

# ICTIMES - 2017

ISBN No: 978-93-85100-13-0

*Proceedings of International Conference on Emerging Technologies  
in Computer Science (ICETCS)*

*Under the Auspices of International Conference on trends in Information,  
Management, Engineering and Sciences (ICTIMES)*



## MALLA REDDY COLLEGE OF ENGINEERING (MRCE)

Permanently Affiliated to JNTUH, Approved by AICTE (New Delhi), Accredited by NBA, An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad, India - 500 100.  
[www.mrce.in](http://www.mrce.in)



Convenor  
Dr. V. Bhoopathy  
Prof. - CSE Dept.

Editor  
Dr. Sunil Tekale  
HOD - CSE Dept.

Editor in Chief  
Dr. P. John Paul  
Principal

## Proceedings of *International Conference on Emerging Technologies in Computer Science (ICETCS)*

### Chief Patron:

Sri Ch. Malla Reddy, Founder Chairman, MRGI  
(Member of Parliament, Govt., INDIA)

### Patrons:

Mr. Ch. Mahender Reddy, Secretary MRGI  
Dr. Ch. Bhadra Reddy, Treasurer MRGI

### International Advisory Committee:

Col. G. Ram Reddy, Director (Admin), MRGI  
Mr. N Sudhir Reddy, Director, Administration, MRCE  
Dr. S. R. C. Murthy, University of Sydney, Australia  
Dr. A.V. Vidya Sagar, BELL, USA  
Dr. K.V.S. S. Narayana Rao, NITIE, Bombay  
Dr. K. Vijay Kumar, CEO, First ESCO India, Vizag  
Dr. Ch. A.V. Prasad, Senior Consultant, TCS  
Dr. A. Govardhan, Principal, JNTUH  
Dr. B. Sudeer Prem Kumar, Chairman, BOS,  
JNTUH  
Dr. K. Venkateswar Rao, JNTUH  
Dr. P. Dasharathan, JNTUH  
Dr. B.N. Bhandari, Director DAP, JNTUH  
Dr. M. Manzoor Hussain, Director Administrations,  
JNTUH  
Dr. M. Madhavi Latha, Former Director, I-Tech,  
JNTUH

Dr. V.C.V. Prathap Reddy, RIT Rochester, USA

Dr. S. Venkateswara Rao, Head- Physics, JNTUH

Mr. N. Shyam Kumar, Group Manager, Tech Mahindra

Mr. S. Goutam, Manager, TCS

Dr. Hussain Reddy, SKU, A.P.

Dr. Seow Ta Wee, University Tun Hussein Onn Malaysia.

Ir. Dr. Goh Hui Hwang, Malaysia

### Chief Guest:

Dr. T.G Thomas, Dean-Admissions (Campus Wide) -  
BITS Pilani, Dubai Campus, Academic City, Dubai, UAE.

### Guest of Honor:

Dr. Balajied Lang Nongrum, Biola University USA.

### Keynote Speakers:

Dr. T.G Thomas, BITS Pilani, Dubai UAE

Dr. Balajied Lang Nongrum, Biola University USA

Dr. C. Krishna Mohan, IIT Hyderabad

Dr. R.Thundil Karuppa Raj, VIT University

Dr. K. Ramulu, Central University, Hyderabad

### Conference General Chair:

Dr. P John Paul, Principal, MRCE



## MALLA REDDY COLLEGE OF ENGINEERING (MRCE)

Permanently Affiliated to JNTUH, Approved by AICTE (New Delhi), Accredited by NBA, An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad, India - 500 100.

[www.mrce.in](http://www.mrce.in)

## *Proceedings of International Conference on Emerging Technologies in Computer Science (ICETCS)*

### **Organizing Committee:**

Dr. M. Thamarai, Dean Academics (Chair)  
Dr. V. Bhoopathy, Dean R&D - CSE  
Dr. P.Velmurgan, Dean R&D - MECH  
Dr. Nikhil Raj, Dean R&D - ECE  
Dr. Ch. Shankar, Dean Academics - MBA  
Dr. J. Gladson, Dean Student Affairs - CSE

### **Co- ordination Committee:**

Prof. Rajesh Durgam  
Prof. M. Shiva kumar  
Prof. Ch. Vijaya Kumari  
Prof. C. Shashi Kanth  
Prof. J. Shashi Kumar

### **Program Committee**

Dr. T. V. Reddy, Vice Principal (Chair )  
Dr. T. Sunil, Dean Academics - CSE  
Dr. S.S Gowda, Dean Academics - MECH  
Dr. G. Sridhar, Dean Student Affairs - ECE  
Dr. A. Karthikeyan, Dean Student Affairs - MECH

### **Information Contact:**

Dr. P John Paul  
Principal, MRCE  
+91-9348161222, 9346162620  
E-mail: principal@mrce.in

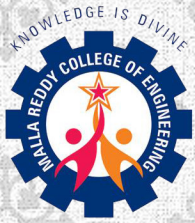


## **MALLA REDDY COLLEGE OF ENGINEERING (MRCE)**

Permanently Affiliated to JNTUH, Approved by AICTE (New Delhi), Accredited by NBA, An ISO 9001:2015 Certified Institution.

Maisammaguda, Hyderabad, India - 500 100.

[www.mrce.in](http://www.mrce.in)



# Malla Reddy College of Engineering

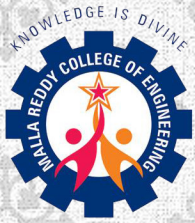
Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



Sri. Ch. Malla Reddy  
Founder Chairman, MRGI  
Member of Parliament

Best Wishes:

I Congrtulate CSE Department on Conducting  
International Conference on " Emerging  
technologies in Computer Science" (ICETCS)



# Malla Reddy College of Engineering

Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



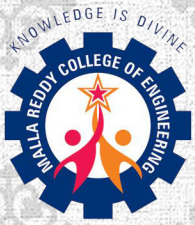
Sri. Ch. Mahender Reddy  
Secretary, MRGI



Sri. Ch. Bhadra reddy  
Treasurer, MRGI

Best Wishes:

We Congrtulate CSE Department on Conducting  
International Conference on " Emerging  
technologies in Computer Science" (ICETCS)



# Malla Reddy College of Engineering

Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



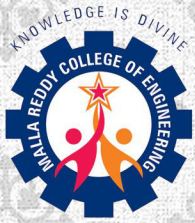
Col. G. Ram Reddy  
Director/Administrations, MRGI



Sri. N. Sudhir Reddy  
Director, MRCE

Best Wishes:

We Congrtulate CSE Department on Conducting  
International Conference on " Emerging  
Technologies in Computer Science" (ICETCS)



# Malla Reddy College of Engineering

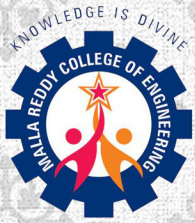
Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



Dr. P. John Paul  
Principal, MRCE  
Editor in Chief

Best Wishes:

Technology has to be invented or adopted.  
My wishes to CSE Department on Conducting  
International Conference on " Emerging  
technologies in Computer Science" (ICETCS)



# Malla Reddy College of Engineering

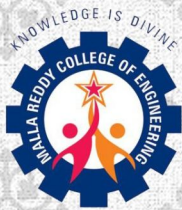
Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



**Dr. T.G Thomas**  
Dean Admissions (Campus Wide) BITS Pilani,  
Dubai Campus, Academic city, Dubai, UAE.

Best Wishes:

My warmest congratulations to you,  
MRCE and all staff on conducting  
International Conference on " Latest Trends in  
Electronics and Communication" (ICLTEC)



# Malla Reddy College of Engineering

Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi),Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)

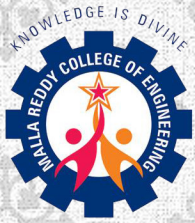


Dr. Balajied Lang Nongrum

Biola University, USA.

*Best Wishes:*

*My warmest congratulations to you,  
MRCE and all staff on conducting  
International Conference on "Latest Trends in  
Electronics and Communication" (ICLTEC)*



# Malla Reddy College of Engineering

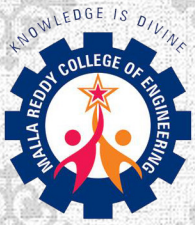
Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



Dr. C. Krishna Mohan  
IIT Hyderabad

Best Wishes:

My warmest congratulations to you,  
MRCE and all staff on conducting  
International Conference on " Emerging  
Technologies in Computer Science" (ICETCS)



# Malla Reddy College of Engineering

Permanently Affiliated to JNTUH, Approved by AICTE(New Delhi), Accredited by NBA  
An ISO 9001:2015 Certified Institution.  
Maisammaguda, Hyderabad - 500 100.  
[www.mrce.in](http://www.mrce.in)



Dr. Sunil Tekale  
HOD - CSE Dept. MRCE



Dr. V. Bhoopathy  
Prof. - CSE Dept. MRCE

Best Wishes:

We Congrtulate CSE Department on Conducting  
International Conference on " Emerging  
Technologies in Computer Science" (ICETCS)

S. No.	Title	Page No.
IC17CS01	IOT SECURITY PROVIDING TO AVOID DEFENSES AND ATTACKS <i>--R. Raja Kumar, Dr. B. Ramasubba Reddy</i>	1
IC17CS02	LIFE TIME ENHANCEMENT FOR WIRELESS SENSOR NETWORK USING FUZZY-ACO COMBINATION ALONG WITH AN ENHANCED CLUSTERING SCHEME <i>--Mr. Arshad Shareef, Mrs. Ursa Sayeed, Dr. P A Abdul Saleem</i>	8
IC17CS03	EVALUATION BASED LOAD BALANCING FOR CLOUD COMPUTING <i>--Dr.Sunil Tekale, Mr.M.Amarnath</i>	22
IC17CS04	TIME-OF-ARRIVAL (TOA), ANGLE-OF-ARRIVAL (AOA) AND HYBRID - TOA &AOA BASED LOCALIZATION IN WIRELESS SENSOR NETWORKS <i>--Dr.V.Bhoopathy, M Aharonu</i>	27
IC17CS05	SPEAKER SEGMENTATION- AN COMPARATIVE STUDY USING SUPPORT VECTOR MACHINES AND AUTO ASSOCIATIVE NEURAL NETWORK <i>--Dr.J.Gladson Maria Britto, Mrs.B.Ananthi</i>	33
IC17CS06	CONSISTENT DATA DELIVERY IN MOBILE ADHOC NETWORKS <i>--R.Nanda Kumar, Dr.Sankara, Malliga G</i>	40
IC17CS07	DATA WEB FOR QUERY FORMULATION LANGUAGE <i>--Dr. G. Silambarasan, Dr. V. Chandrasekar</i>	45
IC17CS08	MOVING HUMAN ACTION RECOGNITION AND IMAGE CLASSIFICATION <i>--Dr. Manikandan.P, Dr.S.P. Anandaraj</i>	51
IC17CS09	DIGITAL IMAGE ENCRYPTION BASED ON FUZZY LOGIC <i>--S R Mahipal, V V Ramanjaneyulu, Dr. Sunil Tekale</i>	59

IC17CS10	RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF GRAPHS	63
	-- <i>M. S. Malchijah Raj</i>	
IC17CS11	TIME SLICING APPROACH FOR RESOURCE ALLOCATION IN CLOUD COMPUTING	68
	-- <i>M.Aharonu, V .DeviPriya, Dr.Sunil Tekale</i>	
IC17CS12	A NOVEL WEIGHTED SCAN-BASED TEST PATTERN FOR BUILT-IN SELF-TEST	73
	-- <i>D.Gaspin Beautly</i>	
IC17CS13	ENERGY-EFFICIENT SECURE DATA AGGREGATION FRAMEWORK (ESDAF) PROTOCOL IN HETEROGENEOUS WIRELESS SENSOR NETWORKS	83
	-- <i>Dr. G. Silambarasan, , Dr. V.Bhoopathy Dr. V. Chandrasekar</i>	
IC17CS14	SECURITY AND PRIVACY ISSUES OF HEALTHCARE APPLICATION AND IMPLICATION OF PREDICTIVE ANALYTICS IN BIG DATA	91
	-- <i>A.S. Gousia Banu, D. Saritha, K Narasimhulu</i>	
IC17CS15	<i>STUDENT LEARNING EXPERIENCE BY DATA MINING &amp; SOCIAL MEDIA</i>	96
	-- <i>S Rasheeduddin, N.Ananth Ram Reddy</i>	
IC17CS16	IMPROVING THE SECURITY ISSUES IN WIRELESS SENSOR USING HETEROGENEOUS ALGORITHM	102
	-- <i>Dr. G. Silambarasan, Dr. V. Chandrasekar</i>	
IC17CS17	IMPROVING ENERGY EFFICIENT IN WIRELESS SENSOR NETWORKS USING PATH ALGORITHM	109
	-- <i>Dr.A.Mummoorthy, Sudha Pavani. K</i>	
IC17CS18	AUTHORIZED AUDITING OF DYNAMIC BIG-DATA ON CLOUD	116
	-- <i>A.S. Gousia Banu, Pramod Kumar Singh</i>	
IC17CS19	SMART AGRICULTURE THROUGH IOT	121

	-- <i>Hari Krishna, Dr.T.Sunil</i>	
IC17CS20	SPECTRUM SENSING FOR PERFORMANCE IMPROVEMENT	126
	-- <i>V.Uday Kiran, Dr.V.Bhoopathy</i>	
IC17CS21	SURVEY ON CRIME ANALYSIS AND PREDICTION USING DATA MINING TECHNIQUES	131
	-- <i>G.Soundarya</i>	
IC17CS22	REACH CENTROID ALGORITHM IDENTIFY THE WIRELESS LOCALIZATION IN WIRELESS SENSOR NETWORKS	139
	-- <i>Dr. G. Silambarasan, Dr. V. Chandrasekar</i>	
IC17CS23	CYBERCRIME: A THREAT TO NETWORK SECURITY	144
	-- <i>Ch.Vijaya Kumari, Ch.Vengaiah</i>	
IC17CS24	ROUTING INFORMATION PROTOCOL FOR WIRELESS SENSOR NETWORKS	149
	-- <i>Dr. V. Chandrasekar, P.Pavani</i>	
IC17CS25	REALISTIC FUTURE FLYING CAR	154
	-- <i>Mr.U.Nagaiah, Ms.Razia sultana, Mr.M.Amarnath, Dr.Sunil Tekale</i>	
IC17CS26	OPTIMAL JAMMING ATTACK DETECTION IN WSN	159
	-- <i>P.Poovizhi, S.Yamuna M.Brindha, V.Poorani</i>	
IC17CS27	VIRTUALIZATION SECURITY FOR CLOUD COMPUTING	165
	-- <i>G.Mamatha, K.HimaBindu</i>	
IC17CS28	DISTRIBUTED TRACKING SYSTEM	175
	-- <i>Ajju P. Benny, Dr.V.Bhoopathy</i>	
IC17CS29	A TECHNIQUE IN COMMUNICATION WITH CLOUD USING RPC	180
	-- <i>K.Bharath, Ch.Vijaya Kumari</i>	
IC17CS30	BLUE DOG	184
	-- <i>D. Namratha, Dr.Chandra Shekar</i>	
IC17CS31	A SURVEY OF NATURE INSPIRED LOAD BALANCING ALGORITHMS IN A CLOUD COMPUTING ENVIRONMENT	188

	-- <i>Priyanka manikonda, Veerender A</i>	
IC17CS32	NEED FOR VARIOUS ENERGY EFFICIENT MECHANISMS IN THE WIRELESS SENSOR NETWORKS	195
	-- <i>Puladas Sandhya Priyanka, Rashmitha</i>	
IC17CS33	CAPTCHA BREAKING USING IMAGE TEMPLATE MATCHING AND MACHINE LEARNING ALGORITHMS	199
	-- <i>G.Madhuri, V.Sandhya</i>	
IC17CS34	UNSTRUCTURAL DATA USING HADOOP	230
	-- <i>K.Madan Mohan, R.Bangari</i>	
IC17CS35	MOBILE AGENT BASED SECURITY USING BIG DATA MANAGEMENT	237
	-- <i>A. S. Gousia Banu, D. Saritha.</i>	
IC17CS36	GOOGLE PROJECT LOON	243
	-- <i>V.Naveena, Dr.T.Sunil,</i>	
IC17CS37	RESOURCE PLANNER	249
	-- <i>P.Priskilla, Dr.T.Sunil,</i>	
IC17CS38	DATA STORAGE SECURITY IN A HOSTED ENVIRONMENT	254
	-- <i>Anitha Bejugama, Shravani Reddy</i>	
IC17CS39	AGRICULTURAL UPDATE VIA SMS	261
	-- <i>Supriya Chinna, Dr.Raja Sekar</i>	
IC17CS40	A SURVEY PAPER ON DATA SECURITY IN CLOUD COMPUTING	267
	-- <i>D.Sravani, Dr.P.Mani Kandan</i>	
IC17CS41	A DISTRIBUTED THREE-HOP ROUTING PROTOCOL TO INCREASE THE CAPACITY OF HYBRID WIRELESS NETWORKS	273
	-- <i>Mounika Dr.Chandra Shekar</i>	
IC17CS42	LI-FI TECHNOLOGY TRANSMISSION OF DATA THROUGH LIGHT	278
	-- <i>R.Sai Chandrika, Dr.Raja Sekar</i>	

IC17CS43	EXPRESSIVE, EFFICIENT, AND REVOCABLE DATA ACCESS CONTROL FOR MULTI-AUTHORITY CLOUD STORAGE	285
	-- <i>E Srinath, Ch.Malleswar Rao</i>	
IC17CS44	ORUTA: PRIVACY-PRESERVING PUBLIC AUDITING FOR SHARED DATA IN THE CLOUD	293
	-- <i>Byreddy Madhavi, Dr.V.Bhoopathy</i>	
IC17CS45	ANDROID SECURITY	297
	-- <i>A.Sindhu, Dr. P. Manikandan</i>	
IC17CS46	SEARCHING AND PREDICTING USING DATA MINING ALGORITHMS	305
	-- <i>P.Pravalika, Dr.Raja Sekar</i>	
IC17CS47	DESIGN AND IMPLEMENTATION OF COMPUTATIONAL DYNAMIC TRUST MODEL FOR INDIVIDUAL AUTHORIZATION	316
	-- <i>Thirumala Vasala, Danuka Nilima Priyadrshini</i>	
IC17CS48	TIME ORIENT SPATIAL TRAFFIC PATTERN BASED MITIGATION OF DISTRIBUTED DENIAL OF SERVICE ATTACKS WITH DJN IN DISTRIBUTED WIRELESS NETWORKS	331
	-- <i>A. Saraswath, Dr.K.Thangadurai, Dr.A.Mummoorthy</i>	
IC17CS49	USING DATA MINING TECHNIQUES ANALYSING ABOUT ROAD ACCIDENTS	339
	-- <i>J.Sofia, P.Sandeep</i>	
IC17CS50	DYNAMIC ACCESS CONTROL POLICIES IN MULTI CLOUD STORAGE BASED NCC CLOUDS	343
	-- <i>K.Divya Bharathi, N.Keerthi</i>	
IC17CS51	A CLOUD BASED SYSTEM TO SENSE SECURITY ULNERABILITIES OF WEB APPLICATION IN OPEN-SOURCE CLOUD IAAS	351
	-- <i>K.Himabindu, G.mamatha, S.Kavitha</i>	
IC17CS52	SENTIMENT ANALYSIS IN HEALTHCARE USING SOCIAL MEDIA	358
	-- <i>Rajasekhar Nennuri, Dr I Surya Prabha</i>	

# IoT Security providing to avoid defenses and attacks

R. Raja Kumar<sup>1</sup> Dr. B. Ramasubba Reddy<sup>2</sup>

<sup>1</sup>Assistant Professor, Dept of CSE , SV Engineering College for Women,, Tirupati, India

Email:raj.rampalli@gmail.com

<sup>2</sup>Professor, Dept of CSE, SV College of Engineering Tirupati, India,

Email:rsreddyphd@gmail.com

**Abstract**—Internet of Things (IoT) has been a massive advancement in the Information and Communication Technology (ICT). It is anticipated that more than 50 billion gadgets will turn out to be a piece of the IoT in the following couple of years. Security of the IoT network should be the foremost priority. In this paper, we assess the security challenges in the four layers of the IoT engineering and their answers proposed from 2010 to 2016. Furthermore, important security technologies like encryption are also analyzed in the IoT context. At last, we talk about countermeasures of the security assaults on various layers of IoT and feature the future research headings inside the IoT engineering. **Keywords**— *Internet of Things; IoT; Security; layer architecture;*

## I. INTRODUCTION

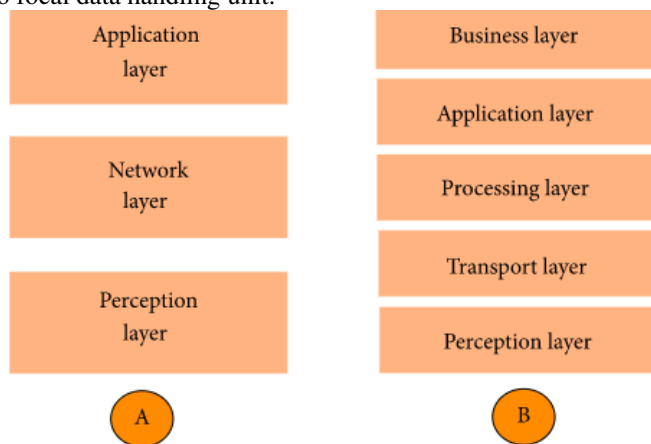
The IoT [1] emerges as a new concept in future Internet when the physical objects become the part of Internet. IoT provides objects the unique identity accessible from the network and its status, location can be track down[2]. Many facilities such as tracking monitoring and controlling become possible with the IoT which changes the human interactions with the physical objects. There are many devices developed which are now used as the IoT such as, RFID (Radio Frequency Identification Devices), infrared sensors, laser scanner, GPS (Global Positioning System) and gas inductors. In IoT various parameters of the objects or processes such as sound light mechanics, chemistry, biology, and position can be monitored and controlled. The great thing about IoT is that all the information is based on real time data.

