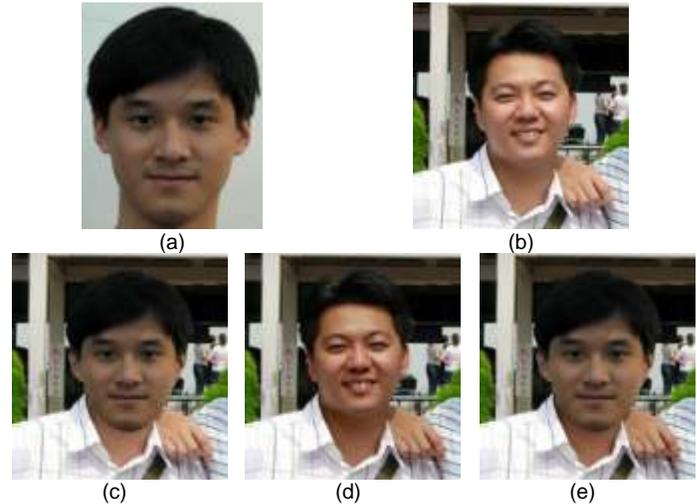


Figure 18. Skin color consistency. (a) Source Face. (b) Target Face. (c) Face replacement without skin color consistency. (d) The mean of histogram of target face is shifted to the same value as the source face. (e) Face replacement with skin color consistency.

Finally, an image blending method is applied to deal with the boundary problem. Though the source skin color is changed so that it is consistent with target face, there is still boundary problem when the source face replaces the target face because of discontinuousness. The objective of image blending is to smooth boundary by using interpolation. The hyperbolic tangent is used as the weight function:

$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}. \qquad (17)$$

The horizontal interpolation is described as:

$$I(x,Y) = w_{h(x)}L(x,Y) + (1 - w_h(x))R(x,Y) \qquad (18)$$

and the vertical interpolation is described as:

$$I(X,y) = w_{v(y)}D(X,y) + (1 - w_h(y))U(X,y) \qquad (19)$$

where $I(x,y)$ is the boundary point; $L(x,Y)$, $R(x,Y)$, $U(X,y)$, and $D(X,y)$ represent the left, right, up, and down image respectively. The result of image blending is exhibited in Fig. 19. These images are not applied with color consistency, so the boundary is sharper because of face replacement. However, it can be seen that the image in Fig. 19b with image blending has smoother boundary than the one in Fig. 19a without image blending.

Figure 19. Image blending. (a) Without Image blending. (b) With image blending.

## V.    EXPERIMENTAL RESULTS

In this section, the results of face pose estimation and face replacement will be shown. Some analyses and comparisons will also be made in this section.

### A.    Face Pose Estimation

To verify the accuracy of the face pose estimation system, face images under various poses are collected and tested. In the database, the face poses are divided into 21 classes according to different yaw angle $\alpha$ and tilt angle $\beta$. The face pose is estimated by multi-class classification based on neural network. The yaw angle is divided into 7 intervals and the tilt angle is divided into 3 intervals, as shown in Fig. 20a and Fig. 20b respectively. Because the face poses of turning right and turning left are the same by applying a reflection matrix, only left profile face poses are considered.
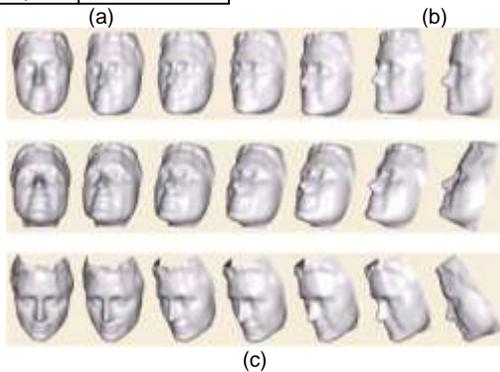


Figure 20. Intervals of different face angles and different face poses. (a) 7 intervals of different $\alpha$. Positive value represents turning left and the magnitude represents the turning degree. (b) 3 intervals of different $\beta$. Positive value represents looking up and the magnitude represents the turning degree. (c) 21 Different face poses.

There are 1680 images from 80 people in the database for training and another 1680 images from other 80 people are used for testing. The accuracy rate of pose estimation, $R_{PE}$, is defined

by:

$$R_{PE} = \frac{\text{Total Photos - Failure estimation}}{\text{Total Photos}}. \tag{20}$$

Therefore, the accuracy rate of face pose estimation can be calculated. There are some other face pose estimation methods, such as pose estimation method based on Support Vector Regression (SVR) using Principal Component Analysis (PCA) [18] and Neural Network (NN) based approach [19]. The comparisons of face pose estimation methods are shown in Fig. 21.

| | The proposed Method | PCA+SVR | Neural Networks |
|---|---|---|---|
| $[\alpha_1, \alpha_4]$ (pure background) | **88.75 %** | **87.5 %** | 86.45 % |
| $[\alpha_5, \alpha_7]$ (pure background) | **89.16 %** | **86.17 %** | 85.67 % |
| $[\alpha_1, \alpha_4]$ (complex background) | **86.94 %** | **85.69 %** | 85.56 % |
| $[\alpha_5, \alpha_7]$ (complex background) | 87.5 % | 82 % | 81.67 % |

Figure 21.  Accuracy comparisons of face pose estimation.

### B.    Face Replacement

In this section, the results of face replacement are shown. Various conditions are considered to test the robustness of this automatic face replacement system, such like wearing glasses, different resolution, different luminance, different skin color, different yaw angle, different roll angle, and different tilt angle. It can be seen from the results that this system performs well while dealing with these conditions.

In Fig. 22, wearing glasses or not is discussed. The target face with glasses, as Fig. 22b, is replaced by the source face without glasses, as Fig. 22a. Since the target face region is replaced by the entire source face, wearing glasses or not will not affect the results.

When the face size and luminance of target face and source face are different, the face size and the skin color will be adjusted. The source face in Fig. 23a will be resized to fit the target face by the affine matrix according to the facial feature matching. Color consistency method is also applied in this case. It can adjust the skin color of the target face in Fig. 23b so that the skin color of target face is similar to the source face and the result would be better after replacement. From the result in Fig. 23c, it can be seen that the skin color of target face are adjusted and shifted to the similar value as source face, especially for the neck region. Since the skin color and the face size are adjusted, the replacement result is more nature.