IoT comprises of a network of highly diverse digital objects interacted with each other and with humans too. It provides a sensor network with communication system, store and manage the information, provides access and also handles the privacy protection and data security problems [3]. Comparing the research aspects on security in IoT to security in Internet, former is the way complex than the later and therefore needs the significant attention of the researcher and a more precise research methodology and tools should be

incorporated. With respect to security, The most basic architecture is a three-layer architecture. The three-layer architecture defines the main idea of the Internet of Things, but it is not sufficient for research on IoT because research often focuses on finer aspects of the Internet of Things. That is why, we have many more layered architectures proposed in the literature. The five layers are perception, transport, processing, application, and business layers (see Figure 1). The role of the perception and application layers is the same as the architecture with three layers. We outline the function of the remaining three layers. As shown in figure 1.

### A. Perception layer:

This layer uses the distinctive sensors, for example, ZigBee, infrared, RFID and QR code to gather data. Data could be temperature, stickiness, vibration, drive, pH level weight, speed and so forth. Transmission of data gathered is brought out through network layer to focal data handling unit.



**Figure 1: Architecture of IoT (A: three layers) (B: five layers).**

Figure 1: IoT Layers

### B. Network Layer:

The transmission of information is carried out in this layer. Information transmitted through the different mediums such as RFTD, Infrared, satellite and Wi-Fi units depending upon the nature of sensors and sensitivity of data. Hence data transmits securely from perception layer to other layers through network layer.

### C. Application layer:

The application layer is in charge of conveying application particular administrations to the client. It characterizes different applications in which the Internet of Things can be deployed, for example, smart homes, smart cities, and smart health. The data gives a stage to such applications which could profit the client from multiple points of view, for example, wellbeing training individual utilize, contraptions, family unit, transportation, correspondence and so on. The data security in IoT ought to be furnished with highlight like classification, recognizable proof, and so on. As IoT is going to be connected in various essential fields like wellbeing, transportation, enterprises, brilliant homes, postal administrations, and so on in this manner the security and protection of IoT ought to be trick evidence. Directed answer for every security factor ought to be characterized.

### D. Processing layer:

This layer attempts to join the system layer and application layer. All the insightful and distributed computing is done in this layer. Bolster layer usefulness incorporates capacity of information from bring down level layers to information base and administration. On the basis of wise registering this layer can figure data and process information naturally.

### E. Transport layer:

The transport layer transfers the sensor data from the perception layer to the processing layer and vice versa through networks such as wireless, 3G, LAN, Bluetooth, RFID, and NFC.

### F. Business layer:

The business layer manages the whole IoT system, including applications, business and profit models, and users' privacy.

In this paper we discuss the three layered architecture of IoT and their security. We discuss different security features and security challenges of these layers. And on the basis of former research we discuss different security aspects like cryptography, communication security, protecting sensor data and outline the challenges briefly. Rest of the paper is organized as follow. Section II provides a brief overview of the security threats of each layer and their countermeasure. In Section III, performance evaluation is done on the basis of the literature. Section IV provides the work directions which are needed to be done in future and in Section V, the work done in the paper is concluded.

## II. LITERATURE REVIEW

In this section, we discuss different security threats and their countermeasure on each layer briefly which have been purposed recently.

### A. Perception Layer Attacks

Hardware attacks are the most common attacks on perception layer. Perception layer generally includes WSN, RFID, zigbee and other kind of sensors. The attacker needs to be in the network or physically close to the nodes of the IoT system. Some of the common attacks on the perception layer are listed below.

- *Hardware Tempering*

Attacker physically close to the nodes can damage the node by changing the parts of its hardware or completely replacing that node [5]. Changing the electronic integration or by capturing the gate way node, the attacker can get all the information on that network including routing table, communication key, cryptographic key, radio key etc and threat all the network including higher layers.

- *Fake node injection*

The attacker can inject a fake or malicious node between the nodes of the network [4], hence the attacker gain access to the network and be able to control all the data flow of the network. It can make the node to stop transmitting the real data and hence destroy the entire network.

- *Malicious code injection*

The nodes of the IoT network can also be compromised by injecting malicious code. Dos attacks in WSN or virus on to the nodes are the most common type of this attack [6]. By this attack the attacker can gain access to the network, and can make the network to lose resources and hence make the services unavailable.

- *Sleep denial attack*

Node on the remote places in IoT network are mostly powered by replaceable batteries, the nodes are programmed to sleep when they are not in use to increase their battery life[7]. In this attack the attacker keep the node awake and prevent them to fall asleep by feeding wrong input to the node which results in power consumption hence the node shutdown.

- *WSN Node Jamming*

Wireless sensor network works on the radio frequency. A denial of service can be created by sending the noise signals over the network or by jamming the signals of WSN. This deny the communication between the nodes of the IoT network. The attacker keep on jamming the signals which result in the denial of services of the IoT [8].

- *RF interference of RFIDs*

RFIDs also works on radio signals as mentioned in WSN network earlier. The difference is that attacker don't need to jam the signals, the attacker can create make the nodes to deny the services just by sending noise signals over the network[9]. This noise interferes the RFID signal which create a hurdle in communication of the nodes.

### B. Perception Layer Attacks Countermeasures

- *Authentication*

To keep malicious devices out of the IoT network authentication of the devices should be done before getting into the network[10]. Without proper authentication the device should not be allowed to communicate with the network which prevents false data flow in the network.

- *Data integrity*

Each device in the IoT network should be provided by error detection mechanism, which minimizes the risk of data tempering. There are different error detection mechanism which being used like parity bit, check sum, etc. To make it more secure, cryptographic hash function should be used[11].

- *Secure booting*

Cryptographic hash algorithm can be used to check integrity and the

authentication of the software on different devices of the IoT network. But in most cases the end devices of the network possess very low computing power. So, most of the hash algorithms cannot be implemented on these devices. WH and NH cryptographic algorithms are the best solution to this problem because they need very low processing to execute [12].

- *IPSec Security channel*

IPSec security provides two kinds of security features, authentication and encryption. Eavesdropping and data tempering can be avoided encoding the data which insure data confidentiality[13]. The sender of the data can be identified by the receiver that's the sender over the IP is real or not.

- *Anonymity*

The inject node in the network by the attacker can hide sensitive information like location, identity etc. These kinds of the nodes are sensed anonymous by the network. Solution to this problem is presented as k-anonymity approach [14]. This approach work very best on low processing devices.

- *Physical secure design*

To resolve most of the attacks of perception layer can be resolved by designing the end devices physically secure. It includes chip selection, radio frequency circuits, data acquisition unit design, etc. These components should be of high quality and should not be easily changeable. The design of antenna for wireless communication should be able to communicate over good distance[15].

### C. Network Layer Attacks

The attacks are targeted to IoT system network. Such attacks can be performed without being close to the network.

- *Traffic analysis attacks:*

The wireless technologies of transmission are sniffed to obtain confidential information such as RFID. In such cases hacker first obtain information related to network by using packet sniffers or port scanning application and then attacks on the targeted information[16].

- *RFID Spoofing*

In RFID spoofing attacker targets the RFID signal to gain access the information imprinted on RFID tag[17]. Once the signal spoofed hacker uses it to transmit his own data using the original id. Now hacker obtained the full access to system.

- *RFID unauthorized access*

As there is no secured authentication system in RFID systems, the tags are accessible to anyone. It means tags can be manipulated easily[18].

- *Sinkhole Attack*

The attack directs all signals from wireless sensor network nodes to a same point. Such attack voids the data safety and drops all the packets instead of delivering to its destination[19].

- *Man in the Middle Attack*

Web attacker interfere the two sensor nodes to access restricted information and violates the privacy of nodes[20].

Such attack doesn't demand the attacker to be physically appeared on the location of network. This can be done by using the communication protocol of the IoT.

- *Routing Information Attack*

These are immediate attacks that the enemy by spoofing, replaying or changing routing data can convolute the system and make routing loops, permitting or dropping movement, sending the false error messages, shortening or amplifying source courses or notwithstanding parceling the network[21].

### D. Network Layer Attacks Countermeassure

- *Data privacy*

Data privacy in the IoT can be achieved by preventing illegal access of the nodes of the network. Different authentication mechanism can be used for this purpose, one of them is point to point encryption[22]. In this method the confidential data is converted immediately to a code which is indecipherable

- *Routing security*

Secure routing is the key to the secure utilization of sensor systems for some applications, yet the larger numbers of routing conventions are unstable. Routing security for the sensor network can be ensured by routing the data through multiple paths which increase error detection of the network[23].

- *Sinkhole attack*

Attack which is from outside the network can be secured by encryption and authentication, so the attacker not be able to join the network. And the attack from inside the network can be secured by security aware ad hoc routing protocol (SAR)[24]. Ssecurity metrics is added to the packets of route request, after analyzing the received data the attacker can be dropped from the network.

- *Spoofing*

Spoofing attack can be encountered by GPS location system. In[25] the GPS system techniques has been described and implemented. It is not the perfect solution but it is one the best solution provided yet.

- *Data integrity*

Data integrity can be achieved by applying cryptographic hash functions on the data[26]. This ensures data is not tempered when it reaches the receiving end. Mitigation problem can be resolved by applying error correction mechanism.

### E. Processing Layer Attacks

Cloud attacks are the most common type of attacks in processing layer of IoT because the data is sent to the cloud at this phase. So we discuss some common attacks which can make the network vulnerable to threats.

- *Application security*

Most of the application on cloud SAAS are delivered through internet i.e. web services. An attacker can easily uses web to get into the IoT network and can steal the data or can perform

malicious activities. Security issues in SAAS are much different from usual web securities issues. OWASP had identified different security issues on SAAS [27].

- *Data security*

Data security is the major concern of a SAAS user. It's the responsibility of the SAAS provider to ensure the security, the data processed and stored on cloud as plain text. The major security issues occur on the facilitation of data backup provided by the service provider [28]. Data back is offered through third party in most of the cases which increase the treat of data theft.

- *Underlying infrastructure security*

In PaaS, lower layers of IoT are not accessible to the developer, underlying layer's security is the responsibility of the provider[29]. Even developer can develop a secure application but its security remain vulnerable due to lower layers of IoT.

- *Third-party relationships*

PaaS not only provides programming language, it also provide third party web service component i.e. mashups[30]. More than one source is combined in mashups which increase security issues like network and data security.

- *Virtualization threats*

Security of virtual machines is as important as the security of the physical machines and any defect in possibly one may influence the other[31]. Virtualization in processing layer is vulnerable to many types of attacks.

- *Shared Resources*

As virtual machines share same resources, this becomes a security threat to the network. Using covert channels an attacker can monitor all the shared resources between the virtual machines so the information might be compromised[32].

#### *F. Processing Layer Attacks Countermeasure*

- *Fragmentation redundancy scattering(FRS)*

In data FRS the sensitive data on the cloud is divided into different fragments and stored on different servers [33]. The fragments of the data don't have any significant information by itself, so the risk of data leakage is minimized.

- *Homomorphic encryption*

This method is based on full homomorphic encryption application. This method allows cipher texts to be computed arbitrarily without being decrypted. This method requires high computation but assure data security[34].

- *Web application scanners*

Web application scanner is program proposed in [35] which detect treats from the front end of the web. There are other web firewall applications which can detect a potential attacker.

- *Hyper Safe*

Hypersafe lockdowns and protects the write protected memory pages from being modified, pointing index is restricted that converts controlled data into the pointer indexes[36].

- *Encryption*

Encryption is used for securing the sensitive data. Data sent or stored on the cloud is in encrypted form. There are different type of encryption mechanisms like Advanced Encryption Standard etc[37]. It can also help in overcoming side channels attack.

#### *G. Application Layer Attacks*

Before In computer security, amenability of security is caused by software attacks. By using Trojan horse programs, worms, viruses, spyware and malicious scripts software attacks can develop the system that can harm IoT System devices, appropriate information, tamper with data and deny service.

- *Phishing Attacks*

From spoofing the user's conformation ID, confidential data can be accessed by attacker through contaminated web site or email[38].

- *Virus, Worms, Trojan Horse, Spyware*

System can be contaminated by opponents with nasty software that can results in pinching information, tampering data or even denial of service[39].

- *Malicious Scripts*

IoT network is generally associated with Internet. Entire system closes up and data stealing is caused by running executable active-x scripts take in by the user that reins the access[40].

- *Denial of Service*

Through application layer, IoT network can be exaggerated by the execution of DoS or DDoS attacks by the attackers that influence the users on the network. Genuine users can be infertile by these attacks and attackers can obtain complete access on application layer, databases and private sensitive data[41]

- *Data Protection and Recovery*

User privacy is involved in communication data. Data can be lost and even catastrophic damage can be caused imperfect algorithms and mechanisms of data processing and data protection[42]. The mass nodes management is also one reason.

- *Software Vulnerabilities*

Vulnerabilities occur due to the non standard code as it was written by programmers, as a result buffer runoff. This technique or method is used by the hackers to accomplish their rational[43].

- *Data security*

User privacy is involved in communication data. Data can be lost and even catastrophic damage can be caused imperfect algorithms and mechanisms of data processing and data

Table 1. Comparative Analysis

Layers	Attack Name	Attack References	Effects	Launch	Countermeasure	Countermeasure reference
Perception Layer	Hardware Tempering	[8]	Data leakage (Keys, routing tables, etc)	2011	Secure Physical Design	[13]
	Fake node injection	[7]	Fake Data Manipulation	2013	Secure Booting	[14]
	Malicious code injection	[9]	Halt Transmission	2012	Intrusion detection Technology(IDT)	[15]
	Sleep denial attack	[10]	Node shutdown	2012	Authentication	[16]
	WSN Node Jamming	[11]	Jam Node Communication	2010	IPSec Security channel	[17]
	RF interference of RFIDs	[12]	Distortion in node Communication	2012	Authentication	[18]
Network Layer	Traffic analysis attack	[19]	Data leakage (about network)	2013	Routing Security	[26]
	RFID Spoofing	[20]	Intrusion in network Data manipulation	2011	GPS Location System	[28]
	RFID unauthorized access	[21]	Node data can be modified (Read, Write & Delete)	2014	Network Authentication	[25]
	Sinkhole Attack	[22]	Data leakage (Data of the Nodes)	2013	Security Aware AdHoc Routing	[27]
	Man in the Middle Attack	[23]	Data Privacy Violation	2011	Point-to-Point Encryption	[29]
	Routing Information Attack	[24]	Routing loops (Network Destruction)	2011	Encrypting Routing Tables	[30]
Processing Layer	Application security	[30]	Privacy Violation	2014	Web Application Scanner	[37]
	Data security	[31]	Data leakage (User data on cloud)	2012	Homomorphic Encryption	[33]
	Underlying infrastructure security	[32]	Service Hijacking	2010	Fragmentation redundancy scattering	[38]
	Third-party relationships	[33]	Data Leakage (User data on cloud)	2013	Encryption	[39]
	Virtualization threats	[34]	Resources destruction	2012	Hyper Safe	[38]
	Shared Resources	[35]	Resources Theft	2011	Hyper Safe	[38]
Application Layer	Phishing Attacks	[40]	Data Leakage (User credentials data)	2016	Biometrics Authentication	[47]
	Virus, Worms, Trojan Horse, Spyware	[41]	Resource Destruction & Hijacking	2012	Protective Software	[48]
	Malicious Scripts	[42]	Hijacking	2011	Firewalls	[49]
	Denial of Service(DoS)	[43]	Resource Destruction	2010	Access Control Lists	[50]
	Data Protection and Recovery	[44]	Data loss & Catastrophic Damage	2011	Cryptographic Hash Functions	[46]
	Software Vulnerabilities	[45]	Buffer over flow	2011	Awareness of security	[51]

protection[44]. The mass nodes management is also one reason

#### H. Application Layer Attacks Countermeasure

- User Authentication

Encryption and Integrity mechanisms are significant for the privacy and protection of system beside data stealing; as a result unauthorized access and data of the system can be protected[45].

- Access Control Lists (ACLs)

For accessing and controlling the IoT system, special policies and permissions are made. Incoming or outgoing traffic and access request of network can be allowed or blocked by ACLs[46].

- Firewalls

If the password is weak, the password of authentication and encryption can be break. Packets can be filtered, blocked that are not required, unfriendly login attempts, and DoS attacks before even authentication process begins by using the firewall[47].

- Anti-virus, Anti-spyware and Anti-adware

For the security, confidentiality, reliability and integrity of IoT system these software are essential.

- Risk Assessment

Application layer can be secured by risk assessment in risk assessment technique continuously detects threats of the system, apply patches and updates of the firmware of the system devices which improve security of the system[48]

### III. PERFORMANCE EVALUATION

In this segment, we have assessed security dangers on the layers of IoT arrange and exhibited their countermeasures. We have consolidated diverse sorts of assaults, their names, the name of layer on which that particular assaults are done. Moreover, we have additionally featured the impacts of these assaults on the IoT arrange or the trade off that we need to make if these assaults are propelled. Propelling times of these assaults are likewise talked about and isolate countermeasures of the considerable number of assaults are introduced, applying which, we can limit the harms that these assaults may do.

The detail of our performance evaluation can be seen in table 1. The table classifies the attacks and preventive measures in such a way that it is easy to identify the type of attack and its solution to limit the attackers from damaging the IoT network.

### IV. DISCUSSION AND OPEN CHALLENGES

By the year 2025 hundred millions gadgets are to be associated in the IoT. So the security of the system ought to be the most critical issue in up and coming days. The security is turning into a test on the grounds that there is no standard design and security techniques actualized for one engineering won't not be achievable for another, therefore likelihood of securityattacks increments. Along these lines, the requirement for a standard design for IoT is obligatory. Hubs in an IoT organize are not sufficiently competent to deal with complex security calculations like Cryptography and so forth henceforth there is a solid requirement for a few calculations that can be executed on these low-preparing gadgets

### CONCLUSION

IoT has been a hot research theme throughout the previous couple of years and like other progressive advances, it likewise faces many difficulties, most huge of which are the security and protection dangers. In this paper, we depicted the working of four layers of IoT (Perception Layer, Network Layer, Processing Layer and Application Layer) and afterward we investigated the security escape clauses that can be misused in these layers. Besides, we clarified the countermeasures that can be embraced to avoid and secure the IoT arrange from the security dangers. Besides, we likewise recommended a few upgrades in the IoT system to influence it more to secure and to defeat the organization issues. As IoT will be a fundamental piece of our day by day lives sooner rather than later, uncommon advances must be taken to guarantee that clients trust and protection.

### References

- [1] F. Xia, L. T. Yang, L. Wang, and A. Vinel, "Internet of Things," *Int. J. Commun. Syst.*, pp. 1101–1102, 2012.
- [2] L. Coetzee and J. Eksteen, "The Internet of Things – Promise for the Future ? An Introduction," in *IST-Africa*, 2011, pp. 1–9.
- [3] I. Andrea, C. Chrysostomou, and G. Hadjichristofi, "Internet of Things : Security Vulnerabilities and Challenges," in *International Workshop on Smart City and Ubiquitous Computing Applications*, 2015, pp. 180–187.
- [4] K. Zhao and L. Ge, "A Survey on the Internet of Things Security," *Ninth Int. Conf. Comput. Intell. Secur.*, 2013.
- [5] D. E. Burgner, "Security of Wireless Sensor Networks," in *Eighth International Conference on Information Technology: New Generations*, 2011.
- [6] T. Halim, "A Study on the Security Issues in WSN," *Int. J. Comput. Appl.*, vol. 53, no. 1, p. 8887, 2012.
- [7] T. Bhattasali, "Sleep Deprivation Attack Detection in Wireless Sensor Network," *Found. Comput. Sci. New York, USA*, 2012.
- [8] M. Li and I. Koutsopoulos, "Optimal Jamming Attack Strategies and Network Defense Policies in Wireless Sensor Networks," *IEEE Trans. Mob. Comput.*, vol. 9, no. 8, 2010.
- [9] L. Li, "Study on Security Architecture in the Internet of Things," in *International Conference on Measurement, Information and Control (MIC) Study*, 2012, no. Mic, pp. 374–377.