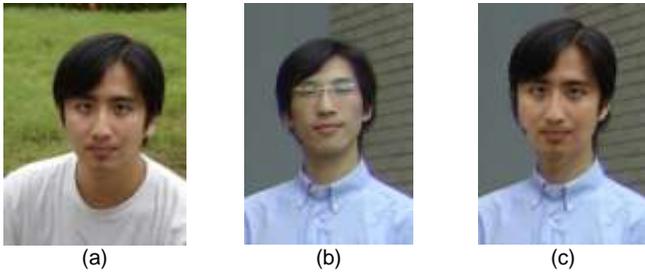(a)                    (b)                    (c)

Figure 22. Face replacement result when considering the glasses. (a) Source face image without glasses. (b) Target face image with glasses. (c) The replacement result.
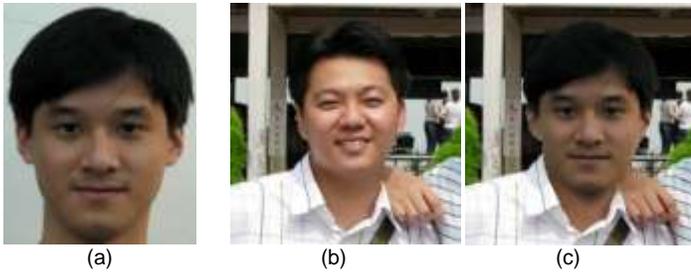


(a)                    (b)                    (c)

Figure 23. Face replacement result when considering different face size and illumance. (a) Source face image with thiner face region and darker skin color. (b) Target face image with wider face region and brighter skin color. (c) The replacement result.

While dealing with the profile pose, such as the face with 90 degrees of yaw angle as Fig. 24a, a face image with similar face pose is chosen from the database to replace the target face. The result would be poor if the replacement is done by only adopting an affine matrix without a proper face pose. From the result shown in Fig. 24b, it can be seen that this system performs well while dealing with profile pose, even if the face has 90 degrees of yaw angle.



(a)                              (b)

Figure 24. Face replacement result when considering the profile face. (a) Target face image. (b) The result of face replacement.

Like the profile pose, when dealing with the face with tilt angle such as Fig. 25b, a proper source face will be found from the database first. According to the face pose estimation system, a face with most similar face pose is chosen, as Fig. 25a. After applying a reflection matrix, the face pose of source face is almost the same as the target face. With color consistency method, the replacement can be done even though there are tilt angles and yaw angles at the same time for a target face. The face replacement result is shown in Fig. 25c.

When considering the target face with a roll angle, such as Fig. 26b, the roll angle is calculated first according to two eyes. After the roll angle is found and a similar pose is chosen from

the database for the source face as Fig. 26a, a rigid transformation is adopted to rotate the source image such that the roll angle of source face is the same as the roll angle of target face. In Fig. 26c, it can be seen that the replacement is done with a rigid transformation.



(a)                    (b)                    (c)

Figure 25. Face replacement result when considering the face with tilt angle. (a) Source face image (b) Target face image with a tilt angle. (c) The replacement result.



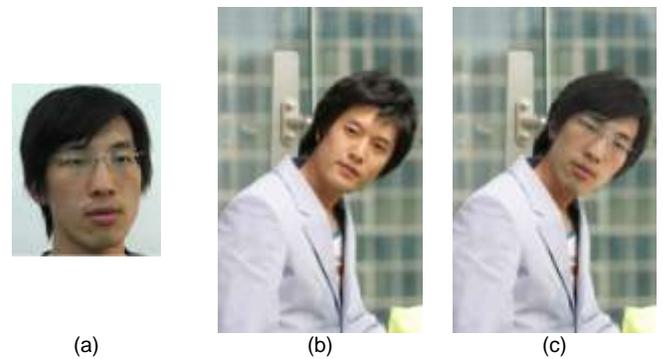(a)                    (b)                    (c)

Figure 26. Face replacement result when considering the face with roll angle. (a) Source face image. (b) Target face image with a roll angle. (c) The replacement result.

## VI.   CONCLUSIONS AND FUTURE WORK

Face replacement system plays an important role in the entertainment industries. In this paper, a face replacement system based on image processing and face pose estimation is described. Various conditions are considered when replacing the face, such as different yaw angle, roll angle, tilt angle, skin color, luminance, and face size. The experiment results show that this face replacement system has good performance while dealing the conditions listed before. In the future, facial expression would be a further challenge task to be considered. Along with the facial expression, the results of face replacement will be realer and the system will be much more powerful and useful in entertainment industries.

### REFERENCES

[1]   J.   Kleiser:   Kleiser-Walczak   on   the   one.   DOI: http://www.kwcc.com/works/ff/the_one.html

[2]  Hsu, R.L., Abdel-Mottaleb, M., Jain, A.K.: Face detection in color images. IEEE Trans. Pattern Analysis and Machine Intelligence. 24(5), 696-706 (2002)

[3]  Zhu, X., Yang, J., Waibel, A.: Segmenting hands of arbitrary color. In: Proceedings of 4th IEEE Conf. Automatic Face and Gesture Recognition, pp. 446-453, 2000

[4]  Jee, H.K., Lee, K., Pan, S.B.: Eye and face detection using SVM. In: Proceedings of IEEE Conf. Intelligent Sensors, Sensor Networks and Information Processing, pp. 577-580, 2004

[5]  Lin, Y.Y., Liu, T.L.: Robust face detection with multi-class boosting. In: Proceedings of IEEE Computer Society Conf. Computer Vision and Pattern Recognition, vol. 1, pp. 680-687, 2005

[6]  Huang, C., Al, H.Z., Wu, B., Lao, S.H.: Boosting nested cascade detector for multi-view face detection. In: Proceedings of 17th IEEE International Conf. Pattern Recognition, vol. 2, pp. 415-418, 2004

[7]  Bingulac S.P.: On the Compatibility of Adaptive Controllers. In: Proceedings of Fourth Ann. Allerton Conf. Circuits and Systems Theory, pp. 8-16, 1994

[8]  MacQueen J.: Some Methods for Classification Analysis of Multivariate Observations. In: Proceedings of Fifth Berkeley Symp. Math. Statistics and Probability, pp. 281-297, 1967

[9]  Fu, H.C., Lai, P.S., Lou, R.S., Pao, H.T.: Face detection and eye localization by neural network based color segmentation. In: Proceedings of IEEE Signal Processing Society Workshop on Neural Networks for Signal Processing X, vol. 2, pp. 507-516, 2000

[10]  Garcia, C., Delakis, M.: Convolutional face finder: a neural architecture for fast and robust face detection. IEEE Trans. Pattern Analysis and Machine Intelligence. **26**(11), 1408-1423 (2004)

[11]  Lin, C.H., Wu, J.L.: Automatic facial feature extraction by genetic algorithms. IEEE Trans. Image Processing. , **8**(6), 834-845 (1999)

[12]  Lin, C., Fan, K.C.: Human face detection using geometric triangle relationship. In: Proceedings of 15th IEEE Int. Conf. Pattern Recognition, vol. 2, pp. 941-944, Barcelona, Spain 2000

[13]  Yokoyama, T., Wu, H., Yachida, M.: Automatic detection of facial feature points and contours. In: Proceedings of 5th IEEE Int. Workshop on Robot and Human Communication. pp. 335-340, Tsukuba, Japan 1996

[14]  Yang, Z.G., Ai, H.Z., Okamoto, T., Lao, S.H.: Multi-view face pose classification by tree-structured classifier. In: Proceedings of IEEE International Conf. Image Processing, vol. 2, pp. 358-361, 2005

[15]  Li, S.Z., Fu, Q.D., Scholkopf, B., Cheng, Y.M., Zhang, H.J.: Kernel machine based learning for multi-view face detection and pose estimation. In: Proceedings of 8th IEEE International Conf. Computer Vision, vol. 2, pp. 674-679, Vancouver, BC, Canada 2001

[16]  Li, Y.M., Gong, S.G., Liddell, H.: Support vector regression and classification based multi-view face detection and recognition. In: Proceedings of 4th IEEE International Conf. Automatic Face and Gesture Recognition, pp. 300-305, Grenble, France 2000

[17]  Zhao, L., Pingali, G., Carlbom, I.: Real-time head orientation estimation using neural networks. In: Proceedings of IEEE International Conf. Image Processing, vol.1, pp. 297-300, 2002

[18]  Hu, Y., Chen, L.B., Zhou, Y., Zhang, H.J.: Estimating face pose by facial asymmetry and geometry. In: Proceedings of 6th IEEE Conf. Automatic Face and Gesture Recognition, pp. 651-656, 2004

[19]  Colbry, D., Stockman, G., Jain, A.: Detection of anchor points for 3D face verification. IEEE Conf. Computer Vision and Pattern Recognition, 3, 118-125 (2005)

[20]  Covell, M., Rahini, A., Harville, M., Darrell, J.: Articulated pose estimation using brightness- and depth-constancy constraints. In: Proceedings of IEEE Conf. Computer Vision and Pattern Recognition, vol. 2, pp. 438-445, 2000

[21]  Harville, M., Rahimi, A., Darrell, T., Gordon, G., Woodfill, J.: 3D pose tracking with linear depth and brightness constraints. In: Proceedings of 7th IEEE International Conf. Computer Vision, vol. 1, pp. 206-213, Kerkyra, Greece 1999

[22]  Chai, D., Ngan, K.N.: Face segmentation using skin-color map in videophone applications. IEEE Trans. on circuits and systems for video technology, 9(4), 551-564 (1999)

[23]  Vapnik, V.: The Nature of Statistical Learning Theory. New York: Springer, 1995

[24]  Eveno, N., Caplier, A. Coulon, P.Y.: A new color transformation for lips segmentation. In: Proceedings of 4th IEEE Workshop on Multimedia Signal Processing, pp. 3-8, Cannes, France 2001

[25]  Nugroho, H., Takahashi, S., Ooi, Y., Ozawa, S.: Detecting human face from monocular image sequences by genetic algorithms. IEEE Int. Conf. Acoustics, Speech, and Signal Processing, 4, 2533-2536 (1997)

AUTHORS PROFILE

Sheng-Fuu Lin (S'84–M'88) was born in Tainan, R.O.C., in 1954. He received the B.S. and M.S. degrees in mathematics from National Taiwan Normal University in 1976 and 1979, respectively, the M.S. degree in computer science from the University of Maryland, College Park, in 1985, and the Ph.D. degree in electrical engineering from the University of Illinois, Champaign, in 1988. Since 1988, he has been on the faculty of the Department of Electrical and Control Engineering at National Chiao Tung University, Hsinchu, Taiwan, where he is currently a Professor. His research interests include image processing, image recognition, fuzzy theory, automatic target recognition, and scheduling.

Shih-Che Chien was born in Chiayi, R.O.C., in 1978. He received the B.E. degree in electronic engineering from the Nation Chung Cheng University, in 2002. He is currently pursuing the M.E. and Ph.D. degree in the Department of Electrical and Control Engineering, the National Chiao Tung University, Hsinchu, Taiwan. His current research interests include image processing, image recognition, fuzzy theory, 3D image processing, intelligent transportation system, and animation.

Kuo-Yu Chiu was born in Hsinchu, R.O.C., in 1981. He received the B.E. degree in electrical and control engineering from National Chiao Tung University, Hsinchu, Taiwan, R.O.C, in 2003. He is currently pursuing the Ph. D. degree in the Department of Electrical and Control Engineering, the National Chiao Tung University, Hsinchu, Taiwan. His current research interests include image processing, face recognition, face replacement, intelligent transportation system, and machine learning.

# A Comprehensive Analysis of Spoofing

P. Ramesh Babu

Dept of Information Technology
Rajamahendri Inst. of Engg &
Technology
Rajahmundry-533103, INDIA
E-mail:rameshbabu_kb@yahoo.co.in

D.Lalitha Bhaskari

Dept of C.S & S.E
AU College of Engineering (A)
Visakhapatnam-530003, INDIA
E-mail:lalithabhaskari@yahoo.co.in

CH.Satyanarayana

Dept of C.S.E
JNTUK College of Engineering
Kakinada – 533003, INDIA
E-mail: ce@jntukakinada.edu.in

***Abstract*--The main intention of writing this paper is to enable the students, computer users and novice researchers about spoofing attacks. Spoofing means impersonating another person or computer, usually by providing false information (E-mail name, URL or IP address). Spoofing can take on many forms in the computer world, all of which involve some type false representation of information. There are a variety of methods and types of spoofing. We would like to introduce and explain following spoofing attacks in this paper: IP, ARP, E-Mail, Web, and DNS spoofing. There are no legal or constructive uses for implementing spoofing of any type. Some of the outcomes might be sport, theft, vindication or some other malicious goal. The magnitude of these attacks can be very severe; can cost us millions of dollars. This Paper describes about various spoofing types and gives a small view on detection and prevention of spoofing attacks.** (Abstract)