- [10] N. Using and E. Curves, "A Secured Authentication Protocol for Wireless Sensor Networks Using Elliptic Curves Cryptography," *Sensors*, pp. 4767–4779, 2011.
- [11] M. Alizadeh, M. Salleh, M. Zamani, J. Shayan, and K. SASAN, "Security and Performance Evaluation of Lightweight Cryptographic Algorithms in RFID," *Recent Res. Commun. Comput.*, pp. 45–50, 2012.
- [12] G. Avoine, M. A. Bingo, X. Carpent, S. Berna, O. Yalcin, and S. Member, "Privacy-Friendly Authentication in RFID Systems : On Sublinear Protocols Based on Symmetric-Key Cryptography," *EEE Trans. Mob. Comput.*, vol. 12, no. 10, pp. 2037–2049, 2013.
- [13] D. Migault, D. Palomares, E. Herbert, W. You, G. Ganne, G. Arfaoui, and M. Laurent, "E2E : An Optimized IPsec Architecture for Secure And Fast Offload," in *Seventh International Conference on Availability, Reliability and Security E2E*, 2012.
- [14] E. Vasilomanolakis, J. Daubert, and M. Luthra, "On the Security and Privacy of Internet of Things Architectures and Systems," *darmstad Univ. J.*, 2015.
- [15] B. Y. Mo, T. H. Kim, K. Brancik, D. Dickinson, H. Lee, A. Perrig, and B. Sinopoli, "Cyber – Physical Security of a Smart Grid Infrastructure," *Proc. IEEE*, vol. 100, no. 1, pp. 195–209, 2012.
- [16] B. S. Thakur and S. Chaudhary, "Content Sniffing Attack Detection in Client and Server Side : A Survey," *Int. J. Adv. Comput. Res.*, no. 2, pp. 4–7, 2013.
- [17] S. Issues, "A Survey of RFID Deployment and Security Issues | Korea Science A Survey of RFID Deployment and Security Issues A Survey of RFID Deployment and Security Issues | Korea Science," *J. Inf. Process. Syst.*, vol. 7, no. 4, pp. 16–17, 2011.
- [18] R. Uttarkar and P. R. Kulkarni, "Internet of Things : Architecture and Security," *Int. J. Comput. Appl.*, vol. 3, no. 4, pp. 12–19, 2014.
- [19] V. Soni, P. Modi, and V. Chaudhri, "Detecting Sinkhole Attack in Wireless Sensor Network," *Int. J. Appl. or Innov. Eng. Manag.*, vol. 2, no. 2, pp. 29–32, 2013.
- [20] R. P. Padhy, "Cloud Computing : Security Issues and Research Challenges," vol. 1, no. 2, pp. 136–146, 2011.
- [21] W. Chen, R. K. Guha, T. J. Kwon, J. Lee, and Y. Hsu, "A survey and challenges in routing and data dissemination in vehicular ad hoc networks," *Wirel. Commun. Mob. Comput.*, no. October 2009, pp. 787–795, 2011.
- [22] F. Baccelli, A. El Gamal, and D. N. C. Tse, "Interference Networks With Point-to-Point Codes," *IEEE Trans. Inf. THEORY*, vol. 57, no. 5, pp. 2582–2596, 2011.
- [23] Z. Xu, Y. Yin, and J. Wang, "A Density-based Energy-efficient Clustering Algorithm for Wireless Sensor Networks," *Int. J. Futur. Gener. Commun. Netw.*, vol. 6, no. 1, pp. 75–86, 2013.
- [24] S. Sharmila, "Detection of sinkhole Attack in Wireless Sensor Networks using Message Digest Algorithms," *IEEE*, pp. 0–5, 2011.
- [25] S. Daneshmand, A. Jafamia-jahromi, A. Broumandan, and G. Lachapelle, "A Low-Complexity GPS Anti-Spoofing Method Using a Multi-Antenna Array," *ION GNSS12 Conf.*, pp. 1–11, 2012.
- [26] C. Chen, Y. Lin, Y. Lin, and H. Sun, "RCDA : Recoverable Concealed Data Aggregation for Data Integrity in Wireless Sensor Networks," *IEEE Trans. PARALLEL Distrib. Syst.*, vol. 23, no. 4, pp. 727–734, 2012.
- [27] A. Razzaq, K. Latif, H. F. Ahmad, A. Hur, Z. Anwar, and P. C. Bloodsworth, "Semantic security against web application attacks," *Inf. Sci. (Ny)*, vol. 254, pp. 19–38, 2014.
- [28] D. H. Patil, "Data Security over Cloud," *Int. J. Comput. Appl.*, pp. 11–14, 2012.
- [29] B. R. Chandramouli and P. Mell, "State of Security Readiness," *Crossroads*, vol. 16, no. 3, pp. 23–25, 2010.
- [30] K. Hashizume, D. G. Rosado, E. Fernández-medina, and E. B. Fernandez, "An analysis of security issues for cloud computing," pp. 1–13, 2013.
- [31] N. Kilari and C. Applications, "A Survey on Security Threats for Cloud Computing," *Int. J. Eng. Res. Technol.*, vol. 1, no. 7, pp. 1–10, 2012.
- [32] K. Dahbur, "A Survey of Risks , Threats and Vulnerabilities in Cloud Computing," in *International Conference on Intelligent Semantic Web-Services and Applications*, 2011.
- [33] Y. Singh, F. Kandah, and W. Zhang, "A Secured Cost-effective Multi-Cloud Storage in Cloud Computing," *IEEE INFOCOM*, pp. 619–624, 2011.
- [34] Z. Brakerski and V. Vaikuntanathan, "Efficient fully homomorphic encryption from (standard) LWE," *SIAM J. Comput.*, vol. 43.2, pp. 831–871, 2014.
- [35] B. L. Suto, "Analyzing the Accuracy and Time Costs of Web Application Security Scanners," *San Fr.*, no. October 2007, 2010.
- [36] S. Kumar, S. Pal, A. Kumar, and J. Ali, "Virtualization , The Great Thing and Issues in Cloud Computing," *Int. J. Curr. Eng. Technol.*, pp. 338–341, 2013.
- [37] D. Koo, J. Hur, and H. Yoon, "Secure and efficient data retrieval over encrypted data using attribute-based encryption in cloud storage q," *Comput. Electr. Eng.*, vol. 39, no. 1, pp. 34–46, 2013.
- [38] A. Tewari, A. K. Jain, and B. B. Gupta, "Recent survey of various defense mechanisms against phishing attacks," *J. Inf. Priv. Secur. ISSN*, vol. 6548, no. Feb, pp. 3–13, 2016.
- [39] J. Wan, N. Cn, and A. No, "Malware detection using pattern classification," 2012.
- [40] H. Tobias and E. Al., "Security Challenges in the IP-based Internet of Things," 2011.
- [41] M. C. M and A. Serbanati, "An overview of privacy and security issues in the internet of things," *Springer*, 2010.
- [42] A. Viejo, "Systems and methods for reducing unauthorized data recovery from solid-state storage devices," *Merry, Jr. al*, vol. 2, no. 12, p. Merry, Jr. et al, 2011.
- [43] Y. Shin, A. Meneely, L. Williams, and J. A. Osborne, "Evaluating Complexity , Code Churn , and Developer Activity Metrics as Indicators of Software Vulnerabilities," *IEEE Trans. Softw. Eng.*, vol. 37, no. 6, pp. 772–787, 2011.
- [44] W. Enck, D. Octeau, and P. McDaniel, "A Study of Android Application Security," *Syst. NTERNET NFRASTRUCTURE Secur.*, no. August, 2011.
- [45] D. William, W. Surrey, and J. F. Benedict, "Authentication using application authentication element," 2012.
- [46] M. Ongtang, S. McLaughlin, W. Enck, and P. McDaniel, "Semantically rich application-centric security in Android," *Secur. Commun. Networks*, no. August 2011, pp. 658–673, 2012.
- [47] S. L. Wiley, O. Park, and U. S. C, "Pin-hole firewall for communicating data packets on a packet network," 2011.
- [48] C. Liu and Y. Zhang, "Research on Dynamical Security Risk Assessment for the Internet of Things Inspired by Immunology," in *8th International Conference on Natural Computation*, 2012, no. Icnc, pp. 874–878.

## Life Time Enhancement for Wireless Sensor Network using Fuzzy-ACO Combination along with an Enhanced Clustering Scheme

**Mr. Arshad Shareef<sup>1</sup>**

Assistant Professor  
Department of CSE

Arkay College of Engg. & Tech.

E-Mail<sup>1</sup>: arshad.shareef323@gmail.com

**Mrs. Ursa Sayeed<sup>2</sup>**

Assistant Professor  
Department of CSE

Arkay College of Engg. & Tech.

E-Mail<sup>2</sup>: ursasayeed@gmail.com

**Dr. P A Abdul Saleem<sup>3</sup>**

Professor & Principal  
Department of CSE

Arkay College of Engg. & Tech.

E-Mail<sup>3</sup>: drsaleemprincipal@gmail.com

**Abstract** - The principle goal of this research is to build up an energy efficient routing technique which bears solid communication by considering routing challenges, and broaden the lifetime of WSNs. This paper has two primary commitments. That is it proposes an altered Energy Efficient Cluster Formation (EECF) technique in which the Cluster Heads (CHs) are chosen in light of parameters like up-keeping energy, distance to the neighbors, density, maximum distance and angle. And the center point of this system is to apportion overwhelming data traffic and high energy consumption load reliably in the network by offering unequal size of clusters in the network. In this, the K-means algorithm is used for the cluster member selection and the cluster head selection process. Then, a Super Cluster Head (SCH) is chosen among the CHs who can just send the information to the mobile BS by picking appropriate fuzzy descriptors. Along with an Ant Colony Optimization (ACO) is used for

the optimal routing selection from the SCH to the Base Station (BS).

**Keywords:** Wireless Sensor Network (WSN), K-means clustering, Cluster head (CH) and Super Cluster Head (SCH), fuzzy descriptors, Ant Colony Optimization (ACO).

### 1. Introduction

Wireless Technologies have shown an alternate dimension to the world of communication Methods. Actually it began with the utilization of radio receivers or transceivers for use in wireless telegraphy during an early stage, and now the term wireless is utilized to depict technologies, for example, the cellular networks and Wireless Broadband Internet etc. In spite of the fact that the wireless medium has restricted spectrum with different imperatives when compared with the guided medium. But wireless medium is the main channel for mobile communication. Remote AdHoc networking is an infrastructure technology which is utilized for random and

quick arrangement of countless networks. Also this technology has utilizations in many fields. For example, tactical communications, disaster relief operations, health care and transitory networking in regions that are not thickly populated.

### **1.1 Wireless Sensor Networks (WSN)**

Wireless Sensor Networks (WSN) are in some cases called Wireless Sensor and Actuator Networks (WSAN) [1]. WSN is spatially distributed autonomous sensors (NODES) in order to screen the physical or natural conditions, and some normal parameters are temperature, sound, pressure, and so on. The improvement of wireless sensor networks was inspired by military and applications, such as, war zone Nuclear Reactor area respectively; today these networks are utilized as a part of numerous mechanical and technological applications, like, modern process checking and control, machine wellbeing observing etc.

The primary characteristics of a WSN include:

- Power consumption requirements for nodes utilizing batteries or energy harvesting
- Ability to adapt to node disappointments (resilience)

- Some mobility of nodes (for very mobile nodes see MWSNs)
- Heterogeneity of nodes
- Scalability to large scale of deployment
- Ability to withstand harsh environmental conditions
- Easy utilization
- Adapting Cross-layer design

The principle use of WSN incorporates Healthcare monitoring, Area monitoring, Environmental/Earth sensing, Air contamination monitoring, Forest fire identification, Landslide recognition, Water quality monitoring, Natural fiasco avoidance, Industrial monitoring, Machine wellbeing monitoring, data logging, Water/Wastewater monitoring, Structural wellbeing monitoring and so on.

Conventional wireless communication networks like Mobile AdHoc Networks (MANET) varies from WSN [2] [3]. WSN has special attributes, for example, the denser level of node deployment, higher unreliability of sensor nodes and extreme energy requirements for computation and capacity requirements which introduce many difficulties in the improvement and utilization of WSN. Research has been done to investigate and discover answers for

different plan engineering and application issues and noteworthy progress has been made in the advancement and the organization of WSNs. WSN regularly contains hundreds or thousands of sensor nodes which takes into consideration sensing over bigger geographical districts with more prominent accuracy [4]. As a rule, the sensor nodes are sent randomly finished the geographical area and these nodes speak with each other to shape a network. Every hub has three essential components: (1) sensing unit, (2) processing unit and (3) transmission unit.

## 1.2 Efficiency in WSN

Late improvements in low-control wireless integrated microsensor technologies have made these sensor nodes accessible in expansive numbers, requiring little to no effort, to be utilized in an extensive variety of uses in military and national security, ecological checking, and numerous different fields [5]. As opposed to traditional sensors, sensor networks offer an adaptable recommendation as far as the simplicity of deployment and various functionalities. In classical sensors, the position of the nodes and the network topology should be foreordained and painstakingly built [6]. Be that as it may, on account of modern

wireless sensor nodes, their smaller physical dimensions allow a huge number of sensor nodes to be randomly sent in difficult to reach landscapes. Furthermore, the nodes in a wireless sensor network are additionally equipped for performing different operations similar to MANET which are data processing and routing. But in traditional sensor networks uncommon nodes with computational abilities must be introduced independently to accomplish such functionalities.

In the immediate transmission protocol, the base station fills in as the destination node to the various nodes in the network where the end user can get to the sensed data [7]. At the point when a sensor node transmits information specifically to the base station, the energy misfortune brought about can be very broad relying upon the area of the sensor nodes with respect to the base station. In such a situation, the nodes that are further far from the base station will have their power sources depleted significantly speedier than those nodes that are nearer to the base station [8]-[10]. Then again, using a conventional multihop routing plan, for example, the Minimum Transmission Energy (MTE) routing protocol will likewise bring about a similarly unwanted impact. In MTE, the Nodes nearest to the

base station will quickly deplete their energy resources since these nodes take part in the routing of an extensive number of information messages (in the interest of different nodes) to the base station [1, 2]. Different routing protocols have been proposed for wireless sensor networks to ease such issues.

In this paper, at first a K-means clustering algorithm uses for the cluster member determination process, the CH among the cluster member is chosen in light of the Euclidean distance between the nodes. The SCH is chosen by using the fuzzy descriptors (Mamdani's rule). At last, the optimal way between the SCH and the base station is assessed by utilizing the ACO algorithm. Here, the lifetime of the WSN is improved in view of the packet delivery, throughput, and the energy consumption.

The rest of the segment of the work is delineated in the area underneath. Area 2 clarified the literature review, the proposed EECF is portrayed in segment 3, the results and the conclusion are delineated in segment 4 and 5.

## 2. Literature Review

A stable and energy-efficient clustering (SEEC) protocol for heterogeneous wsns is proposed by Farouk et al. [11]. What's more,

the expansion to multi-level of SEEC is introduced. It relies upon the network structure that is partitioned into clusters. In the multi-level architectures, all the more effective super nodes are relegated to covering removed sensing territories. At each level of heterogeneity, the optimum number of capable nodes that accomplishes the base energy utilization of the network is acquired. Reenactment results demonstrate that the proposed protocol gives a more drawn out stability period, more energy productivity and higher normal throughput than the current protocols.

Chandet al. execution of HEED for a heterogeneous network [12]. Contingent on the kind of nodes, it characterizes one-level, two-level, and three-level heterogeneity and as needs be the usage of HEED is alluded to as hetHEED-1, hetHEED-2, and hetHEED-3, individually. One extra parameter is considered in this work, i.e., remove and apply fuzzy logic to decide the cluster heads and in like manner the hetHEED-1, hetHEED-2, and hetHEED-3 are named as HEEDFL, hetHEED-FL-2, hetHEED-FL-3, individually. The simulation results demonstrate that as the level of heterogeneity increments in the network, the nodes stay alive for longer time and the rate of energy scattering diminishes. And

furthermore, expanding the heterogeneity level causes sending more packets to the base station and builds the network lifetime.

An improved harmony search based energy efficient routing algorithm (IHSBEER) for WSNs is exhibited by Zeng and Dong [13], which depends on harmony search (HS) algorithm (a meta-heuristic). To address the WSNs routing issue with HS algorithm, a few key enhancements have been advanced: First of all, the encoding of harmony memory has been enhanced in light of the qualities of routing in WSNs. Furthermore; the ad lib of another harmony has likewise been making strides. A dynamic adjustment for the parameter HMCR is acquainted with stay away from the rashness in early ages and fortify its nearby pursuit capacity in later ages. In the interim, the change procedure of HS algorithm has been disposed of to make the proposed ROUTING algorithm containing less parameter. Thirdly, a successful local search strategy is proposed to upgrade the local search capacity, in order to enhance the merging pace and the ACCURACY of routing algorithm.

A powerful multi-level cluster algorithm utilizing link correlation is proposed for heterogeneous WSN was created by Javaid

et al. [14]. The level-k hierarchy with single-hop communication between nodes inside a cluster is accomplished utilizing link correlation. The heterogeneous nodes are embraced as level-k cluster heads and actualizing network coding on those NODES builds network lifetime altogether. In the mean time, actualizing time division multiple access (TDMA) system inside a cluster makes a sorted out cluster design enhancing the energy productivity.

Rejina Parvin and Vasanthanayaki proposed E-OEERP [15] which lessens/takes out individual node formation and enhances the general network lifetime when contrasted with the current protocols. It can be accomplished by applying the ideas of Particle Swarm Optimization (PSO) and Gravitational Search Algorithm (GSA) for cluster formation and routing individually. For each cluster head, a strong node called cluster assistant (CA) node is chosen to lessen the overhead of the cluster head. With the assistance of PSO, clustering is performed until the point when every one of the nodes turns into a member of any of the cluster. This disposes of the individual node formation which brings about a similarly better network lifetime.

Zahedi et al.[16] proposed a Swarm Intelligence Based Fuzzy Routing Protocol (named SIF). In SIF, fuzzy C-means clustering algorithm has used to cluster all sensor nodes into adjusted clusters, and after that fitting cluster heads are chosen through Mamdani fuzzy inference system. This methodology not just certifications to produce adjusted clusters over the network, yet additionally can decide the exact number of clusters. Since tuning the fuzzy rules extremely influences on the execution of the fuzzy system, a half breed swarm intelligence algorithm in light of firefly algorithm is used and recreated toughening to advance the fuzzy rule base table of SIF.

An approach of choosing Super Cluster Head (SCH) the CHs was proposed by Nayak and Devulapalli [17]. The SCH can just send the data to the mobile BS by picking appropriate fuzzy descriptors, for example, remaining battery power (RBP), mobility of BS and centrality of the clusters. Fluffy inference engine (Mamdani's rule) is utilized to choose the opportunity to be the SCH. The outcomes have been gotten from NS-2 simulator and demonstrates that the proposed protocol performs superior to LEACH protocol regarding First node bites the dust, half node alive, better stability and better lifetime.

An energy efficient cluster formation algorithm called active node cluster formation (ANCF) was proposed by Faheem et al. [18]. The center expect to propose ANCF algorithm is to disseminate substantial data traffic and high energy consumption load equitably in the network by offering unequal size of clusters in the network. A lightweight sensing mechanism called active node sensing algorithm (ANSA) is additionally introduced. The key intend to propose the ANSA algorithm is to maintain a strategic distance from high sensing covering data repetition by selecting an arrangement of active nodes in each cluster with fulfill scope close to the occasion. An active node routing algorithm (ANRA) is additionally proposed to address complex inter and intra cluster routing issues in profoundly thick deployment in view of the node ruling values.