***Keywords: Spoofing, Filtering, Attacks, Information, Trust***

## I. INTRODUCTION

Spoofing can take on many forms in the computer world, all of which involve some type false representation of information. There are a variety of methods and types of spoofing. We would like to introduce and explain following types in this paper:

- IP Spoofing
- ARP Spoofing
- E-Mail Spoofing
- Web Spoofing
- DNS Spoofing

There are no legal or constructive uses for implementing spoofing of any type. Some of the outcomes might be sport, theft, vindication or some other malicious goal. The gravity of these attacks can be very severe, can cost us millions of dollars and should not be overlooked by the Internet security community.
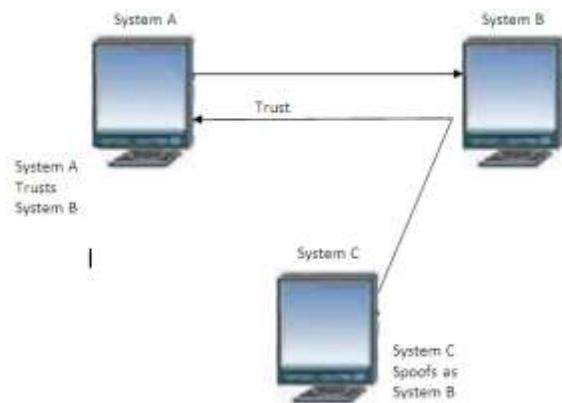
## II. IP SPOOFING

IP spoofing is used to gain unauthorized access to a computer. The attacker forwards packets to a computer with a source address indicating that the packet is coming from a trusted port or system. Attackers must go through some complicated steps to accomplish the task [1]. They must:

- Obtain a target.
- Obtain an IP address of a trusted machine.
- Disable communication of the trusted machine (e.g. SYN flooding).
- Sample a communication between the target and trusted hosts
- Guess the sequence numbers of the trusted machine.
- Modify the packet headers so that it appears that the packets are coming from the trusted host.
- Attempt connection to an address authenticated service or port.
- If successful, the attacker will plant some kind of backdoor access for future reference

System A impersonates system B by sending B's address instead of its own. The reason for doing this is that systems tend to function within groups of other ``trusted'' systems. This trust is implemented in a one-to-one fashion; system A trusts system B. IP spoofing occurs in the following manner: if system A trusts system B and system C spoofs system B, then system C can gain otherwise denied access to system A. This is all made possible by means of IP address authentication, and if the packets are coming from external sources- poorly configured routers [2].
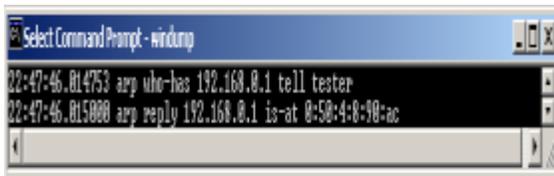
One of the major drawbacks with IP spoofing is that C never "sees" the responses from A. This is completely blind attack, much experience and knowledge of what to expect from the target's responses is needed to successfully carry out his attack.

Some of the most common ways to avoid this type of attack are to disable source-routed packets and to disable all external incoming packets with the same source address as a local host.

## III. ARP SPOOFING

ARP stands for Address Resolution Protocol. ARP is used to map IP addresses to hardware addresses [6]. A table, usually called the ARP cache, is used to maintain a correlation between each MAC address and its corresponding IP address. ARP provides the protocol rules for making this correlation and providing address resolution in both directions [3]. When an incoming packet sent to a host machine on a network arrives at a router, it asks the ARP program to find a MAC address that matches the IP address. The ARP program looks in the ARP cache and, if it finds the address, provides it so that the packet can be converted to the right packet length and format and sent to the machine. If no entry is found for the IP address, ARP broadcasts a request packet in a special format to all the machines on the network to determine if any machine knows who has that IP address. A machine that recognizes the IP address as its own returns a reply so indicating. ARP updates the ARP cache for future reference and then sends the packet to the MAC address that replied. Here is a sample ARP broadcast query:



One might deduct that this addressing scheme could also be spoofed to provide a host with incorrect information "ARP Spoofing involves constructing forged ARP request and reply packets. By sending forged ARP replies, a target computer could be convinced to send frames destined for computer A to instead go to computer B." This referred to as ARP poisoning. There are currently programs that automate the process of ARP poisoning – ARPoison, Ettercap, and Parasite. All three have the capability to provide spoofed ARP packets and therefore redirect transmission, intercept packets, and/or perform some type of man in the middle attack. Either enabling MAC binding at a switch or implementing static ARP tables achieves prevention of ARP spoofing. MAC binding makes it so that once an address is assigned to an adapter; it cannot be changed without authorization. Static ARP management is only realistically achieved in a very small network. In a large dynamic network, it would be impossible to manage the task of keeping the entries updated. ARPWATCH, for UNIX based systems, monitors changes to the ARP cache and alerts administrator as to the changes.
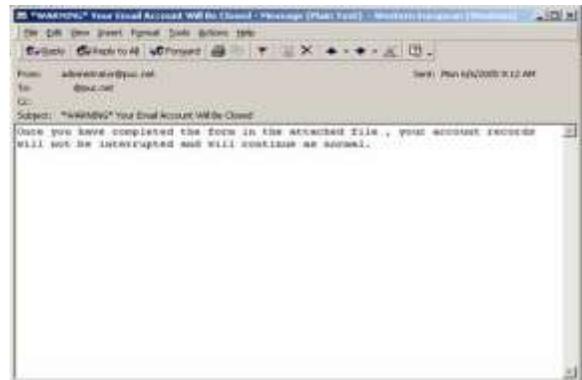
## IV. E-MAIL ADDRESS SPOOFING

Spoofing is when an e-mail message appears to come from a legitimate source but in fact is from an impostor. E-mail spoofing can be used for malicious purposes such as spreading viruses, trawling for sensitive business data and other industrial espionage activities [8].