### **3. Proposed Methodology**

The center point this system is to administer overwhelming data traffic and high energy consumption load reliably in the network by offering unequal size of clusters in the network. The created strategy settles each CH close to the sink and sensing occasion while the staying set of the CHs is designated amidst each cluster to

accomplish the largest amount of ENERGY proficiency in dense deployment. Second, a super cluster head (SCH) is chosen among the CHs who can just send the information to the mobile BS by picking reasonable fuzzy descriptors, for example, remaining battery power, mobility of BS and the centrality of the clusters [18]. We likewise think of one as more parameter, i.e., the energy dissipation ratio. Fluffy inference engine (Mamdani's rule) is utilized to choose the opportunity to be the SCH. We apply Ant Colony Optimization (ACO) for optimizing the routing in the proposed approach. By focusing on inalienable properties of routing, it will be reasonable to be settled by the ant colony algorithm. This determination of a path is viewed as optimal among the diverse paths. The possibility of energy efficient cluster formation alongside SCH logic builds the network lifetime significantly.

### 3.2. Cluster Member Selection based on K-Means Clustering

In the proposed conspire, the proficient clustering of nodes with minimum repetitive (un-clustered) node in the network is done in light of the k-means clustering algorithm. Here, the object's clusters are shaped by the Euclidean distance between the object. The

K-means algorithm merges three stages to assess the cluster head.

#### (a) Initial Clustering

For cluster formation with the target WSN, the K-means algorithm is utilized. Presume, the number of nodes in the WSN is represented as  $n$  and the Euclidean distance between the nodes are assessed to isolate the  $k$  clusters. At first, the cluster heads (CHs) are randomly selected from the  $k$  out of  $n$  nodes. As per the Euclidean distance, the rest of the nodes in the network are considered as nearest cluster of CH.

#### (b) Reclustering

The CENTROID of the every CLUSTER is assessed, after the CLUSTER FORMATION for each NODE in the NETWORK. Consider the TWO DIMENSIONAL SPACE, every CLUSTER CENTROIDS of every NODE  $S$  are computed as takes after,

$$Centroid(U, V) = \left( \frac{1}{S} \sum_{i=1}^n u_i, \frac{1}{S} \sum_{i=1}^n v_i \right)$$

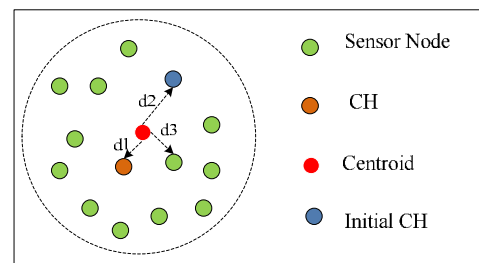


Figure.1: A cluster with the centroid

The center position of the cluster is meant as the centroid and it is known as a virtual node. The illustration, clustering process is portrayed in the Figure1, which has 15 nodes, here, the sensor node is distinguished by arbitrarily chosen CH before all else round. The nearest cluster of the centroid is considered as the new CH and the at first financed CH isn't a nearest to the centroid. The procedure is rehashed until there is no progressions may happen in the CH choice.

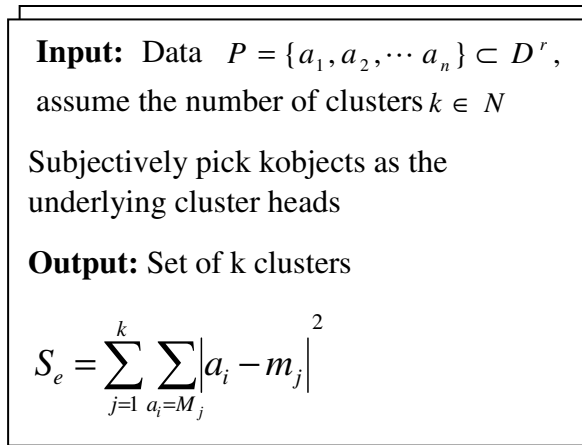


Figure.2: K-means Clustering Algorithm

### 3.3. Cluster Head Selection

After the cluster formation, an altered energy efficient CH determination system relies upon the parameter like residual energy, distance to the neighbor, density, maximum distance and angle. In the WSN network, the overwhelming data traffic and the and the high energy consumption are lessened by utilizing these strategies. According to the separation from the

centroid, one ID number is doled out to every node in the cluster and which node has less ID number it is considered as the nearest one. In the CH node determination, the ID number assumes an imperative part. The CH are chosen in view of the accompanying advances,

- **Residual energy ( $RE_n$ )**

To hold the availability of the network, the CH's lingering energy is checked each round. It is used for the packet transmission.

$$CH_{RE_n} > Threshold$$

- **Maximum distance**

In this progression, the greatest distance between the two nodes are estimated.

- **Angle**

**In this step, the angle between the two sensor nodes are estimated.**

- **Density ( $\rho$ )**

In the given locale, the NUMBER OF NODES is assessed by utilizing the DENSITY FUNCTION. In the thick district the value of  $\rho$  being one and the thin locale the value of  $\rho$  is 2.

- **Distance to the Neighbor**

Here, the distance to the neighbor node of the CH is computed as follows,

$$AN = \frac{RE_n}{\left( \sum_{i=1}^n d_{gh} / d_{gh(max)} \right)^2 + \left( \rho - \left( \frac{d_i}{100} \right) \right)^2} \forall SN \quad (2)$$

In the above equation, AN is the active node, SN is the sensor node,  $RE_n$  is the residual energy and  $\rho$  is the density function.

### 3.4. Mamdani's fuzzy approach for SCH Selection

The SCH choice process is critical to expand the lifetime of the cluster and cluster head. So we have used a Mamdani's fuzzy approach for choosing super cluster head in the considered sensing region. By utilizing this procedure, the SCH are straightforwardly sent to the BS after the SCH are chosen among the already chose CH. There are four models incorporated into the fuzzy descriptors; they are (1) fuzzification (2) rule evaluation (3) aggregation of the rule outputs (4) defuzzification.

#### • Fuzzification

In the fuzzification procedure, the crisp values are going about as the input and it is changed over into fuzzy set.

#### • Rule Evaluation

In the rule evaluation process, the rashly changed over fuzzy input is given to the

before the fuzzy rule and afterward it is sent to the ensuing membership function.

#### • Aggregation of the administer yields

Each rule's combining yield is incorporated into this procedure.

#### • Defuzzification

The defuzzification is the reverse of the fuzzification procedure. Here, the fuzzy sets are changed over into the crisp values.

The Mamdani's fuzzy approach used the remaining battery power, the mobility of bs, the centrality of the cluster, and energy dissipation ratio for the SCH choice process. In the SCH choice process, the accompanying formula is used to infer the fuzzy rule,

$$SCH = RBP + Mobility + Centrality + EDR \quad (3)$$

In the above equation, the remaining battery power (RBP) of each node is computed by using the following equation,

$$RBP = (BP - 1) \quad (4)$$

In the above equation,  $RBP$  is the remaining battery power,  $EDR$  is the energy dissipation ratio and  $BP$  is the battery power. At the season of SCH choice, the energy and the RBP of every node is diminished by the information drive process. Here, the mobility and the centrality are

going about as the additive factors, since the detachment of SCH from the BS augmentations or diminishes as the BS moves.

### 3.5. Ant Colony Optimization (ACO) approach for path selection

After the SCH determination process, the optimum path between the BS and the SCH nodes are chosen by using the ACO algorithm. Ant Colony Optimization [19] is a standard answer to finding optimal paths (from source to destination). ACO depends on the intellectual foundation that can without much of a stretch be portrayed in one sentence: ants select the best path among the current barriers and requirements in nature to accomplish food. This choice of a path is viewed as optimal among the distinctive paths. The possibility of energy efficient cluster formation alongside SCH rationale builds the network lifetime drastically.

In the ACO algorithm, the probability with which ant  $A$  in sensor  $S$ , selects to move to the sensor  $R$  is given viz,

$$R = \begin{cases} \arg \max_{v \in J_A(S)} \{ [\lambda(S, R)] \cdot [\beta(S, R)]^\mu \} & \text{if } p \leq p_0 \\ r, & \text{otherwise} \end{cases} \quad (5)$$

In the above equation, the random number is represented as  $p$  which is distributed within the range  $(0.....1)$ , the parameter is

represented as  $p_0$  which ranges from 0 to 1 (i.e.,  $0 \leq p_0 \leq 1$ ), the pheromone is represented as  $\lambda$ , the inverse of the distance  $\chi(S, R)$  is represented as  $\beta = 1/\chi$ , the set of node is represented as  $J_A(S)$  and it is visited by ant  $A$ . The relative importance of the pheromone versus distance ( $\mu > 0$ ) is determined by the parameter  $\mu$  and the positioned on the node  $S$ .

$$Q_A(S, R) = \begin{cases} \frac{[\lambda(S, R)] \cdot [\beta(S, R)]^\mu}{\sum_{v \in J_A(S)} [\lambda(S, R)] \cdot [\beta(S, R)]^\mu} & \text{if } R \in J_A(S) \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

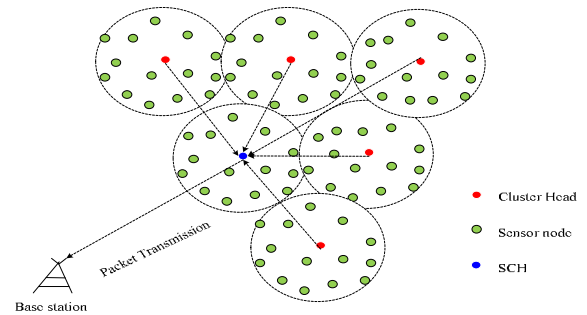


Figure.3: Packet Transmission Using EERRCUF Protocol

After every ant has completed its tour, the global updating is done. The pheromone was deposited by the every globally best ant. The global updating rule is given in the

$$\lambda(S, R) \leftarrow (1 - \delta) \cdot \lambda(S, R) + \delta \cdot \Delta \lambda(S, R) \quad (7)$$

In the above equation,

$$\Delta\lambda(S,R)=\begin{cases} (G_{in})^{-1}, & \text{if } (S,R) \in \text{global best tour} \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

In the above equation, the length of the globally best tour from the beginning of the trail is represented as  $G_{in}$  and the pheromone decay parameter is represented as  $\delta$ .

In the process of finding the solution, when an ant visits a particular edge, the pheromone level changes according to the local updating rule, which given by

$$\lambda(S,R) \leftarrow (1 - \vartheta) \cdot \lambda(S,R) + \vartheta \cdot \Delta\lambda(S,R) \quad (9)$$

Where  $0 < S < 1$  is a parameter and the initial pheromone level being represented as  $\Delta\lambda$ . Finally the best path is detected by using this ACO algorithm.

#### 4. Simulation Results

The validity of the proposed protocol is checked by using the NS-2 simulator. The execution of the network in light of the lifetime improvement of the system. For expanding the lifetime of a WSN, the proposed energy efficient cluster formation (EECF) method joined by fuzzy logic utilizing IEEE 802.11 MAC layer.

#### 4.1. Performance Analysis

For the execution calculation, the proposed EECF strategy is contrasted and the current FLEACH (Fuzzy LEACH) and the LEACH strategies. The execution is evaluated in light of the packet delivery ratio, energy consumption, end to end delay, network lifetime and throughput.

##### (a) Packet Delivery Ratio

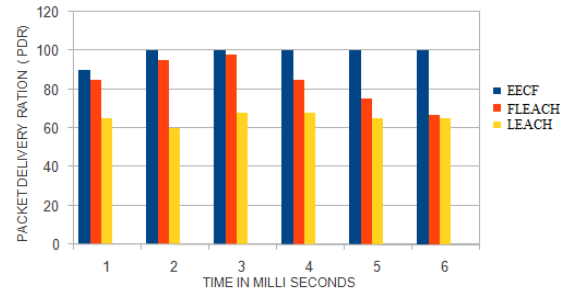


Figure.4 : Packet Delivery Ratio

Figure 4 demonstrates the packet delivery ratio of the proposed EECF technique and it contrasted and the FLEACH and LEACH strategies. The number of sending and got Packet's ratio is known as the packet delivery ratio. At the point when the packet delivery ratio was most extreme VALUE around then just the network plays a capable network. The reliability and trustworthiness of the WSN network are expanded, if the packet delivery is high. The proposed EECF procedure has better packet delivery ratio when contrasted and the other two methods.

### (b) Energy Consumption

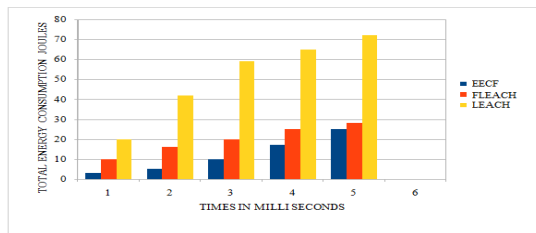


Figure.5 : Total Energy Consumption

Energy consumption is the total amount of energy required to transmit the data packet from the node to the base station. Here, the energy consumption of the proposed work estimated by the measurement of the various WSN time slots is compared with the existing techniques. The energy consumption of the proposed EERRCUF method is depicted in the Figure 5. When compared with the existing FLEACH and LEACH technique, the proposed EERRCUF technique required less energy to transmit the packet. Due to this reason, the lifetime of the proposed EERRCUF technique is enhanced.

### (c) Network Lifetime

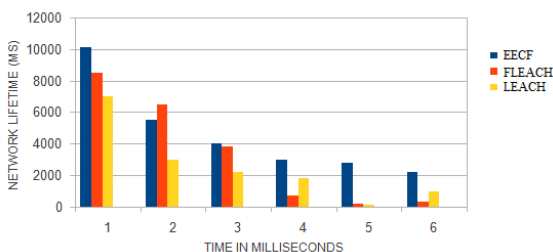


Figure.6: Network Lifetime

The network lifetime of the proposed EERRCUF technique is shown in the Figure 6 and it is compared with the existing techniques. Here, the lifetime of the WSN network is enhanced by utilizing the EERRCUF protocol. The proposed method enhances the WSN lifetime proficiently when compared with the other techniques.

### (d) Throughput

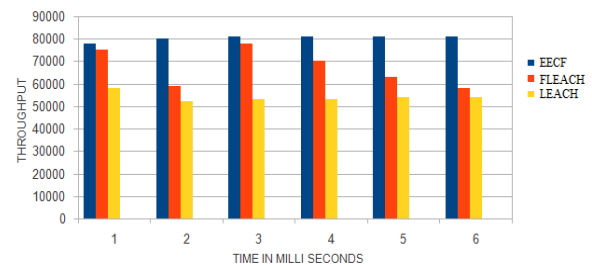


Figure.7: Throughput

The throughput of the proposed EERRCUF protocol is depicted in the Figure 7 and it is compared with the existing techniques. The efficiency of the data transmits from the node to the base station is known as the throughput. Here, the packet is delivered in bit per second and compared with the existing FLEACH and LEACH protocol the proposed EERRCUF protocol has high throughput.

## 5. Conclusion

This paper proposes an EECF protocol for upgrading the lifetime of the WSN. The proposed routing protocol executed utilizing

an ACO algorithm and which assesses the optimum paths from the every accessible Path. The proposed protocol is contrasted and the current FLEACH and the LEACH protocol to assess the execution of the network. The trial result demonstrated that the proposed EECF protocol worked proficiently to upgrade the lifetime of the WSN in view of the packet delivery ratio, energy consumption, throughput.

## References

- [1] Singh, Tejpreet, Jaswinder Singh, and Sandeep Sharma. "Energy Efficient Secured Routing Protocol For Manets". Wireless Netw (2016).
- [2]Sohrabi, K., Gao, J., Ailawadhi, V., Pottie, G., "Protocols for Self-Organization of a Wireless Sensor Network," IEEE Personal Communications Mag., Vol.7, No.5, pp.16-27, Oct. 2000.
- [3]. F. Akyildiz et al, (March 2002), "Wireless sensor networks: a survey", Computer Networks, V Vol. 38, pp. 393-422.
- [4]R. Min, et al., (January 2001), "Low Power Wireless Sensor Networks," International Conference o on VLSI Design, Bangalore, India
- [5]Shio Kumar Singh, M P Singh and D K Singh, (August 2010), "A Survey of Energy-Efficient Hierarchical Cluster-Based Routing in wireless Sensor Network," Int J. of Advanced Networking and Applications Volume: 02, Issue: 02, Pages: 570-580.
- [6]R Ramanathan, R Hain, "Topology Control of Multihop Wireless Networks Using Transmit Power Adjustment", Proceeding Infocom 2000.
- [7]S. Lindsey and C. Raghavendra, "PEGASIS: Power-Efficient gathering in Sensor Information Systems", International Conference on Communications, 2001.
- [8]Rabiner, W., Chandrakasan, A., Balakrishnan, H., "Energy-Efficient Communication Protocol for Wireless Microsensor Networks," Hawaii International Conference on System Sciences, Maui, HI, pp.10-19, Jan. 2000
- [9]I. F. Akyildiz et al., "Wireless Sensor Networks: A Survey," Elsevier Sci. B. V. Comp. Networks, vol. 38, no. 4, Mar. 2002, pp. 393–422.
- [10]W. R. Heinzelman, A. P. Chandrakasan, and H. Balakrishnan, "EnergyEfficient Communication Protocol for Wireless Microsensor Networks," Proc. 33rd Hawaii Int'l. Conf. Sys. Sci., Jan. 2000.

- [11]F. Farouk, F. Zaki and R. Rizk, "Multi-level stable and energy-efficient clustering protocol in heterogeneous wireless sensor networks", IET Wireless Sensor Systems, vol. 4, no. 4, pp. 159-169, 2014.
- [12]S. Chand, S. Singh and B. Kumar, "Heterogeneous HEED Protocol for Wireless Sensor Networks", Wireless Pers Commun, vol. 77, no. 3, pp. 2117-2139, 2014.
- [13]B. Zeng and Y. Dong, "An improved harmony search based energy-efficient routing algorithm for wireless sensor networks", Applied Soft Computing, vol. 41, pp. 135-147, 2016.
- [14]N. Javaid, M. Rasheed, M. Imran, M. Guizani, Z. Khan, T. Alghamdi and M. Ilahi, "An energy-efficient distributed clustering algorithm for heterogeneous WSNs", J Wireless Com Network, vol. 2015, no. 1, 2015.
- [15]J. RejinaParvin and C. Vasanthanayaki, "Particle Swarm Optimization-Based Clustering by Preventing Residual Nodes in Wireless Sensor Networks", IEEE Sensors J., vol. 15, no. 8, pp. 4264-4274, 2015.
- [16]Z. Zahedi, R. Akbari, M. Shokouhifar, F. Safaei and A. Jalali, "Swarm intelligence based fuzzy routing protocol for clustered wireless sensor networks", Expert Systems with Applications, vol. 55, pp. 313-328, 2016.
- [17]P. Nayak and A. Devulapalli, "A Fuzzy Logic-Based Clustering Algorithm for WSN to Extend the Network Lifetime", IEEE Sensors J., vol. 16, no. 1, pp. 137-144, 2016.
- [18]M. Faheem, M. Abbas, G. Tuna and V. Gungor, "EDHRP: Energy efficient event driven hybrid routing protocol for densely deployed wireless sensor networks", Journal of Network and Computer Applications, vol. 58, pp. 309-326, 2015.
- [19] A. Colorni, M. Dorigo, and V. Maniezzo, "Distributed optimization by ant colonies," in Proceedings of the 1st European Conference on Artificial Life, pp. 134–142, 1991.

# Evaluation Based Load Balancing for Cloud Computing

Dr.Sunil Tekale  
Professor  
Department of CSE  
MRCE  
Hyderabad, Telangana  
[sunil.tekale2010@gmail.com](mailto:sunil.tekale2010@gmail.com)

Mr.M.Amarnath  
Assistant Professor  
Department of CSE  
MRCE  
Hyderabad, Telangana  
[amaranth.cse@mrce.in](mailto:amaranth.cse@mrce.in)

## Abstract:

Cloud computing is the best technology today for all those people who want to go with minimum investment on infrastructure and want to outsource the burden of handling technical issues to third party by paying the charges for the services utilized.

Today there is huge amount of demand from the clients to make use of cloud technology as it provides multiple features and take off the load of maintaining infrastructure. This has created a huge amount of load on servers. So it is must to handle issues related to load balancing. This is basically to see that the load on a particular server is kept maximum to its threshold level. So that it can handle the task and also can complete it in a faster manner. It minimizes the cost and time involved in the major computational models and helps to improve proper utilization of resources and system performance. Many algorithms are recommended by various researchers from all over the world to solve the problem of load balancing.

In this paper, we present a new algorithm named as combo algorithm to address the issue of load balancing in a cloud environment.

**Keywords** - Cloud Computing optimization Load Balancing Network

## 1. Introduction

Cloud computing is a newly progressing technique which offers online computing resources, storage and permits users to organize applications with enhanced scalability, availability and fault tolerance. Cloud computing is about storing the stuff on remote servers instead of on own computers or other devices[1].

This information can be retrieved using the internet with any device, everywhere in the world as long as that device can support cloud computing systems. The cloud computing system is comprised of a front-end, which is the client side and a back-end which is a collection of the servers and computers owned by a third party which stores the data. A central server which is a fragment of the back-end follows protocols and uses middleware to communicate between networked computers. Cloud computing accumulates all the computing resources and manages them automatically. Its characteristics[4]

describe a cloud computing system: on-need self-service, pooling of resources, access to the internet, the elasticity of service availability and measurement of services utilized by individual users[6]. Cloud computing is everywhere with tools like Google Drives replacing Microsoft Office, Amazon Web Services replacing traditional enterprise data storage, banking websites replacing branch offices and Dropbox storing all our data and files. The cloud even provides different deployment models and service models[8].

The four deployment models present in cloud computing are:[2]

**1. Public cloud:** In the public cloud, the cloud provider provides resources for free to the public. Any user can make use of the resources; it is unrestricted. The public cloud is connected to the public internet for anyone to leverage.

**2. Private cloud:** In a private cloud, the planning and provisioning of the cloud are operated and owned by the organization or the third party. Here the hosted services are provided to a restricted number of people or group of individuals.

**3. Community cloud:** These type of cloud infrastructures exists for special use by a group of users. These are a group of users who share a common mission or have specific regulatory requirements, and it may be managed by the third party or organizations.

**4. Hybrid Cloud:** Hybrid Cloud provides the best of above worlds. It is created by combining the benefit of different types of cloud (private cloud & public cloud). In these clouds, some of the resources are provided and managed by public cloud and others as a private cloud[7].

The three different service models present in cloud computing are[8]:

**1.Infrastructure as a Service (IaaS):** IaaS model provides just the hardware and the network. It allows users to develop and install their operating system, software and run any application as per their needs on cloud hardware of their own choice.

**1. Platform as a Service (PaaS):** In PaaS model, an operating system, hardware, and network are provided to the user. It enables users to build their applications on cloud making use of supplier specific tools and languages

**2. Software as a Service (SaaS):** In SaaS model, a pre-built application together with any needed software, hardware, operating system and the network is provided to the user.

## 2. Load Balancing

Load balancing is a serious concern in cloud computing. With the increase in attractiveness of cloud computing among users, the load on the servers and the quantity of processing done is surging drastically. There are multiple nodes in the cloud, and due to the random allocation of a request made by the client to any node, the nodes become unevenly loaded[10]. So to avoid the condition where some nodes are either severely loaded or under loaded, the load balancer will evenly divide the workload among all the nodes <sup>3</sup>. Thus load balancing will equally distribute the workload among the nodes, and it can help in minimizing delays in communication, maximizing the throughput, minimizing execution time and maximizing resource utilization [3].

### 2.1 Goals of load balancing:

Some of the key purposes of a load balancing algorithm as pointed out by are:

1. It should possess fault tolerance.
2. It should be capable of modifying itself according to any change or expansion in the distributed system configuration <sup>3</sup>.
3. Regarding system performance, it should give greater overall improvement at a minimal cost.
- 4.Regardless of the origin of job it must treat all jobs in the system equally.
4. It should also maintain system stability.

### 2.2 Issues of Load Balancing

The issues of load balancing are described below[4]:

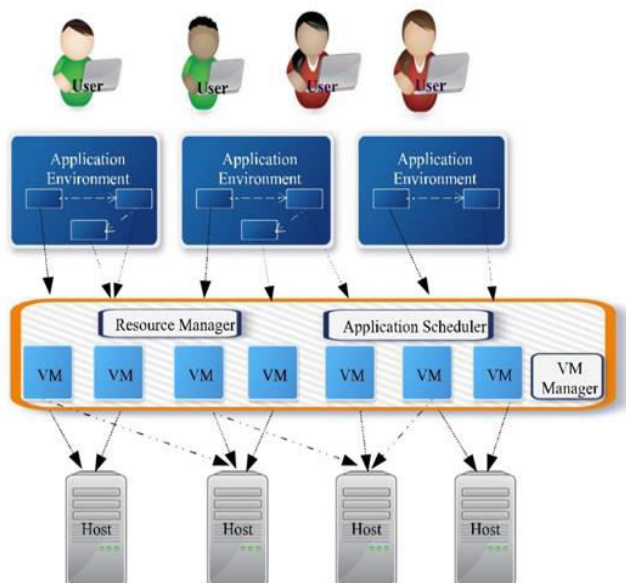
1. Load balancing becomes critical because, in the middle of execution, the processes may shift amongst nodes to ensure equal workload on the system <sup>5</sup>.
2. For a load balancing scheme to be good it should be scalable, general and stable and should add minimal overhead to the system. These requirements are interdependent <sup>6</sup>.
3. One of the critical aspects of the scheduling problem is load balancing <sup>7</sup>. The challenge for a scheduling algorithm is to avoid the conflict between prerequisites: fairness and data locality.
4. Algorithms for load balancing have to be dependent on the hypothesis that the on hand information at each node is accurate to avoid processes from being continuously circulated the system without any progress<sup>5</sup>.
5. How to accomplish a balance in load distribution amongst processors such that the computation can be done in the minimum possible time is one of the important problems to resolve.

6. Load balancing and task scheduling in distributed operating systems is a vital factor in gross system efficiency because the distributed system is not pre-emptive and non-uniform, that is, the processors may be different <sup>7</sup>.

### 2.3 Components of Load Balancing Algorithms[1]:

A load balancing algorithm has five major components <sup>8</sup>

1. Transfer Policy: The portion of the load balancing algorithm that picks a job for moving from a local node to a remote node is stated as Transfer policy or Transfer strategy.
  2. Selection Policy: In this policy, it specifies the processors involved in the load exchange (processor matching) so that the overall response time and throughput may be improved.
  3. Location Policy: The portion of the load balancing algorithm that is responsible for choosing a destination node for a task to transfer is stated as location policy or Location strategy[2].
  4. Information Policy: The part of the dynamic load balancing algorithm that is in charge of gathering information about the nodes present in the system is started to as Information Policy or Information strategy.
  5. Load Estimation Policy: In this policy, it determines the total workload of a node in a system.
- Fig:1 shows the VM, Application, Host relationship in DC



### Benefits of Cloud Load Balancing

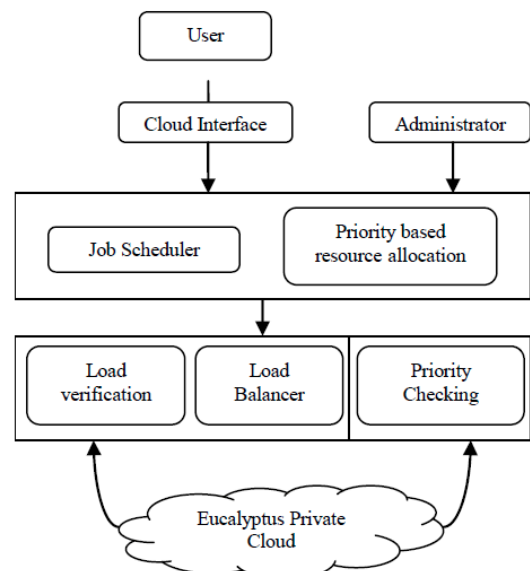
The benefits of cloud load balancing in particular arise from the scalable and global character of the cloud itself. The ease and speed of scaling in the cloud means that companies can handle traffic spikes (like those on Cyber Monday) without degraded performance by placing a cloud load balancer in front of a group of application instances, which can quickly autoscale in reaction to the level of demand. The ability to host an application at multiple cloud hubs around the world can boost reliability.

If a power outage hits the northeastern U.S. after a snowstorm, for example, the cloud load balancer can direct traffic away from cloud resources hosted there to resources hosted in other parts of the country.

### 4. Proposed Algorithm.

The Evaluation based load balancing can be done in a much more easier and simplest manner with high efficiency. In this we have to just calculate the capacity of each server in terms of its memory, computing speed and so on. Based on the evaluation of this we can give a score in terms of marks on a scale of 100. So all data centers will have server and each server will have a score based on the evaluation of its capacity[7].

Fig:2 Represents the proposed load balancing.



The centralized control system will allocate the task that comes to the system by evaluating the amount of resources required in a tentative manner and based on that its allocates the server with either high or low scoring. The maximum the resources are required the server with highest score will be allocated and where ever the low resources are required as server with less score will be assigned or allocated. With this an advantage is the work or the task that comes will be given to appropriate server and there by the task completes faster and efficiently. The server then will be ready to take up the next task.

All server will also have threshold level and this is also taken into account for allocation of task. The threshold level will help us to know the amount of load that needs to be assigned to that server, once the server reaches that level the other server with next level of score will be given the task.

The sum of loads of all virtual machines is defined as

$$L = \sum_{i=1}^k l_i$$

where  $i$  represents the number of VMs in a data center.

The load per unit capacity is defined as

$$LPC = \frac{L}{\sum_{i=1}^m c_i}$$

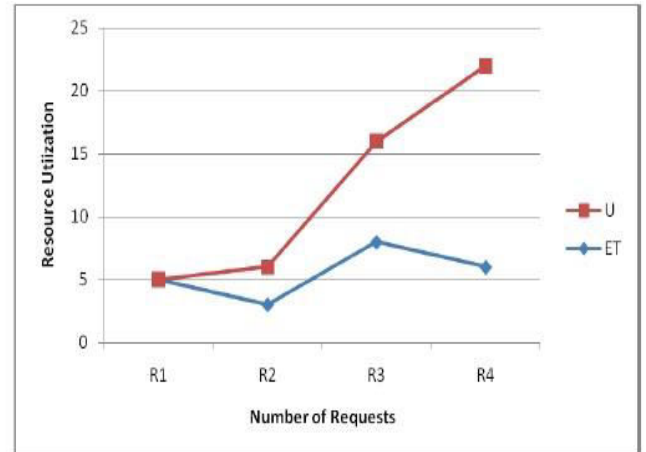
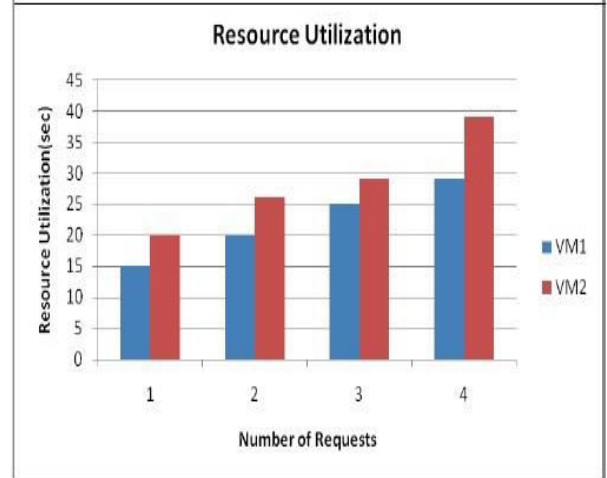
$$\text{Threshold } T_i = LPC * c_i$$

where  $c_i$  is the capacity of the node.

The load imbalance factor of a particular virtual machine is given by

$$\text{If VM } \begin{cases} < \left| T_i - \sum_{v=1}^k L_v \right|, & \text{Underloaded,} \\ > \left| T_i - \sum_{v=1}^k L_v \right|, & \text{Overloaded,} \\ = \left| T_i - \sum_{v=1}^k L_v \right|, & \text{Balanced.} \end{cases}$$

Fig: 3Shows Resource Utilization



### Simulation toolkits

Thinking of unpredicted network environment and laboratory resource scale (like servers), sometimes it is helpful to and more convenient for developing and running simulation tools to simulate large-scale experiments. The research on dynamic and large-scale distributed environment can be fulfilled by constructing data center simulation system, which offers visualized modeling and simulation for large-scale applications in cloud infrastructure[9].

The data center simulation system can describe the application workload statement, which includes user information, data center position, the amount of users and data centers, and the amount of resources in each data center.<sup>10</sup> Under the simulated data centers, load balancing algorithms can be easily implemented and evaluated.

**CloudSim:** CloudSim is an event-driven simulator implemented in Java. Because of its object-oriented programming feature, CloudSim allows extensions and definition of policies in all the components of the software stack, thereby making it a suitable research tool that can mimic the complexities arising from the environments.<sup>11</sup>

**CloudSched:** CloudSched enables users to compare different resource scheduling algorithms in Infrastructure as a Service (IaaS) regarding both hosts and workloads. It can also help the developer identify and explore appropriate solutions considering different resource scheduling algorithms.<sup>9</sup>

## Conclusion:

Load Balancing is a necessary task in Cloud Computing environment to attain maximum use of resources. In this paper, we talk about Evaluation based load balancing method which helps in providing maximum efficiency in terms of execution and allocation of various tasks. The advantage of this algorithm is most sought after resources task will get the maximum scored server to perform computation and least resource required task gets the least scored server. As high rated server is performing high rated task and low rated server is performing low rated task the resource utilization will happen in most appropriate manner.,

## References

- 1 D. Saranya et.al, "Load Balancing Algorithms in Cloud Computing: A Review," International Journal of Advanced Research in Computer Science and Software Engineering, vol. 5, Issue 7, July 2015.
- 2 S. Sethi et.al, "Efficient Load Balancing in Cloud Computing using Fuzzy Logic," IOSR Journal of Engineering (IOSRJEN) ISSN: 2250-3021 vol. 2, pp. 65-71, July 2012.
- 3 T. Desai et.al, "A Survey of Various Load Balancing Techniques and Challenges in Cloud Computing," International Journal of Scientific & Technology Research, vol. 2, Issue 11, November 2013.
- 4 S. Rajoriya et.al, "Load Balancing Techniques in Cloud Computing: An Overview," International Journal of Science and Research (IJSR), vol. 3, Issue 7, July 2014
- 5 Sharma S. et.al, "Performance Analysis of Load Balancing Algorithms," World Academy of Science, Engineering and Technology, 38, 2008.
- 6 Gross D. et.al, "Noncooperative load balancing in distributed systems", Elsevier, Journal of Parallel and Distributed Computing, No. 65, pp. 1022-1034, 2005.
- 7 Nikravan M. et.al, "A Genetic Algorithm for Process Scheduling in Distributed Operating Systems Considering Load Balancing", Proceedings 21st European Conference on Modelling and Simulation (ECMS), 2007.
- 8 Sunil.T Accountability Clouds for Serving Auditing Cloud computing
- 9.Tian W, Xu M, Chen A, et al. Open-source simulators for cloud computing: Comparative study and challenging issues. Simul Modell Pract Theory. 2015;58:239–254.
- 10.Calheiros RN, Ranjan R, Beloglazov A, De Rose CA, Buyya R. Cloudsim: A toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms. Software: Pract Exper. 2011;41(1):23–50.
- 11.Tian W, Zhao Y, Xu M, Zhong Y, Sun X. A toolkit for modeling and simulation of real-time virtual machine allocation in a cloud data center. IEEE Trans Autom Sci Eng. 2015;12(1):153–161.

# Time-of-Arrival (TOA), Angle-of-Arrival (AOA) and Hybrid - TOA & AOA based Localization in Wireless Sensor Networks

**Dr.V.Bhoopathy<sup>1</sup>,**

*Professor,*

*Department of Computer Science and Engineering,*

*Malla Reddy College of Engineering*

E-mail: [v.bhoopathy@gmail.com](mailto:v.bhoopathy@gmail.com)

**M Aharonu<sup>2</sup>.**

*Assistant Professor,*

*Department of Computer Science and Engineering,*

*Malla Reddy College of Engineering*

E-mail: [aharon.mattakoyya@gmail.com](mailto:aharon.mattakoyya@gmail.com)

**Abstract** - This paper mainly focuses to reduce the localization error of a sensor node which is deployed in a Wireless Sensor Network (WSN). Here, time-of-arrival (TOA) and angle-of-arrival (AOA) based random transmission directed localization (RTDL) technique is considered. This technique can be applied in wireless sensor networks, especially suitable for network with low frequency range in the wireless sensor network. In this work, localization Error could be improved via TOA-AOA during node communicating with each other. Ad hoc On-Demand Distance Vector (*AODV*) routing protocol to be implemented in this work. The node categories in two ways: those nodes equipped with known vector position of the Omni-direction antenna (Source-nodes) and those nodes equipped with unknown vector position of the Omni-direction antenna (Sink-nodes). All nodes are capable of communicating with other nodes. Source-nodes are capable of positioning (TOA-AOA estimation) the other nodes located in their coverage area. The system estimates the error rate of distance between the nodes, by measuring TOA-AOA based RTDL at an appropriate number of nodes. This paper shows the proficiency of the proposed method to reduce the localization error and determine the attacker's nearest location in the network.

**Keywords:** Wireless Sensor Networks (WSN), Time-of-Arrival (TOA), Angle-of-Arrival (AOA), Ad-hoc On-Demand Distance Vector (*AODV*)

## I. INTRODUCTION

### 1.1 Wireless Sensor Networks

Wireless sensor networks comprises of the upcoming technology that has attained noteworthy consideration from the research community. Sensor networks comprise of many small, low cost devices and are naturally self-organizing ad hoc systems. The function of the sensor network is monitoring the physical environment, collect and transmit the information to other sink nodes. In general the range of the radio transmission for the sensor networks are in the orders of the magnitude which is smaller than the geographical extent of the intact network [1].

Wireless sensor network comprises of a great number of minute electromechanical sensor devices which possess the sensing, computing and communication abilities. These devices can be utilized for gathering sensory

information, like measurement of temperature from an extended geographical area [2].

Many of the features of the wireless sensor networks give rise to challenging problems [3-7]. The most important three characteristics are:

- Sensor nodes are the ones which are prone to maximum failures.
- Sensor nodes make use of the broadcast communication pattern and have severe bandwidth restraint.
- Sensor nodes have limited amount of resources.

Despite the huge research effort, still a well-accepted approach on how to solve the localization issue is being realized. Since the sensor nodes are inexpensive and are in huge number it is not practical to equip these sensors with a Global Positioning System (GPS) receiver. Various localization approaches have been proposed and can be seen in the literature [8] and there is not a single approach which is simple, distinct and gives decentralized solution for WSNs. The Ultra-Wide Band (UWB) techniques [9] give very decent localization accuracy but the systems are expensive.

The commonly used approaches for measuring position estimate in WSN are Time of Arrival (TOA) [10], Time Difference of Arrival (TDOA)[16], Received Signal STRENGTH (RSS)[17] AND ANGLE OF ARRIVAL (AOA) A.K.A., Direction of Arrival (DOA)[18]. Where, the TOA, TDOA, and RSS measurement gives the distance calculation between the source sensor and the receiver sensors while DOAs provide the information of the angle and the distance measurements from the source and the receiver. Calculating these distance and angle measurements is not simple because of the nonlinear relationships with the source.

Given the TOA, TDOA, RSS and DOA information, the main focus of this paper is based on TOA positioning algorithms. We consider a two dimensional (2D) rectangular area where the sensors are deployed in Line-of-Sight (LOS) transmission, i.e., there is a direct path between the source and each receiver [19]. Also, we conclude that the measurements are well inside the expected range in order to obtain reliable location estimation.

## II. RELATED WORK

In the view of localization infrastructure, [11] they utilized infrared strategies, and [12] ultrasound to perform localization. Both of them have to deploy specific framework for localization. Then again, disregarding its few meter-level accuracy [13], utilizing RSS [14, 15], the work in [16] is an alluring methodology, in light of the fact that it can reuse the current wireless infrastructure. Dealing with ranging approach, range based algorithms include distance estimation to milestones by utilizing the estimation of different physical properties [17], for example, RSS [14, 18], time of arrival (TOA) [19], and time difference of arrival (TDOA) [12].

Mobile location with TOA/AOA data at single base station is initially proposed in [20]. The authors in [21] examine the location exactness of TOA/AOA hybrid algorithm with a single base station in the LOS situation. Deng and Fan [22] present TOA/AOA location algorithm with various base stations. Moreover, the velocity of the mobile station is thought to be low and the relative movement between the base station and the mobile station is not recognized [23]. Utilizes an obliged nonlinear optimization technique, when range estimations are accessible from three base stations. Limits on the Non-Line-Of-Sight (NLOS) error and the relationship between the real ranges are obtained from the geometry of the cell design and the measured range circles to serve as demands [24]. Introduces two hybrid TOA/AOA systems, Enhanced Time of Arrivals (E-TOA) and Enhanced Angle of Arrival (E-AOA), so as to advance the estimation of location positioning [25]. Proposes a residual test (RT) that can in turn all the while focus the amount of LOS base stations and recognize them such that localization can move ahead with just those LOS base stations.

Hybrid location approaches by joining time and angle estimations which can minimize the amount of time taken by the receiving base stations and enhance the scope of location based service all the while. Exhaustive overviews of design difficulties and newly proposed hybrid positioning algorithms for wireless networks could be found in [26–28].