If you receive a snail mail letter, you look to the return address in the top left corner as an indicator of where it originated. However, the sender could write any name and address there; you have no assurance that the letter really is from that person and address. E-mail messages contain return addresses, too – but they can likewise be deliberately misleading, or "spoofed." Senders do this for various reasons, including:

- The e-mail is spam and the sender doesn't want to be subjected to anti-spam laws
- The e-mail constitutes a violation of some other law (for example, it is threatening or harassing)
- The e-mail contains a virus or Trojan and the sender believes you are more likely to open it if it appears to be from someone you know
- The e-mail requests information that you might be willing to give to the person the sender is pretending to be (for example, a sender might pose as your company's system administrator and ask for your network password), as part of a "social engineering" attack
- The sender is attempting to cause trouble for someone by pretending to be that person (for example, to make it look as though a political rival or personal enemy said something he/she didn't in an e-mail message)

Here is an example of a spoofed email made out to look like it originated from administrator@puc.net



## V. WEB SPOOFING

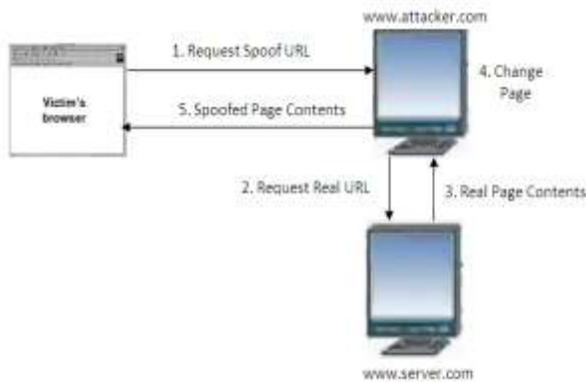As with the other forms of spoofing Web or Hyperlink spoofing provides victims with false information. Web Spoofing is an attack that allows someone to view and modify all web pages sent to a victim's machine. They are able to observe any information that is entered into forms by the victim. This can be of particular danger due to the nature of information entered into forms, such as addresses, credit card

numbers, bank account numbers, and the passwords that access these accounts [4].

Web Spoofing works on both Internet Explorer and Netscape and is not necessarily prevented by secure connections. This is due the way that the SSL protocol uses certificates to authenticate websites. The attacker can observe and modify all web pages and form submissions, even when the browser is indicating that there is a secure connection. The attack can be implemented using JavaScript and Web server plug-ins, and works in two parts. First, the attacker causes a browser window to be created on the victim's machine, with some of the normal status and menu information replaced by identical-looking components supplied by the attacker. Then, the attacker causes all Web pages destined for the victim's machine to be routed through the attacker's server. On the attacker's server, the pages are rewritten in such a way that their appearance does not change at all, but any actions taken by the victim (such as clicking on a link) would be logged by the attacker. In addition, any attempt by the victim to load a new page would cause the newly loaded page to be routed through the attacker's server, so the attack would continue on the new page. The attack is initiated when the victim visits a malicious Web page, or receives a malicious email message.

Current browsers do not completely prevent Web Spoofing, and there seems to be little movement in the direction of addressing this problem. I believe that there can be no fully secure electronic commerce on the Web until the Spoofing vulnerability has been addressed.



## VI. DNS SPOOFING

A DNS spoofing attack can be defined as the successful insertion of incorrect resolution information by a host that has no authority to provide that information. It may be conducted using a number of techniques ranging from social engineering through to exploitation of vulnerabilities within the DNS server software itself. Using these techniques, an attacker may insert IP address information that will redirect a customer from a legitimate website or mail server to one under the attacker's control – thereby capturing customer information through common man-in-the-middle mechanisms [9].

According to the most recent "Domain Health Survey" (Feb 2003), a third of all DNS servers on the Internet are vulnerable to spoofing.

Operating normally, a customer can expect to query their DNS server to discover the IP address of the named host they wish to connect to. The following diagram reflects this process.
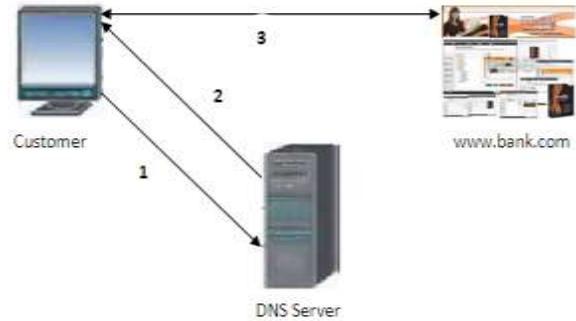


Fig 1: The Normal DNS Motion Process

1. The customer queries the DNS server – "What is the IP address of www.bank.com?"

2. The DNS responds to the customer query with "The IP address of www.bank.com is 150.10.1.21"

3. The Customer then connects to the host at 150.10.1.21 – expecting it to be www.bank.com.

However, with a successful DNS spoofing attack, the process has been altered. The following diagram reflects this process.



Fig 2: The DNS motion process having fallen victim to a DNS spoofing attack

1. The attacker targets the DNS service used by the customer and adds/alters the entry for www.mybank.com – changing the stored IP address from 150.10.1.21 to the attacker's fake site IP address (200.1.1.10).

2. The customer queries the DNS server "What is the IP address of www.bank.com"

3. The DNS responds to the customer query with "The IP address of www.bank.com is 200.1.1.10" – not the real IP address.

4. The Customer then connects to the host at 200.1.1.10 – expecting it to be www.bank.com, but in fact reaching the attackers fake site.

## VII.   AVOIDANCE OF SPOOFING

The Packet filtering is the best method to avoid various spoofing attacks. In this section we have described three packet filtering methods which are used to filter the spoofed packets, they are

A.   Ingress Filtering Method - IFM

B.   Egress Filtering Method - EFM

C.   Spoofing Prevention Method - SPM

### A.   Ingress Filtering Method

- Ingress filtering is a technique used to make sure that incoming packets are actually from the networks that they claim to be from [7].

- Networks receive packets from other networks. Normally a packet will contain the IP address of the computer that originally sent it. This allows other computers in the network to know where it came from, which is needed for things like sending a packet back to the sending computer [7].

- In certain cases, the sending IP address will be spoofed. This is usually done as part of an attack, so that the attacked computer does not know where the attack is really coming from [7].

- Filtering a packet is when the packet is not processed normally, but is denied in some way. The computer processing the packet might simply ignore the packet completely, or where it is possible it might send a packet back to the sender saying the packet is denied.

- In ingress filtering, packets coming into the network are filtered if the network sending it should not send packets from IP addresses of the originating computer.

- In order to do ingress filtering, the network needs to know which IP addresses each of the networks it is connected to may send. This is not always possible. For instance, a network that has a single connection to the Internet has no way to know if a packet coming from that connection is spoofed or not [7].

- Ingress filtering is a packet filtering technique used by many Internet service providers to try to prevent source address spoofing of Internet traffic, and thus indirectly combat various types of net abuse by making Internet traffic traceable to its source [7].

- Ingress filtering is a "good neighbor" policy which relies on mutual cooperation between ISPs for their mutual benefit.