Range-free algorithms [29–31] use coarser measurements to place limits on applicant positions. An alternate technique for classification depicts the system for mapping a node to a location. Lateral approaches [31–33] use separations to milestones, which is as much angulation as it uses the angles from points of interest.

In [37], proposed a novel technique for minimizing the localization error in MANET. TOA/AOA based Random Transmission Directed Location (TAR) approach used for locating nodes. And also they used hybrid TOA/AOA which effortlessly segregates the nodes in the same vector with the help of the magnitude. Moreover, Ad hoc On-Demand Multipath Distance Vector has been used for reliable routing and to identify the presence of an attacker. Their simulation results have shown that TOA/AOA based on random transmission directed localization technique reduces the delay, drop rate, location error, estimation error, localization time and increases the throughput and network efficiency.

In [38], presented the problem of position estimation of a sensor node in a Wireless Sensor Network, using TOA measurements in LOS environments. The CRLB for the position estimation problem has been derived first and later four methods namely LLS, SA, WLLS and Two Step WLS methods of linear approach have been derived and presented. Their research work can be extended to the nonlinear approaches also and shall be reported in a future communication. They have not discussed in detail about TOA and combined solution for both TOA and AOA.

In [39], proposed a novel technique for an improved algorithm of the node localization in ad hoc network and decrease in location accuracy due to transmission delay is mitigated. They claimed that the positioning of the node according to current routing table will be incorrect if the nodes are moving fast. They proposed a node localization forecasting algorithm to plot the current position of the node. The proposed method reduces only the error impact of node localization due to transmission delay but no precise location is achieved.

## III. MATHEMATICAL MEASUREMENT MODEL

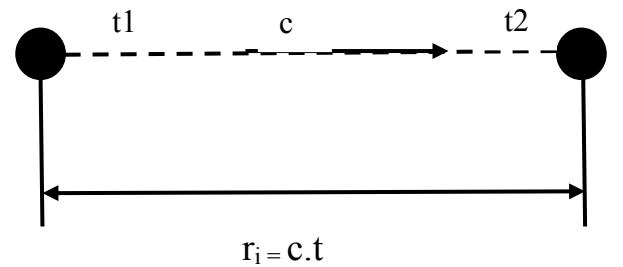
### 3.1 Time of Arrival (TOA) Ranging Measurements

Time of Arrival (TOA) is used method for positioning and ranging. The vital idea behind TOA is to estimate delay between the source nodes and sink. To attain this, source nodes and sink clock should be synchronized.

Then distance can be estimated from the speed of the signal from source to destination and also from the time taken by the signal from the source to destination.

$$r_i = c \cdot t_i = \sqrt{(x - x_i)^2 + (y - y_i)^2}$$

where  $r_i$  is the distance, where  $c$  represent signal propagation speed between nodes while  $t_i$  represents the time taken by the signal to travel from the source to destination.



**Figure 1: TOA approach for distance estimation**

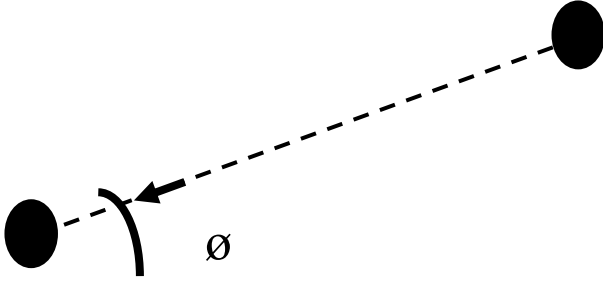
The distance estimation using TOA with two nodes, as shown in Fig. 1. It is shown that the distance between the nodes is proportional to the signal propagation time between two nodes.

### 3.2 Angle of Arrival (AOA) Ranging Measurement

Angle of Arrival (AOA) is another method to estimate the position of the target node based on the angle measured. The AOA technique can be further subdivided into two classes: those that make use of the response phase of receiver antenna and those that make use of amplitude phase of receiver antenna. Beam forming is done using anisotropy in the response pattern of the antenna. The beam of the antenna at the receiver is rotated and the

direction in which maximum signal strength is obtained is considered as the direction of the transmitter. But problem arises when there is varying transmitted signal strength.

Because of this, the receiver cannot differentiate variation in signal strength. To address this issue, a non-rotating omnidirectional antenna is used.



**Figure 2: AOA Approach for angle measurement**

AOA estimation at the target node provides information about the direction over which the target node has been located as shown in Fig. 2.

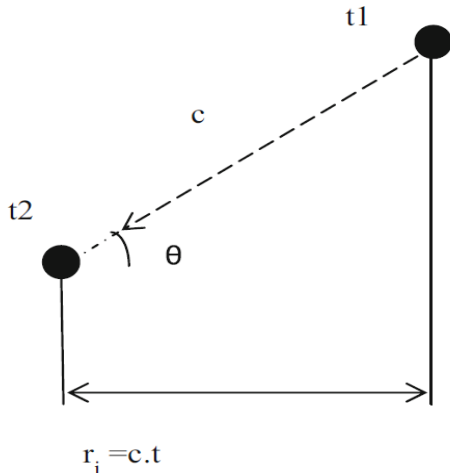
### 3.3 Hybrid TOA/AOA Ranging Measurement

Hybrid approach plays a vital role in providing location estimation. Hybrid techniques are combination of distance or range and AOA measurements to locate unknown nodes which are useful when the network coverage is poor. To locate a node with high accuracy, it is necessary to measure TOA and AOA to obtain the distance and direction estimates. With the estimated distance  $r_i$  and direction  $h_i$ , the position of the unknown node  $(x_i, y_i)$  from the base node  $(x, y)$  is computed as

$$x_i = x + r_i = \begin{bmatrix} \cos\theta_i \\ \sin\theta_i \end{bmatrix}$$

$$y_i = y + r_i = \begin{bmatrix} \cos\theta_i \\ \sin\theta_i \end{bmatrix}$$

which defines the location of the unknown node  $(x_i, y_i)$  from the TOA and AOA measurements with respect to the base node.

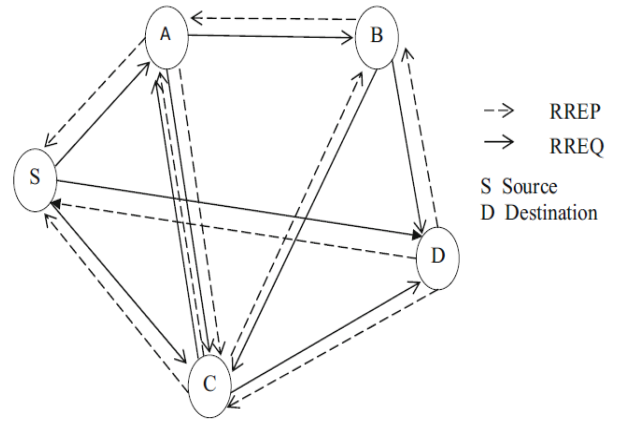


**Figure 3: TOA/ AOA approach for distance and angle measurement**

The TOA/AOA estimation is used to identify the position of the target node with known distance and angle measurements, as shown in Fig. 3.

## IV. TAR APPROACH AND AODV PROTOCOL

A novel TOA/AOA based on random transmission directed localization (TAR) technique is proposed to identify the location of the target node. Initially the network consists of random number of nodes. If the number of intermediate nodes between the source node and the destination node increases, the position accuracy decreases where by resulting in the increase of the error rate. To solve this issue, the proposed technique makes accurate identification of the target node by estimating the distance and the direction of the target node. In addition, the system uses an Ad hoc On-Demand Distance Vector Routing protocol (AODV) to determine shortest path between the source and destination that is reliable and secure (Fig. 3).



**Figure 4: RREQ and RREP propagation in AODV**

### 4.1 TOA/AOA Based on Random Transmission Directed Localization

The initial assumption is that the packet data size and the number of intermediate nodes between the source and the destination are obtained. The localization process starts by forwarding the data packet to the destination. The departure time and the arrival time of the data packets are recorded. Next, packet travelling time is estimated, in addition packet waiting time or regeneration time if any are added to packet travelling time to obtain TOA for measuring the distance using Eq. 1. The direction of the node is then obtained by using AOA measuring technique. The TOA/AOA measurements are then used to obtain the position of the destination node. In this approach, TOA and AOA are computed for two cases: two nodes with same magnitude and vector and also for two nodes with different magnitude and vector. Finally the accuracy and error factor of the TOA is computed using Eq. 4 and 5 respectively.

$$AF = (AT_j * AD_j) - (ST_j * AD_j)$$

where  $AT_j$  is the packet arrival time at the  $j$ th round,  $AD_j$  is the angular velocity at the destination at the  $j$ th data transfer cycle,  $ST_j$  is the start time of the packet while  $AF$  is the Accuracy Factor.

The error factor is estimated as,

$$E = \sqrt{(ET - ST)} / \sqrt{(AD - AS) + SC}$$

where ET is the arrival time of the packet at the destination, ST is departure time of the packet at the source, AD angular velocity at the destination, AS is the angular velocity at the source, SC is the packet sequence count while E is the error factor.

#### 4.2 Ad hoc On-Demand Distance Vector Routing protocol (AODV)

Ad hoc On-Demand Distance Vector (AODV) is the most widely used ad hoc routing protocol on the on-demand basis. AODV uses a hop-by-hop routing instead of source routing. The key idea in AODV is to compute multiple paths during route discovery.

AODV uses the concept of advertised hop count in order to maintain multiple hops with the same destination sequence number. The advertised hop count field contains the total number of paths to the destination, next hop ip has the list of next hop nodes, hop count field gives the corresponding hop count value, while the entry expiration time field provides the expiry time after which a request reply packet will not be received and in turn discards the entry.

When the source node needs to forward packet to the destination, it initializes route discovery by broadcasting Route Request Message (RREQ) to the destination. The intermediate upon receiving the RREQ packets, sets up the reverse path to the source using the previous hop count as the next hop value. If the intermediate node contains valid route to the destination, a RREP packet is generated and forwards to the source. If the intermediate node does not contain route to the destination, it rebroadcasts the RREQ packets. The duplicate RREQ packets received by the node define an alternate path to the destination.

Any two RREQ packets arriving at the intermediate node via unique neighbor of the source would not have travelled the same node, since node does not forward duplicate RREQ packets. This ensures node disjoint path. When the destination node receives the intermediate node, it generates RREQ packet and forwards it to the source via the reverse path. As the RREP packet traverses to the source, a forward path is established to the destination node. Here a loop free path is assured by only accepting the alternate path whose hop count is less than the advertised hop count of the destination path.

### V. SIMULATION RESULTS

#### 5.1. Simulation Setup

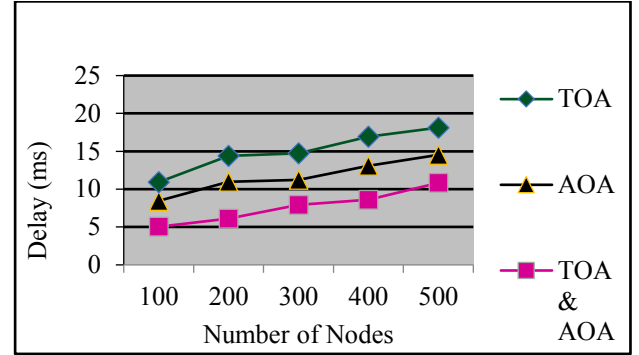
The performance of Time of arrival and Angel of arrival is evaluated through NS2 simulation [12]. A random network deployed in an area of 1000 X 1000 m is considered. Initially the nodes are placed randomly in the specified area. The base station is assumed to be situated 100 meters away from the above specified area. The IEEE 802.11b MAC layer is used for a reliable and single hop communication among the devices, providing access to the physical channel for all types of transmissions and appropriate security mechanisms.

**Table 1:** Simulation Parameters

No. of Nodes	100
Area Size	1000 m X 1000 m
MAC	802.11b
Routing protocol	AODV
Simulation Time	600 sec
Traffic Source	CBR
Packet Size	512 bytes
Transmission Range	250m
No. of events	4
Speed of events	20 m/s

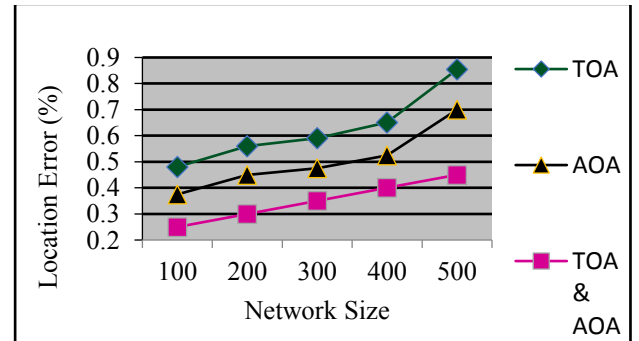
#### 5.2. Simulation Result

The simulation results are presented in below



**Figure 5:** Nodes Vs Delay

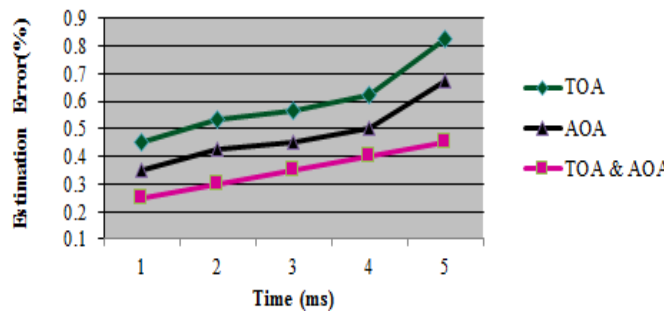
Figure 5 gives delay analysis with respect to the number of nodes. The delay for TOA, AOA and HYB-TOA & AOA when the number of nodes is increased. From the figure, it can be seen that the Time of arrival (TOA), Angle of arrival (AOA) has high delay when compared with Hybrid time of arrival and hybrid angle of arrival (HYB-TOA & AOA).



**Figure 6:** Network Size Vs Location Error

Figure 6 presents the location error of for TOA, AOA and method. The location error and it is defined as the process of estimating the location of mobile nodes in

the WSN. The error rate for each transmission is computed in milliseconds (ms) and it is shown that there is reduced number of errors in proposed approach.



**Figure 7: Time Vs Estimation Error**

Figure 7 shows the graphical representation of estimation error for a particular time and the estimation error factor of the proposed technique in case of both direct and n-hop neighbors is reduced. The error rate for each transmission is computed in milliseconds (ms) and it is shown that there is reduced number of errors in proposed approach.

## VI. CONCLUSION

This paper mainly focused to reduce the localization error of a sensor node which is deployed in a Wireless Sensor Network (WSN). Hence, time-of-arrival (TOA) and angle-of-arrival (AOA) based random transmission directed localization (RTDL) technique was implemented. In this work, Ad hoc On-Demand Distance Vector (AODV) routing protocol implemented to improve localization Error via TOA–AOA. The system estimated the error rate of distance between the nodes, by measuring TOA–AOA based RTDL at an appropriate number of nodes. The research work shows the proficiency of the proposed method to reduce the localization error and determine the attacker's nearest location in the network.

## References

- [1] Dorottya Vass, Attila Vidaacs, "Distributed Data Aggregation with Geographical Routing in Wireless Sensor Networks", Pervasive Services, *IEEE International Conference* on July 2007.
- [2] Jukka Kohonen, "Data Gathering in Sensor Networks", Helsinki Institute for Information Technology, Finland. Nov 2004.
- [3] Gregory Hartl, Baochun Li, "Loss Inference in Wireless Sensor Networks Based on Data Aggregation", IPSN 2004.
- [4] Bhoopathy, V. and Parvathi, R.M.S. "Energy Efficient Secure Data Aggregation Protocol for Wireless Sensor Networks", *European Journal of Scientific Research*, Vol. 50, Issue 1, pp.48-58, 2011.
- [5] Bhoopathy, V. and Parvathi, R.M.S. "Secure Authentication Technique for Data Aggregation in Wireless Sensor Networks" *Journal of Computer Science*, Vol. 8, Issue 2, pp 232-238, 2012.
- [6] Bhoopathy, V. and Parvathi, R.M.S. "Energy Constrained Secure Hierarchical Data Aggregation in Wireless Sensor Networks" *American Journal of Applied Sciences*, Vol. 9, Issue 6, pp. 858-864, 2012.
- [7] Bhoopathy, V. and Parvathi, R.M.S. "Securing Node Capture Attacks for Hierarchical Data Aggregation in Wireless Sensor Networks" *International Journal of Engineering Research and Applications*, Vol. 2, Issue 2, pp. 458-466, 2012.

- [8] H. Wymeersch, J. Lien and M. Z. Win, (2009), "Cooperative localization in wireless networks," *IEEE Signal Processing Mag.*, vol.97, no.2, pp.427-450.
- [9] Gezici S, Tian Z, Giannakis GB, Kobayashi H, Molisch AF, Poor HV, et al. (2005), "Localization via ultra-wideband radios," *IEEE Signal Process Mag.*; 22:70–84.
- [10] A. Jagoe (2003), *Mobile Location Service: The Definitive Guide*, Upper Saddle River: Prentice- Hall
- [11] Want, R., Hopper, A., Falcao, V., & Gibbons, J. (1992). The active-badge location system. *ACM Transactions on Information Systems*, 10(1), 91–102.
- [12] Priyantha, N., Chakraborty, A., & Balakrishnan, H. (2000). The cricket location support system. In *Proceedings of ACM MobiCom*, pp. 32–43.
- [13] Krishnakumar, A., & Krishnan, P. (2005). On the accuracy of signal-strength based location estimation techniques. In *Proceedings of IEEE INFOCOM*, pp. 642–650.
- [14] Bahl, P., & Padmanabhan, V. N. (2000). RADAR: An in-building RF-based user location and tracking system. In *Proceedings of IEEE INFOCOM*, pp. 775–784.
- [15] Chen, Y., Francisco, J., Trappe, W., & Martin, R. P. (2006). A practical approach to landmark deployment for indoor localization. In *Proceedings of 3rd IEEE SECON*, pp. 365–373.
- [16] Yang, J., & Chen, Y. (2008). A theoretical analysis of wireless localization using RF-based fingerprint matching. In *Proceedings of 4th SMTSP*, pp. 1–6.
- [17] Patwari, N., Ash, J. N., Kyperountas, S., Hero, A. O., Moses, R. L., & Correal, N. S. (2005). Locating the nodes. *IEEE Signal Processing Magazine*, 22(4), 54–69.
- [18] Chen, Y., Kleisouris, K., Li, X., Trappe, W., & Martin, R. P. (2006). The robustness of localization algorithms to signal-strength attacks: A comparative study. In *Proceedings of DCSS*, pp. 546–563.
- [19] Enge, P., & Misra, P. (2001). *Global positioning system: Signals, measurements, and performance*. Lincoln, MA: Ganga-Jamuna.
- [20] Cesbron, R., & Arnott, R. (1998). Locating GSM mobiles using antenna array. *Electronics Letters*, 34, 1539–1540.
- [21] So, H. C., & Shiu, E. M. K. (2003). Performance of TOA–AOA hybrid mobile location. *IEICE Transactions on Fundamentals*, E86-A, 2136–2138.
- [22] Deng, P., & Fan, P. -Z. (2000). An AOA assisted TOA positioning system. In *Proceedings of International Conference on Communication Technology*, Beijing, China, pp. 1501–1504.
- [23] Venkatraman, S., Caffery, J, Jr, & You, H.-R. (2004). A novel TOA location algorithm using LOS range estimation for NLOS environments. *IEEE Transactions on Vehicular Technology*, 53, 1515–1524.
- [24] Deligiannis, N., Louvros, S., & Kotsopoulos, S. (2007). Optimizing location positioning using hybrid TOA-AOA techniques in mobile cellular networks. In *Proceedings of Mobimedia'07*, Nafpaktos, Greece, pp. 1–7.
- [25] Chan, Y.-T., Tsui, W.-Y., So, H.-C., & Ching, P.-C. (2006). Time-of-arrival based localization under nlos conditions. *IEEE Transactions on Vehicular Technology*, 55, 17–24.
- [26] Guvenc, I., & Chong, C.-C. (2009). A survey on TOA based wireless localization and NLOS mitigation techniques. *IEEE Communications Surveys & Tutorials*, 11, 107–124.
- [27] Gustafsson, F., & Gunnarsson, F. (2005). Mobile positioning using wireless networks. *IEEE Signal Processing Magazine*, 22, 41–53.
- [28] Tang, H., Park, Y.-W., & Qiu, T.-S. (2008). A TOA–AOA-based NLOS error mitigation method for location estimation. *EURASIP Journal on Advances in Signal Processing*, 8, 1–14.
- [29] Shang, Y., Ruml, W., Zhang, Y., & Fromherz, M. P. J. (2003). Localization from mere connectivity. In *Proceedings 4th ACM MobiHoc*, pp. 201–212.
- [30] He, T., Huang, C., Blum, B., Stankovic, J. A., & Abdelzaher, T. (2003). Range-free localization schemes in large-scale sensor networks. In *Proceedings of 9th ACM MobiCom*, pp. 81–95.
- [31] Niculescu, D., & Nath, B. (2001). Ad hoc positioning system (APS). In *Proceedings of IEEE GLOBECOM*, pp. 2926–2931.

- [32] Enge, P., & Misra, P. (2001). Global positioning system: Signals, measurements, and performance. Lincoln, MA: Ganga-Jamuna.
- [33] Langendoen, K., & Reijers, N. (2003). Distributed localization in wireless sensor networks: A quantitative comparison. *Computer Networks*, 43(4), 499–518.
- [34] Bahl, P., & Padmanabhan, V. N. (2000). RADAR: An in-building RF-based user location and tracking system. In *Proceedings of IEEE INFOCOM*, pp. 775–784.
- [35] Youssef, M., Agrawal, A., & Shankar, A. U. (2003). WLAN location determination via clustering and probability distributions. In *Proceedings 1st IEEE PerCom*, pp. 143–150.
- [36] Chen, Y., Kleisouris, K., Li, X., Trappe, W., & Martin, R. P. (2006). The robustness of localization algorithms to signal-strength attacks: A comparative study. In *Proceedings of DCOSS*, pp. 546–563.
- [37] Kalpana & Baskar, “TAR: TOA–AOA Based Random Transmission Directed Localization” (2016), *Wireless Pers Commun.*
- [38] Ravindra & Jagadeesha, (2015) “Time of Arrival Based Localization in Wireless Sensor Networks: A Non-Linear Approach” *Signal & Image Processing: An International Journal (SIPIJ)*.
- [39] Wang Anbao and Zhu Bin, (2014) “An Improved Algorithm of the Node Localization in Ad Hoc Network”, *JOURNAL OF NETWORKS*, VOL. 9, NO. 3, pp. 549 -557.

# Speaker Segmentation- an Comparative study using Support Vector Machines and Auto Associative Neural Network

<sup>1</sup> Dr.J.Gladson Maria Britto\* and <sup>2</sup> Mrs.B.Ananthi

<sup>1</sup>Professor, Department of Computer Science & Engineering, Malla Reddy College of Engineering, Secunderabad  
gmbritto@gmail.com

<sup>2</sup>Assistant Professor, CSE, Vivekananda College of Engineering for Women, Thiruchencode, India.

## ABSTRACT

In this paper we propose a classification based method to identify speaker turn point detection and segmenting speech contains individual speaker using support vector machines (SVM) and Auto associative neural network (AANN). Speaker turn point detection is important for automatic segmentation of multi speaker speech data into homogenous segments with each segment containing the data of one speaker only. Existing approach for speaker turn point detection are based on the dissimilarities of the distribution of data before and after a speaker turn point. Patterns extracted from the data around the speaker turn points are used as positive examples. Patterns extracted from the data between the speaker turn point are used as negative examples. The linear predictive cepstral coefficients (LPCC) and Mel frequency cepstral coefficients (MFCC) and extracted from the speech signal, the positive and negative examples are used in training a SVM and AANN separately for speaker turn point detection. The extraction of fixed length pattern from speaker are given as input to SVM and AANN models are used to classify the speaker turn points and speaker no turn points using specific features. Experiments are carried out on different audio databases and the proposed method is better for detecting speaker turn point changes with sort duration of speech.

## Keywords:

*Linear Predictive Coefficients (LPC), Linear Predictive Cepstral Coefficients (LPCC), Mel-Frequency Cepstral Coefficients (MFCC), Weighted Linear Predictive Cepstral Coefficients (WLPCC), Support Vector Machines (SVM) , Autoassociative Neural Network (AANN).*

## INTRODUCTION

Speaker turn point detection involves determining the points at which there is a speaker turn changes in the multi speaker speech data as in audio recordings of conversation, broadcast news and movie. Speaker turn point detection is the first step in the speaker based segmentation of multi speaker only. Speaker segmentation is important for tasks such as audio indexing, speaker tracking and speaker adaptation in automatic transcription of conversational speech. Speaker turn point detection should do without the knowledge of the number of speakers and the identity of speakers. Therefore, a Speaker turn point[12],[13],[14] detection systems should be speaker independent.

The existing approaches for Speaker turn point detection are based on the dissimilarity in the distributions of data before and after the points of speaker change. Dissimilarity measurement is commonly based on comparison of the parametric statistic model of the

distribution such as Mahalanobis distance, Weighted Euclidean distance, Bayesian information criteria. In these approaches for Speaker turn point detection, the dissimilarity is measured for the data between two adjacent windows of fixed length. The points at which the dissimilarity is above a threshold are hypothesized as the speaker turn points. We propose an approach in which a classification model is trained to detect the Speaker turn point points and segment the data for according to the speakers.

## 1. ACOUSTIC FEATURE EXTRACTION

Acoustic features representing the speaker information can be extracted from the speech signal at the segmental features are the features[2],[3],[1] extracted from the short (20 milliseconds) segments of the speech signal. These features represent the short time spectrum speech signal. The short time spectrum envelop of the speech signal is attributed primarily to the shape of the same sound uttered by two persons may differ due to change in the shape of the individual's vocal tract[4]system and the manner of speech production. For acoustic feature extraction, the differenced speech signal is divided in to frame of 20 milliseconds, with a shift of 10 milliseconds. Feature extraction is done by using LPC and LPCC.

### 1. a) Data Collection:

Two speaker Conversation speech signal is recorded.

- i. Male-male conversation.
- ii. Male-Female conversation.
- iii. Female-Female conversation.

The recording rate for speech is 8 KHz. The sampling bit rate is 16 bits. The recording mode is mono. The sample values vary from -32,768 to +32,767. We used unidirectional microphone. For 1 second 8000 samples will be recorded. Frame size is 20 milliseconds (160 samples).Frame shift is 10 milliseconds (80 samples).The file is stored as .wav extension format.

### 1. b) LPC Model:

Linear Predictive Coefficients (LPC) is a powerful speech analysis technique. It is predominant technique for estimating the basic speech parameters. e.g. pitch, formants and vocal tract area function and for representing speech for low bit rate transmission or storage[15],[16],[17].

- i. Linear prediction coefficients.
- ii. Linear prediction (LP) analysis.

Each sample is predicted as linear weighted sum of past p samples, where p is the order of LP analysis[8],[9],[10].

$$x(n) = \sum_{k=1}^p a_k x(n-k) \quad (1)$$

The predicted signal value is given in the above Equation (1)

$x(n)$  is the predicted signal value.

$x(n-k)$  is the previous observed values.

$a_k$  is the predictor coefficients.

For example p = 14

$$x(15) = a_1 x(14) + a_2 x(13) + \dots a_{14} x(1)$$

$$x(16) = a_1 x(15) + a_2 x(14) + \dots a_{14} x(2)$$

.....

$$x(160) = a_1 x(159) + a_2 x(158) + \dots a_{14} x(146) \quad (2)$$

Linear prediction coefficients (LPC)  $\{a_k\}$  are determined from the above equations.

The basic idea behind the LPC model is that given a speech sample at time  $n$ ,  $s(n)$  can be approximated as linear combination of the past p speech samples.

The linear prediction analysis is to determine the set of predictor coefficient  $\{a_k\}$ , directly from the speech signal. So that the spectral properties of the digital from the below match those of the speech wave from with in the analysis window. Since the spectral characteristics of the speech vary over time, the predictor coefficients at a given time N must be estimated from a short segment of the speech signal occurring around time n. Thus the basic approach is to find a set of prediction coefficients that minimize the mean squared prediction error over a short segment of the speech wave form.

Durbin's recursive solution for the autocorrelation equation is used for finding LPC Coefficients.

### 1.c) Autocorrelation method:

A simple and straight forward way of defining the limits on m in this summation is to assume that the speech segment  $S_n(m)$ , is identically zero outside the intervals  $0 \leq m \leq (n-1)$ . Thus the speech sample for the minimization can be expressed as  $0 \leq m \leq (n-1)$

$$S_n(m) = s(m+n).w(n), \quad 0 \leq m \leq (n-1)$$

= 0 other wise

### Autocorrelation analysis:

Each frame of windowed signal is next autocorrelated to give

$$r_n(i-k) = \sum_{m=0}^{N-1-(i-k)} S_n(m) S_n(m+i-k) \quad (3)$$

where,

$r_1$  is the energy of the 1<sup>st</sup> frame using Equation (3).

Where the highest autocorrelation value P is the order of LPC. P values from 8 to 20 are used.

$R_1(0)$  = energy of the 1<sup>st</sup> frame.

In this paper, p=16 gives better performance compared to other order of LPC. The autocorrelation function is symmetric,

i.e.  $r_n(-k) = r_n(k)$ , the LPC equations can be expressed as

Autocorrelation equation

$$\sum_{k=1}^p \hat{a}_k r_n(|i-k|) = r_n, \quad 1 \leq i \leq p \quad (4)$$

where  $a_k$  is the predictor coefficients.

This is expressed in matrix form

$$\begin{bmatrix} r_n(0) & r_n(1) & r_n(2) & \dots & r_n(p-1) \\ r_n(1) & r_n(0) & r_n(1) & \dots & r_n(p-2) \\ r_n(2) & r_n(1) & r_n(0) & \dots & r_n(p-3) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_n(p-1) & r_n(p-2) & r_n(p-3) & \dots & r_n(0) \end{bmatrix} \begin{bmatrix} \hat{a}_1 \\ \hat{a}_2 \\ \hat{a}_3 \\ \vdots \\ \hat{a}_p \end{bmatrix} = \begin{bmatrix} r_n(1) \\ r_n(2) \\ r_n(3) \\ \vdots \\ r_n(p) \end{bmatrix} \quad (5)$$

### 1. d) Durbin's recursive solution for autocorrelation Equation:

The most efficient method for solving this particular system of equation is Durbin's recursive procedure which can be stated as. The process of solving for the predictor coefficients for the predictor of an order p. the solution for the predictor coefficient of all order less than p have also been obtained (i.e.) the predictor coefficient for a predictor of order 2.

$$E^{(0)} = r(0)$$

$$k_i = \left\{ r(i) - \sum_{j=1}^{L-1} \alpha_j^{(i-1)} r(|i-j|) / E^{(i-1)} \right\}, 1 \leq i \leq p$$

(6)

$$\alpha_i^{(i)} = k_i$$

$$\alpha_j^{(1)} = \alpha_j^{(i-1)} - k_i \alpha_{i-j}^{(i-1)}$$

$$E^{(i)} = (1 - k_i^2) E^{(i-1)} \quad (6)$$

where,

$$\text{LPC coefficients} = a_m = \alpha_m^{(p)}, 1 \leq m \leq p$$

$k_m$  = PARCOR coefficients.

#### LPC parameter conversion to Linear Predictive cepstral coefficients (LPCC)

A very important LPC parameter set, which can be derived directly from the LPC coefficient set, is the LPC cepstral coefficients,  $c(m)$ . The recursion used is

$$C_0 = \ln \sigma^2 \quad (7)$$

$$c_m = a_m + \sum_{k=1}^{m-1} \left( \frac{k}{m} \right) c_m a_{m-k}, 1 \leq m \leq p \quad (8)$$

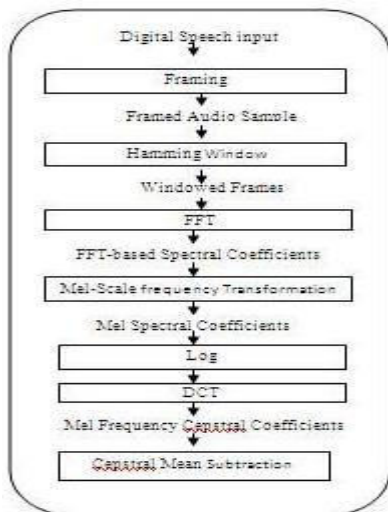
$$c_m = + \sum_{k=1}^{m-1} \left( \frac{k}{m} \right) c_m a_{m-k}, m > p, \quad (9)$$

where,

$\sigma^2$  is the gain term in the LPC model.

$C_0$  to  $c_m$  are LPCC coefficients

Suppose,  $p=19$  means 19<sup>th</sup> LPCC extracted from conversation speech signal. The cepstral coefficients, which are the coefficients of the Fourier transform representation of the log magnitude spectrum.

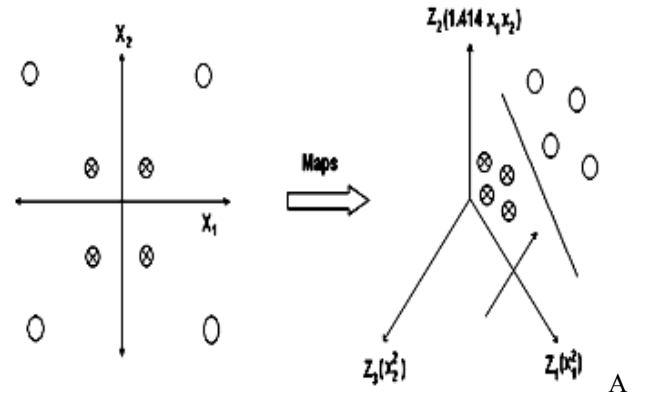


**Fig.1** MFCC feature extraction

#### 1. e) Mel-frequency Cepstral Coefficients (MFCCs):

MFCCs have been widely used in the field of Speaker turn point detection system and are able to represent the dynamic features of a signal as they extract both linear and non-linear properties. MFCC can be a useful tool of feature extraction in vibration signals as vibrations contain both linear and non-linear features. The Mel-frequency Cepstral Coefficients (MFCC) is a type of wavelet in which frequency scales are placed on a linear scale for frequencies less than 1 kHz and on a log scale for frequencies above 1 kHz. The complex cepstral coefficients obtained from this scale are called the MFCC. The MFCC contain both time and frequency information of the signal and this makes them useful for feature extraction. The following steps are involved in MFCC computations[5],[6],[7].

#### 2. MODELING FOR SPEAKER TURN POINT DETECTION



**Fig.2:** SVM Maps 2-Dimensional input space to 3-Dimensional input Space

support vector machine (SVM) [11] is a machine learning technique that learns the decision surface through a process of discrimination and has good generalization characteristics. SVM is based on the principle of structural risk minimization. Like RBFNN (Radial Basis Function Neural Network), support vector machines can be used for pattern classification and non linear regression. Support vectors are used to find hyper plane between two classes. Support vectors are close to the hyper plane. Support vectors are the training samples that define that optimal separating hyper plane and are difficult patterns to classify. For linearly separable data, SVM finds a separating hyper plane, which separates the data with the largest margins. For linearly separable data, it maps the data in the input space into high dimension space  $x \in R^1 \rightarrow \Phi(x) \in R^H$  with kernel function

$\Phi(x)$ , to find the separating hyper plane. SVM was originally developed for two class classification problems. The N class classification problem can be solved using NSVMs. Each SVM separates a single class from all the remaining classes (one-vs.-rest approach).

Given a set of features corresponding to N subjects for training, N SVMs are created. Each SVM is trained to distinguish between features of a single subject and all other features in the training set. During testing, the distance from

$x$  to the SVM hyper plane is used to accept or reject the identity claim of the subject.

Inner product kernel maps input space to higher dimensional feature space. Inner product kernel represents

$$K(x, x_i) = \Phi(x) \cdot \Phi(x_i). \quad (10)$$

where  $x$  is input patterns,  $x_i$  is support vectors.

For example,

$$\text{assume } x = [x_1, x_2]^T \text{ is input patterns} \quad (11)$$

$$x_i = [x_{i1}, x_{i2}]^T \text{ is support vectors} \quad (12)$$

$$K(x, x_i) = (x^T x_i)^2 \\ = (x_1 x_{i1} + x_2 x_{i2})^2 \quad (13)$$

$$= x_1^2 x_{i1}^2 + x_2^2 x_{i2}^2 + 2 x_1 x_2 x_{i1} x_{i2} \quad (14)$$

where,

$$\Phi(x) = (x_1^2, x_2^2, 1.414 x_1 x_2)$$

$$\Phi(x_i) = (x_{i1}^2, x_{i2}^2, 1.414 x_{i1} x_{i2})$$

$$K(x, x_i) = \Phi(x) \cdot \Phi(x_i)$$

SVM maps two-dimensional input space to three-dimensional feature space that is shown in Fig.2

## 2. a) SVM principle:

It is one of the popular classification models. Support vectors are used to find hyper plane between two classes. Maps input space to higher dimensional feature space. Support vectors are the training samples that define the optimal separating hyper plane. Support vectors are close to the hyper plane.

## 2. b) SVM modeling for Speaker turn point detection:

The overlapping frame is calculated manually.

Each frame in LPCC (19<sup>th</sup> order) must ends with +1 or -1. +1 indicates the overlapping frame (Speaker turn point points). -1 indicates the non-overlapping frame.

### • SVM Training

The manually calculated overlapping frames are appended with +1 and non-overlapping frames are appended with -1 using SVM Torch.

### • SVM Testing

Test conversation LPCC values are given to SVM test .The result is stored in the result.dat.

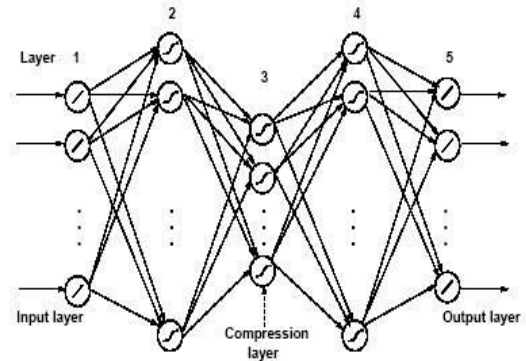
The result.dat file contains positive (+1) and negative (-1) values. Positive values indicate overlapping frames in the conversation file. That is, Speaker changes points.

## 3. AANN MODEL FOR CAPTURING THE DISTRIBUTION OF ACOUSTIC FEATURE VECTORS

Autoassociative neural network models are feed forward neural networks performing an identity mapping of the input space, and are used to capture the distribution of the input data [13].

A five layer autoassociative neural network model, as shown in Figure 3, is used to capture the distribution of the feature vectors in our study. The second and fourth layers of the network have more units than the input layer. The third layer has fewer units the first or fifth. The processing units in the first and third hidden layers are nonlinear, and the units in the second compression/hidden layer can be linear or nonlinear

The structure of the AANN model used in our study is 19L 38N 5N 38N 19L, where L denotes a linear and N denotes a nonlinear units. The nonlinear output function for each unit is  $\tanh(s)$ , where  $s$  is the activation value of the unit. The standard back propagation learning algorithm is used to adjust the weights of the network to minimize the mean square error for each feature vector. As the error between the actual and the desired output vectors is minimized, the cluster of points in the input space determines the shape of the hyper surface obtained by the projection onto the lower dimensional space. The AANN captures the distribution of the input data



**Fig. 3:** A five layer AANN model

In order to visualize the distribution capturing ability, one can plot the error for each input data point in the form of some probability surface. The error  $e_i$  for the data point  $i$  in the input space is plotted as  $p_i = \exp(-e_i/\alpha)$ , where  $\alpha$  is a constant. Note that  $p_i$  is not strictly a probability density function, but we call the resulting surface as probability surface. The plot of the probability surface shows large amplitude for smaller error  $e_i$  indicating better match of the network for that data point.

One can use the probability surface to study the characteristics of the distribution of the input data captured by the network. Ideally, one would like to achieve the best probability surface, best defined in terms of some measure corresponding to a low average error.

#### 4. THE PROPOSED SPEAKER TURN POINT DETECTION ALGORITHM

This paper purposes a novel Speaker turn point detection algorithms using AANN. The basic concept of the proposed method is illustrated in Figure 2. We begin with the assumption that there is a Speaker turn point located in the data stream at the centre of the analysis window under consideration. If the speech signal of this analysis window comes from different speakers, all the feature vectors in the right half of the window may not fall into the distribution of the feature vectors from the left half window. On the contrary, if the speech signal of this analysis window comes from only one speaker then the feature vectors is in the right half of the window falls into the distribution of feature vectors of the left half window.