- There are many possible ways of implementing this policy; one common mechanism is to enable reverse path forwarding on links to customers, which will indirectly apply this policy based on the provider's route filtering of their customers' route announcements [7].

### B.   Egress Filtering Method

- Egress filtering is the practice of monitoring and potentially restricting the flow of information outbound from one network to another. Typically it is information from a private TCP/IP computer network to the Internet that is controlled [7].

- TCP/IP packets that are being sent out of the internal network are examined via a router or firewall. Packets that do not meet security policies are not allowed to leave - they are denied "egress" [7].

- Egress filtering helps ensure that unauthorized or malicious traffic never leaves the internal network [7].

- In a corporate network, typically all traffic except that emerging from a select set of servers would be denied egress. Restrictions can further be made such that only select protocols such as http, email, and DNS are allowed. User workstations would then need to be set to use one of the allowed servers as a proxy. Direct access to external networks by the internal user workstation would not be allowed [7].

- Egress filtering may require policy changes and administrative work whenever a new application requires external network access. For this reason egress filtering is an uncommon feature on consumer and very small business networks [7].

### C.   Spoofing Prevention Method (SPM)

A new approach for filtering spoofed IP packets, called Spoofing Prevention Method (SPM). The method enables routers closer to the destination of a packet to verify the authenticity of the source address of the packet. This stands in contrast to standard ingress filtering which is effective mostly at routers next to the source and is ineffective otherwise. In the proposed method a unique temporal key is associated with each ordered pair of source destination networks (AS's, autonomous systems). Each packet leaving a source network $S$ is tagged with the key $K(S;D)$, associated with $(S;D)$, where $D$ is the destination network. Upon arrival at the destination network the key is verified and removed. Thus the method verifies the authenticity of packets carrying the address $s$ which belongs to network $S$. An efficient implementation of the method, ensuring not to overload the routers, is presented [5]. The major benefits of the method are the strong incentive it provides to network operators to implement it, and the fact that the method lends itself to stepwise deployment, since it benefits networks deploying the method even if it is implemented only on parts of the Internet. These two properties, not shared by alternative approaches, make it an attractive and viable solution to the packet spoofing problem.

*D. Some other plans to avoid spoofing in web applications:*

- Use cryptographic signatures to exchange authenticated email messages. Authenticated email provides a mechanism for ensuring that messages are from whom they appear to be, as well as ensuring that the message has not been altered in transit. Similarly, sites may wish to consider enabling SSL/TLS in their mail transfer software. Using certificates in this manner increases the amount of authentication performed when sending mail [7].

- Configure your mail delivery daemon to prevent someone from directly connecting to your SMTP port to send spoofed email to other sites [7].

- Ensure that your mail delivery daemon allows logging and is configured to provide sufficient logging to assist you in tracking the origin of spoofed email [7].

- Consider a single point of entry for email to your site. You can implement this by configuring your firewall so that SMTP connections from outside your firewall must go through a central mail hub. This will provide you with centralized logging, which may assist in detecting the origin of mail spoofing attempts to your site [7].

- Educate your users about your site's policies and procedures in order to prevent them from being "social engineered," or tricked, into disclosing sensitive information (such as passwords). Have your users report any such activities to the appropriate system administrator(s) as soon as possible [7].

## VIII. CONCLUSION

With the current implementations of spoofing, the network security community needs to be aware of the magnitude and potential cost of these types of attacks. People can effectively maintain patching and monitoring of logs to minimize the potential damage.

Professionals must remain current with the Operating Systems that we use in our day to day activities. A steady stream of changes and new challenges is assured as the hacker community continues to seek out vulnerabilities and weaknesses in our systems and our networks.

The authors have stated that the paper they presented will cater the needs of novice researchers and students who are interested in information security.

## REFERENCES

[1] Daemon, Route, Infinity, "IP Spoofing Demystified", Phrack Magazine; 1996;

[2] "IP Address Spoofing and Hijacked Session Attacks"; 1/23/95
http://ciac.llnl.gov/ciac/bulletins/f-08.shtml;

[3] Neil B. Riser "Spoofing: An Overview of Some Spoofing Threats" from SANS Reading room.

[4] Felten, Balfanz, Dean, Wallach D.S., "Web Spoofing, An Internet Con Game"; http://bau2.uibk.ac.at/matic/spoofing.htm;

[5] A. Bermlerand H. Levy. "Spoofing prevention Method," INFOCOM'05, 2005.

[6] Whalen, Sean; "An Introduction to ARP Spoofing";packetstorm.security.com/papers/protocols/intro_to_arp_spoofing.pdf;6/25

[7] www.wikipedia.com

[8] http://www.puc.net/email_spoofing.htm

[9] http://www.technicalinfo.net/papers//index.htm

[10] Whalen, Sean; "An Introduction to ARP Spoofing"; packetstorm.securify.com/papers/protocols/intro_to_arp_spoofing.pdf; 6/25/01.

[11] DNSAbuse; http://packetstorm.securify.com/papers/protocols/mi004en.htm; 6/30/01.

## AUTHORS PROFILE

Ms. Dr D. Lalitha Bhaskari is an Associate Professor in the Department of Computer Science and Engineering of Andhra University. She did her PhD from JNTU Hyderabad in the area of Steganography and Watermarking. Her areas of interest include Theory of computation, Data Security, Image Processing, Data communications, Pattern Recognition and Cyber Forensics. Apart from her regular academic activities she holds prestigious responsibilities like Associate Member in the Institute of Engineers, Member in IEEE, Associate Member in the Pentagram Research Foundation, Hyderabad, India. She is also the recipient of "Young Engineers" Award from the prestigious Institution of Engineers (INDIA) for the year 2008 in Computer Science discipline.

Dr. Ch. Satyanarayana working as Associate Professor in the Department of Computer Science & Engineering, University College of Engineering, JNTU Kakinada for the last 12 years. He obtained his Ph.D from JNTU Hyderabad. He published 22 research papers in various International Journals and Conferences. Under his guidance 12 Research scholars working on different areas like Image Processing, Speech Recognition, Pattern Recognition.

Mr. P. Ramesh babu is an Assistant Professor in the Department of Information Technology of Rajamahendri Institute of Engineering & Technology - Rajahmundry. His research interests include Steganography, Digital Watermarking, Information security, Network communications and Cyber Forensics.

Mr.Ramesh babu did his M.Tech in Computer Science & Engineering from JNTU Kakinada University. He has 6 years of teaching and industrial experience.

# Key Management Techniques for Controlling the Distribution and Update of Cryptographic keys

T.Lalith
Senior Lecturer, Department of MCA
Sona College of Technology
Salem, Tamilnadu., India
lalithasrilekha@rediffmail.com

R.Umarani
Reader, Department of Comp. Science
Sri Saradha College for Women, Salem, Tamilnadu, India
umainweb@gmail.com

G.M.Kadharnawaz
Director,Department of MCA
Sona College of Technology
Salem, Tamilnadu., India

*Abstract-***Key management plays a fundamental role in cryptography as the basis for securing cryptographic techniques providing confidentiality, entity authentication, data origin authentication, data integrity, and digital signatures. The goal of a good cryptographic design is to reduce more complex problems to the proper management and safe-keeping of a small number of cryptographic keys, ultimately secured through trust in hardware or software by physical isolation or procedural controls. Reliance on physical and procedural security (e.g., secured rooms with isolated equipment), tamper-resistant hardware, and trust in a large number of individuals is minimized by concentrating trust in a small number of easily monitored, controlled, and trustworthy elements.**

## I. INTRODUCTION

Systems providing cryptographic services require techniques for initialization and key distribution as well as protocols to support on-line update of keying material, key backup/recovery, revocation, and for managing certificates in certificate-based systems.

Key management [1] is the set of techniques and procedures supporting the establishment and maintenance of keying relationships between authorized parties.

Key management encompasses techniques and procedures supporting:

- Initialization of system users within a domain

- Generation, distribution, and installation of keying material

- Controlling the use of keying material.

- Update, revocation, and destruction of keying material and

- Storage, backup/recovery, and archival of keying material.

## II. CLASSIFYING KEYS BY ALGORITHM TYPE AND INTENDED USE

The terminology of Table I is used in reference to keying material. A symmetric cryptographic system is a system

involving two transformations – one for the originator and one for the recipient – both of which make use of either the same secret key (symmetric key) or two keys easily computed from each other. An asymmetric cryptographic system is a system involving two related transformations – one defined by a public key (the public transformation), and another defined by a private key (the private transformation) with the property that it is computationally infeasible to determine the private transformation from the public transformation.

TABLE I  PRIVATE, PUBLIC, SYMMETRIC AND SECRET KEYS

| Term | Meaning |
|---|---|
| private key, public key | paired keys in an asymmetric cryptographic system |
| symmetric key | key in a symmetric (single-key) cryptographic system |
| secret | adjective used to describe private or symmetric key |

Table II indicates various types of algorithms commonly used to achieve the specified cryptographic objectives. Keys associated with these algorithms may be correspondingly classified, for the purpose of controlling key usage .The classification given requires specification of both the type of algorithm (e.g., encryption vs. signature) and the intended use (e.g., confidentiality vs. entity authentication).

TABLE II. TYPES OF ALGORITHMS COMMONLY USED TO MEET SPECIFIED OBJECTIVES

| Cryptographic objective | Algorithm type | |
|---|---|---|
| | public-key | symmetric-key |
| confidentiality | encryption | encryption |
| data origin authentication | signature | MAC |
| key agreement | Diffie-Hellman | various methods |

| entity authentication | 1.signature<br><br>2.decryption<br><br>3.Customized | 1.MAC<br><br>2.encryption |
|---|---|---|

## III. KEY MANAGEMENT OBJECTIVES, THREATS AND POLICY

Keying relationships in a communications environment involve at least two parties (a sender and a receiver) in real-time. In a storage environment, there may be only a single party, which stores and retrieves data at distinct points in time.

The objective of key management is to maintain keying relationships and keying material in a manner which counters relevant threats, such as:

- Compromise of confidentiality of secret keys.

- Compromise of authenticity of secret or public keys. Authenticity requirements include knowledge or verifiability of the true identity of the party a key is shared or associated with.

- Unauthorized use of secret or public keys.

### A. Security policy and key management

Key management is usually provided within the context of a specific security policy. A security policy explicitly or implicitly defines the threats a system is intended to address. The policy may affect the stringency of cryptographic requirements, depending on the susceptibility of the environment in question to various types of attack. Security policies typically also specify:

- Practices and procedures to be followed in carrying out technical and administrative aspects of key management, both automated and manual responsibilities and accountability of each party involved and the types of records (audit trail information) to be kept, to support subsequent reports or reviews of security-related events.

## IV. TRADE OFFS AMONG KEY ESTABLISHMENT PROTOCOLS

In selected key management applications, hybrid protocols involving both symmetric and asymmetric techniques offer the best alternative. More generally, the optimal use of available techniques generally involves combining symmetric techniques for bulk encryption and data integrity with public-key techniques for signatures and key management.

### A. Public-key vs. symmetric-key techniques (in key Management)

Primary advantages offered by public-key (vs. symmetric-key) techniques for applications related to key management include:

- Simplified key management. To encrypt data for another party, only the encryption public key of that party need be obtained. This simplifies key management as only authenticity of public keys is required, not their secrecy.

- On-line trusted server not required. Public-key techniques allow a trusted on-line server to be replaced by a trusted off-line server plus any means for delivering authentic public keys (e.g., public-key certificates and a public database provided by an entrusted on-line server). For applications where an on-line trusted server is not mandatory, this may make the system more amenable to scaling, to support very large numbers of users.

- Enhanced functionality. Public-key cryptography [2] offers functionality which typically cannot be provided cost-effectively by symmetric techniques (without additional online trusted third parties or customized secure hardware). The most notable such features are non-repudiation of digital signatures, and true (single-source) data origin authentication.

## V. TECHNIQUES FOR DISTRIBUTING PUBLIC KEYS

Protocols involving public-key cryptography are typically described assuming a priori possession of (authentic) public keys of appropriate parties. This allows full generality among various options for acquiring such keys. Alternatives for distributing explicit public keys with guaranteed or verifiable authenticity, including public exponentials for Diffie-Hellman key agreement (or more generally, public parameters), include the following.

- Point-to-point delivery over a trusted channel. Authentic public keys of other users are obtained directly from the associated user by personal exchange, or over a direct channel, originating at that user, and which (procedurally) guarantees integrity and authenticity (e.g., a trusted courier or registered mail). This method is suitable if used infrequently (e.g., one-time user registration), or in small closed systems. A related method is to exchange public keys and associated information over an untrusted electronic channel, and provide authentication of this information by communicating a hash thereof (using a collision-resistant hash function) via an independent, lower bandwidth authentic channel, such as a registered mail..

- Use of an off-line server and certificates. In a one-time process, each party A contacts an off-line trusted party referred to as a *certification authority* (CA), to register its public key and obtain the CA's signature verification public key (allowing verification of other users' certificates). The CA certifies A's public key by binding it to a string identifying A, thereby creating a certificate. Parties obtain authentic public keys by exchanging certificates or extracting them from a public directory.

- Use of systems implicitly guaranteeing authenticity of public parameters. In such systems, including identity-based systems and those using implicitly certified keys

by algorithmic design, modification of public parameters results in detectable, non-compromising failure of cryptographic techniques.

## VI. PUBLIC KEY CERTIFICATES

Public-key certificates are a vehicle by which public keys may be stored, distributed or forwarded over unsecured media without danger of undetectable manipulation. The objective is to make one entity's public key available to others such that its authenticity (i.e., its status as the true public key of that entity) and validity are verifiable. In practice, X.509 certificates are commonly used.

### A. Definition

A public-key certificate [4] is a data structure consisting of a data part and a signature part. The data part contains clear text data including, as a minimum, a public key and a string identifying the party (subject entity) to be associated there with. The signature part consists of the digital signature of a certification authority over the data part, thereby binding the subject entity's identity to the specified public key.

The Certification Authority (CA) is a trusted third party whose signature on the certificate vouches for the authenticity of the public key bound to the subject entity. The significance of this binding (e.g., what the key may be used for) must be provided by additional means, such as an attribute certificate or policy statement. Within the certificate, the string which identifies the subject entity must be a unique name within the system (distinguished name), which the CA typically associates with a real-world entity. The CA requires its own signature key pair, the authentic public key of which is made available to each party upon registering as an authorized system user.

### B. Creation of public-key certificates

Before creating a public-key certificate for a subject entity A, the certification authority should take appropriate measures (relative to the security level required, and customary business practices), typically non-cryptographic in nature, to verify the claimed identity of A and the fact that the public key to be certified is actually that of A. Two cases may be distinguished.

- Trusted party creates key pair. The trusted party creates a public-key pair, assigns it to a specific entity, and includes the public key and the identity of that entity in the Certificate. The entity obtains a copy of the corresponding private key over a secure (authentic and private) channel after proving its identity (e.g., by showing a passport or trusted photo-id, in person). All parties subsequently using this certificate essentially delegate trust to this prior verification of identity by the trusted party.

- Entity creates own key pair. The entity creates its own public-key pair, and securely transfers the public key to the trusted party in a manner which preserves authenticity.(e.g., over a trusted channel, or in person). Upon verification of the authenticity (source

of the public key, the trusted party creates the public-key certificate the signer.

### C. Use and verification of public-key certificates

The overall process whereby a party B uses a public-key certificate to obtain the authentic public key of a party A may be summarized as follows:

- (One-time) acquire the authentic public key of the certification authority.

- Obtain an identifying string which uniquely identifies the intended party A.

- Acquire over some unsecured channel (e.g. from a central public database of certificates, a public-key certificate corresponding to subject entity A and agreeing with the previous identifying string.

### D. Attribute certificates

Public-key certificates bind a public key and an identity, and include additional data fields necessary to clarify this binding, but are not intended for certifying additional information. Attribute certificates are similar to public-key certificates, but specifically intended to allow specification of information (attributes) other than public keys (but related to a CA, entity, or public key), such that it may also be conveyed in a trusted (verifiable) manner. Attribute certificates may be associated with a specific public key by binding the attribute information to the key by the method by which the key is identified, e.g., by the serial number of a Corresponding public-key certificate, or to a hash-value of the public key or certificate. Attribute certificates may be signed by an attribute certification authority, created in conjunction with an attribute registration authority, and distributed in conjunction with an attribute directory service More generally, any party with a signature key and appropriate recognizable authority may create an attribute certificate.

## VII. KEY LIFE CYCLE ISSUES

Key management is simplest when all cryptographic keys are fixed for all time. Crypto periods [3] necessitate the update of keys. This imposes additional requirements, e.g., on certification authorities which maintain and update user keys. The set of stages through which a key progresses during its existence, referred to as the life cycle of keys, is discussed in this section.

### A. Lifetime protection requirements

Controls are necessary to protect keys both during usage and storage. Regarding long-term storage of keys, the duration of protection required depends on the cryptographic function (e.g., encryption, signature, data origin authentication/integrity) and the time-sensitivity of the data in question.

Security impact of dependencies in key updates: Keying material should be updated prior to crypto period expiry. Update involves use of existing keying material to establish

new keying material, through appropriate key establishment protocols and key layering .To limit exposure in case of compromise of either long term secret keys or past session keys, dependencies among keying material should be avoided. For example, securing a new session key by encrypting it under the old session key is not recommended (since compromise of the old key compromises the new).

### B. Key management life cycle

Except in simple systems where secret keys remain fixed for all time, crypto periods associated with keys require that keys be updated periodically. Key update necessitates additional procedures and protocols, often including communications with third parties in public-key systems. The sequence of states which keying material progresses through over its lifetime is called the key management life cycle.

Life cycle stages may include:

- User registration
- User initialization
- Key generation
- Key installation
- Key registration
- Normal use
- Key backup
- Key update
- Archival
- key de-registration and destruction
- Key recovery
- Key revocation

## VIII. CONCLUSION

Key management plays a fundamental role in cryptography as the basis for securing cryptographic techniques providing confidentiality, entity authentication, data origin authentication, data integrity, and digital signatures. The goal of a good cryptographic design is to reduce more complex problems to the proper management and safe-keeping of a small number of cryptographic keys, ultimately secured through trust in hardware or software by physical isolation or procedural controls. Reliance on physical and procedural security (e.g., secured rooms with isolated equipment), tamper-resistant hardware, and trust in a large number of individuals is minimized by concentrating trust in a small number of easily monitored, controlled, and trustworthy elements.

### REFERENCES

[1] National Institute of Standards and Technology.

[2] R. Rivest, A. Shamir, L. Adleman. A Method for Obtaining Digital Signatures and Public-Key Cryptosystems. Communications of the ACM, Vol. 21 (2), pp.120–126. 1978. Previously released as an MIT "Technical Memo" in April 1977, and published in Martin Gardner's Scientific American.

[3] http://www.discretix.com/PDF/Using%20Public%20Key%20Cryptography%20in%20Mobile%20Phones.pdf

[4] http://dlc.sun.com/pdf/316194901A/316194901A.pdf

[5] http://www.safecomprogram.gov/NR/rdonlyres/7C31664D-5B0B-4128-A85B-DC79B1D734ED/0/Security_Issues_Analysis_Report.pdf