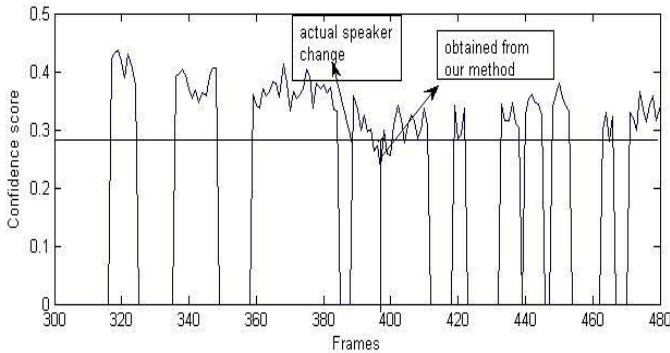
Given the speech feature vectors  $S = s_i$ ,  $i = 1, 2, \dots, n$  where  $i$  is the frame index and  $n$  is the total number of feature vectors in the speech signal. The proposed algorithm for detecting Speaker turn point is given below:

$m$  number of feature vectors ( $m \bmod 2 = 1$ ) are considered for  $k^{\text{th}}$  analysis window  $W_k$  and is given by

$$W_k = \{S_j\}, k \leq j < m + k \quad (15)$$

It is assumed that the Speaker turn point occurs at the middle feature vector ( $c$ ) of the analysis window.

$$c = k + \frac{m}{2} \quad (16)$$



**Fig 4 :** Basic concept of the proposed algorithm

- (1) We consider all the feature vectors in the analysis window  $W_k$  that are located left of  $c$  as left half window ( $Lk$ ).

$$Lk = \{S_j\}, k \leq j \leq c - 1 \quad (17)$$

Similarly, all the feature vectors that are located right of  $c$  is in right half window ( $Rk$ ).

$$Rk = \{S_j\}, c + 1 \leq j \leq m + k \quad (18)$$

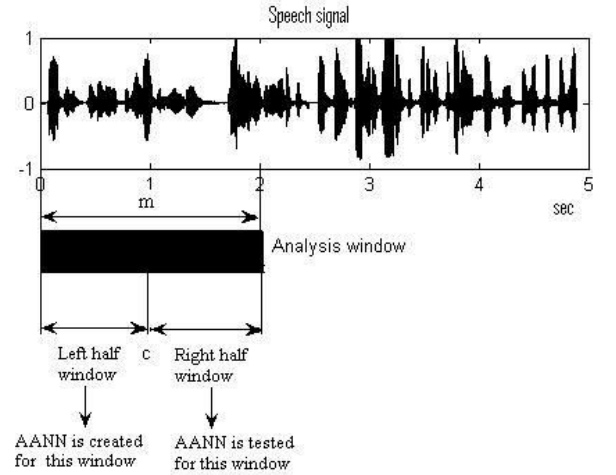
- (2) AANN is trained using the feature vectors in  $Lk$  and the model captures the distribution of this block of vectors. Then feature vectors in  $Rk$  are given as input to the AANN model and the output of model is compared with the input to compute the normalized squared error  $e_k$ . The normalized squared error ( $e_k$ ) for the feature vector  $y$  is given by

$$e_k = \frac{\|y - 0\|^2}{\|y\|^2} \quad (19)$$

where  $0$  is the output vector given by the model. The error  $e_k$  is transformed into a confidence score  $S$  using

$$S = \exp(e_k) \quad (20)$$

If true Speaker turn point occurs at  $c$ , then  $Lk$  and  $Rk$  will be from different speakers and the confidence score  $s$  for this  $c$  will be low. Likewise, if  $c$  is not the true Speaker turn point and both  $Lk$  and  $Rk$  are from the same speaker then



**Fig 5 :** Analysis window

the confidence score  $s$  will be high. The next possibility is either  $Lk$  or  $Rk$  may have the speech feature vectors from both the speakers. If this is the case, the confidence score  $s$  will be in between the above two values.

The value  $k$  is incremented by one and the steps from 1 to 4 are repeated until  $m + k$  reaches  $n$ .

It is not possible to obtain the same confidence score for all true Speaker turn points. The confidence score of Speaker turn point will be low when compared to the confidence scores of the frames on either side of the Speaker turn point. So the local minimum of the confidence score is considered instead of global minimum. To avoid the false alarms, the local minima which are less than the threshold value are considered. Hence, after obtaining the confidence score for the entire speech signal

the hypothesized Speaker turn point is validated by using a threshold. The threshold (t) is calculated from the

Window size(frames)	Missed Detection Rate (MDR)
5	5.13
10	8.99
15	11.23
20	3.24

confidence score as

**Table 1:** Performance of the Speaker turn point detection system at various stages

$$T = s \min + a s \min \quad (21)$$

Where s min is the global minimum confidence score and  $a$  is the adjustable parameter. The proposed method is unsupervised because it can detect the speaker changes without any knowledge of the identity of speakers and there is no need for training speaker models beforehand.

## 5. EXPERIMENTAL RESULTS

For our experiments on Speaker turn point detection, we use the extended data consist of two-speaker conversation. A total dataset of 9 conversations is used in our studies. This dataset includes three conversations for each of male-male, male-female and female-female speaker conversations. The Speaker turn point in the conversation is manually marked. The total dataset is divided in to training dataset a validation dataset and test dataset.

The speech data is processed using a frame size of 20 milliseconds. Each frame size is represented by a19 dimensional LPCC feature vector and MFCC feature vector. The speech data of the conversations in the training dataset is processed to obtain positive and negative examples for training the Speaker turn point detection SVM and AANN. The speech data of conversations in the validation dataset is used for obtaining the negative examples to train the false alarm reduction SVM. The speech data of the conversation in the test dataset is used for evaluating the performance of the Speaker turn point detection system.

The speech data of conversation is given as the input to the Speaker turn point detection system. The sliding window method is used to obtain hypotheses from the Speaker turn point detection SVM. The output of the SVM is smoothed to eliminate the short duration speaker turns. The hypotheses after removal of short speaker turns are processed by the false alarm reduction SVM to give the Speaker turn point detection. The Speaker turn point detection performance is measured as the missed detection rate (MDR) and the false alarm rate (FAR). The missed detection rate is defined as the ratio of the number of Speaker turn point missed (M) and the number of actual Speaker turn point points (A) are given in equation (22) and (23).

$$MDR = \frac{M}{A} \times 100 \quad (22)$$

$$FAR = \frac{F}{T-A} \times 100 \quad (23)$$

where F is the number of false hypotheses and T is the number of test patterns.

The MDR and FAR are determined at different stages of the Speaker turn point detection system.

The performance of the different window length is given in the Table 1.

The performance of speaker segmentation is assessed in terms of two types of error related to Speaker turn point detections namely false alarms and missed detections. A false alarm ( $\alpha$ ) of Speaker turn point detection occurs when a detected Speaker turn point is not a true one. A missed detection ( $\beta$ ) occurs when a true Speaker turn point cannot be detected. The false alarm rate ( $\alpha_r$ ) and detection rate ( $\beta_r$ ) are defined as [16], [17].

$$\alpha_r = \frac{\text{Number of false alarmed speaker changes}}{\text{Number of detected speaker changes}} \quad (24)$$

$$\beta_r = \frac{\text{Number of missed detection}}{\text{Number of true speaker changes}} \quad (25)$$

Two other measures namely precision ( $p$ ) and recall ( $r$ ) can also be used, which are closely related to  $\alpha_r$ ,  $\beta_r$  [14], [15]. They are defined as

$$p = \frac{\text{Number of correctly found speaker changes}}{\text{Total number of changes found}} \quad (26)$$

$$r = \frac{\text{Number of correctly found speaker changes}}{\text{number of actual speaker changes}} \quad (27)$$

In order to compare the performance of different systems, the f-measure is often used and is given by

$$f = 2 \frac{pr}{p+r} \quad (28)$$

The f – measure carries from 0 to 1, with a higher f – measure indicating better performance.

Table.2 shows Performance of the Speaker turn point detection system comparison (LPCC and MFCC). Using SVM and AANN

Classifier	$\alpha_r$	$\beta_r$	f-measure
AANN	15.77%	4.65%	89.03%
SVM	27.01%	25.11%	70.78%

## 6. CONCLUSION:

In this paper we have proposed an alternative method for speaker segmentation using LPCC and MFCC features and SVM and AANN.

The proposed algorithm can achieve effective speaker segmentation with data collection and it is capable of detecting speaker segments of shorten duration the algorithm can be applied for real time application and it does not require any prior knowledge about the speaker identity and their model. The proposed method was carried out for two speaker conversations and by using clean speech signals. The proposed algorithm shows that better performance using with SVM than AANN.

## REFERENCES

1. S. Jothilakshmi, S. Palanivel, V. Ramalingam, "Unsupervised Speaker Segmentation using Autoassociative Neural Network" 2010, International Journal of Computer Applications (0975 – 8887) Volume 1 – No. 7.
2. Shilei zhang shuwu zhang,Bo Xu. "A Two-Level Method For Unsupervised Speaker-based Audio Segmentation". The 18<sup>th</sup> International Conference on pattern Recognition(ICPR'06) IEEE, 2006.
3. A.Malegaonkar, A.Ariyaceinia, P.SivaKumaran and J.Fortuna. "Unsupervised Speaker turn point Detection Using Probabilistic Pattern Matching" IEEE SIGNAL PROCESSING, LETTERS, VOL.13, NO.8, AUGUST 2006.
4. Amit S.Maleganonkar, Aladdin M.Ariyaceinia and Perasiriyn Sivakumaran. "Efficient Speaker turn point detection using adapted Gaussian mixture models" IEEE TRANSACTIONS ON AUDIO, SPEECH AND LANGUAGE PROCESSING, VOLUME.15.No.6, Aug 2007.
5. Jitendra Ajmera,Iain Mccowan and Herve Bourslard, "Robust Speakers change Detection." IEEE SIGNAL PROCESSING LETTERS, VOL.2, NO.8, AUGUST 2004.
6. Andre G.Adami,Sachin S.Kajarekar, Hynek Hermansky, "A New Speaker turn point Detection Method For Two-Speaker Segmentation" IEEE, 2002.
7. Po-Chuan Lin, Jia-Chingwang, Jhing-Fa Wang and Hao-Ching Sung. "Unsupervised Speaker turn point detection using SVM Training Misclassification Rate" IEEE TRANSACTIONS ON COMPUTERS, VOL.56, No.9, Sep 2007.
8. Noureddine ELLOUZE. "Robustness Improvement Of Speaker Segmentation Techniques Based on the Bayesian Information Criterion", IEEE, 2006.
- 9.V.Kartik and D.Srikrishna Sathish and C.Chandra Sekar."Speaker turn point detection using Support Vector Mechines", "Speech and Vision Laboratory, Indian Institute of Technology Madras, Page No. 1-5, 2006.
- 10.Guillaume Lathoud Iain A.McCowan. "LOCATION BASED SPEAKER SEGMENTATION", IEEE, 2003.
- 11.Margarita Kotti, Emmanouil Benetos, Jaime S.Cardoso "Automatic Speaker Segmentation Using Multiple Features And Distance Measures: Comparison Of Three Approaches"1-4244-0367-7/06 IEEE, 2006.
- 12.V.Vapnik, "Statistical learning theory", John Wiley and Sons, New York, 1998.
- 13.B. Yegnanarayana, S. P. Kishore. 2002. AANN: " An alternative to GMM for pattern recognition Neural Networks." 15,459–469.
- 14.H. Kim, D. Elter, T. Sikora. 2005." Hybrid speaker based segmentation system using model level clustering". In Proceedings of the IEEE International conference on Acoust. Speech, Signal Processing (ICASSP 05). 745–748.
- 15.J. Ajmera, I. McCowan, H. Boursland. 2004. "Robust speaker change detection". IEEE Signal Process. Lett. 11, 8, 649–651.
16. P. Delacourt, C. J. Wellekens. 2000. DISTBIC: "A speaker based segmentation for audio data indexing". Speech comm 32,111–126.
17. S.Cheng, H. Wang. 2004. Metric SEQDAC: "A hybrid approach for audio segmentation". In Proceedings of the 8<sup>th</sup> International conference on spoken language processing 1617– 1620.

# Consistent Data Delivery in Mobile Adhoc Networks

R.Nanda Kumar<sup>#1</sup>, Dr.Sankara Malliga G<sup>#2</sup>

Research Scholar, Department of Information and Communication Engineering<sup>#1</sup>

Professor and Head, Department of Electronics and Communication Engineering<sup>#2</sup>

Oasys Institute of technology, Trichy, Tamil Nadu, India. <sup>#1</sup>

Anand Institute of Higher Technology, Chennai, Tamil Nadu, India. <sup>#2</sup>

nanthartr@gmail.com<sup>#1</sup>, sankaramalligag.ece@aiht.ac.in<sup>#2</sup>

**Abstract-** Mobile adhoc Networks (MANETs) technologies have been gradually impacting almost all spheres of our lives. Recently, the Mobile adhoc Networks (MANETs) are gaining increased attention for generating extensive wireless communication. Routing protocols for ad hoc wireless networks must be able to perform efficient and effective mobility management. Most existing ad hoc routing protocols are susceptible to node mobility, especially for large-scale networks. Also in Mobile adhoc network, the node mobility and the location update interval are main factors leading to packet forwarding failure due to the receiver moving from one position to another. Focused this concern Spot based Opportunistic Routing Protocol (SOR) has been proposed to face issues of the high node mobility. This proposed scheme is designed for the dynamic nature of MANETs associated with various constraints. SOR routing protocol which takes advantage of find an effective routing which is used to transmit information from source to destination across the whole network topology. End of the results shows that SOR achieves excellent performance even under high node mobility with acceptable overhead.

**Keywords** — MANETs, SOR, Mobility, Data Delivery, Location Update Interval.

## I. INTRODUCTION

Realizing the necessity of open standards in this emerging area of computer communication, a working group within the Internet Engineering Task Force (IETF), termed the name Mobile ad hoc networks (MANETs) working group was formed to standardize the routing protocols and functional specifications of the Mobile adhoc networks. Mobile adhoc networks routing functionality is to mainly support for the self-organizing mobile networking infrastructure. Even though ad hoc networks are expected to work in the absence of any fixed infrastructure. Mobile Ad hoc networks are defined as the category of wireless networks that utilize multi-hop radio relaying and are capable of operating without the support of any fixed infrastructure. MANETs are self-organizing networks and they have been made up of mobile nodes, which are using their neighbors as a mean of communication with other nodes in the network.

Existing routing protocols in ad-hoc networks utilize the single route that is built for source and destination node pair. Due to node mobility, node failures and the dynamic characteristics of the radio channel, links in a route may become temporarily unavailable, making the route invalid.

The overhead of finding alternative routes mounts along with additional packet delivery delay. Mobile Adhoc networks change their topology, expressed by the node connectivity over time, as the nodes change their position in space. [1] Mobile ad hoc networks are characterized by dynamic topology due to node mobility, limited channel bandwidth and limited battery power of nodes.

The key challenge here is to be able to route with low overheads even in dynamic conditions [2]. Adhoc mobile networks are very dynamic, self-organizing, self-healing distributed networks which support data networking without an infrastructure. [4] Mobile Ad hoc Networks consists of a set of wireless mobile nodes communicating to each other without any centralized control or fixed network infrastructure and can be deployed quickly [5]. Adhoc mobile networks are very dynamic, self-organizing, and self-healing distributed networks which support data networking without an infrastructure. Due to lack of trusted nodes, Mobile adhoc networks require specialized authenticated protocol [6]. Mobile Adhoc Network is an autonomous system consisting of a set of mobile hosts that are free to move without the need for a wired backbone or a fixed base station. [7] Mobile Adhoc Networks having self-organizing and self-configuring network without the need of any centralized base station and physical connections of mobile devices. Mobile ad-hoc network has no fixed topology due to mobility of nodes, interference, path loss and multipath propagation. Mobile Ad hoc Networks is a robust infrastructure less wireless network having mobile nodes. A MANETs can be created either by mobile nodes or by both static and dynamic mobile nodes. A mobile node has arbitrarily associated with each other forming uniformed topologies. They serve up as both routers and hosts. [8] A Mobile Adhoc Network (MANETs) is a collection of autonomous wireless mobile nodes forming frequently changing network topology. Which nodes can communicate with each other through wireless links that needs efficient this is enough for into dynamic routing protocols [9].

To summarize, this paper is organized as follows. Section II provides details of the various classifications of Routing protocols, Section III Existing MANETs Routing Protocol, Section IV Literature Survey, V. Proposed System, and the conclusion of the paper in Section VI.

## II. CLASSIFICATION ON ROUTING PROTOCOLS

The routing protocols in Mobile adhoc networks can be broadly classified into four categories. They are

- Routing information update mechanism
- Use of temporal information for routing
- Routing topology
- Utilization of specific resources

### A. Based on the Routing Information Update Mechanism

Ad hoc wireless network routing protocols can be classified into three major categories as follows based on the routing information update mechanism.

- Proactive or table-driven routing protocols
- Reactive or on-demand routing protocols
- Hybrid routing protocols

### B. Based on the Use of Temporal Information for Routing

This classification of routing protocols is based on the use of temporal information used for routing. Since ad hoc wireless networks are highly dynamic and path breaks are much more frequent than in wired networks, the use of temporal information regarding the lifetime of the wireless links and the lifetime of the paths selected assumes significance. The protocols that fall under this category can be listed as follows.

- Routing protocols using past temporal information
- Routing protocols that use future temporal information

### C. Based on the Routing Topology

Routing topology being used in the Internet is hierarchical in order to reduce the state information maintained at the core routers. Ad hoc wireless networks, due to their relatively smaller number of nodes, can make use of either a flat topology or a hierarchical topology for routing.

- Flat topology routing protocols
- Hierarchical topology routing protocols

### D. Based on the Utilization of Specific Resources:

Based on the Utilization of specific resources in MANETs has been classified in the following ways. They are listed below.

- Power-aware routing
- Geographical information assisted routing

## III. EXISTING MANET'S ROUTING PROTOCOLS

Reactive or On-Demand Routing protocols execute the path finding process and exchange routing information only when a path is required by a node to communicate with a destination. Ad hoc on-demand distance vector (AODV) and dynamic source routing (DSR) are well-known examples of Reactive routing protocols. This section explores some of the existing on-demand routing protocols in detail.

- Dynamic source routing protocol (DSR)
- Ad hoc on-demand distance vector (AODV)
- Location-aided routing protocol (LAR)

### A. Dynamic Source Routing Protocol

Dynamic source routing protocol (DSR) is an on-demand protocol designed to restrict the bandwidth consumed by control packets in ad hoc wireless networks by eliminating the periodic table-update messages required in the table-driven approach. The major difference between this and the other on-demand routing protocols is that it is beacon-less and hence does not require periodic hello packet (beacon) transmissions, which are used by a node to inform its neighbors of its presence. The basic approach of this protocol (and all other on-demand routing protocols) during the route construction phase is to establish a route by flooding RouteRequest packets in the network. The destination node, on receiving a RouteRequest packet, responds by sending a RouteReply packet back to the source, which carries the route traversed by the RouteRequest packet received.

Consider a source node that does not have a route to the destination. When it has data packets to be sent to that destination, it initiates a RouteRequest packet. This RouteRequest is flooded throughout the network. Each node, upon receiving a RouteRequest packet, rebroadcasts the packet to its neighbors if it has not forwarded already or if the node is not the destination node, provided the packet's time to live (TTL) counter has not exceeded.

Each RouteRequest carries a sequence number generated by the source node and the path it has traversed. A node, upon receiving a RouteRequest packet, checks the sequence number on the packet before forwarding it. The packet is forwarded only if it is not a duplicate RouteRequest. The sequence number on the packet is used to prevent loop formations and to avoid multiple transmissions of the same RouteRequest by an intermediate node that receives it through multiple paths. Thus, all nodes except the destination forward a RouteRequest packet during the route construction phase. A destination node, after receiving the first RouteRequest packet, replies to the source node through the reverse path the RouteRequest packet had traversed. This protocol uses Route cache that stores all possible information extracted from the source route contained in a data packet. Also which is used at the time of routing construction phase.

#### Advantages

DSR protocol uses the reactive approach which eliminates the need to periodically flood the network with table update messages. This on-demand routing approach route is established only when it is required and hence the need to find routes to all other nodes in the network as required. The intermediate nodes also utilize the route cache information efficiently to reduce the control overhead.

#### Disadvantages

Drawbacks of DSR protocol is that the route maintenance mechanism does not locally repair a broken link. Even though this protocol performs well in static and low-mobility environments, the performance degrades rapidly with increasing mobility. Also, considerable routing overhead is involved due to the source-routing mechanism employed in DSR. This routing overhead is directly proportional to the path length.

### B. *Adhoc On-Demand Distance-Vector Routing Protocol*

Ad hoc on-demand distance vector (AODV) routing protocol uses an on demand approach for finding routes, that is, a route is established only when it is required by a source node for transmitting data packets. It employs destination sequence numbers to identify the most recent path. In AODV, the source node and the intermediate nodes store the next-hop information corresponding to each flow for data packet transmission.

In an on demand routing protocol, the source node floods the RouteRequest packet in the network when a route is not available for the desired destination. It may obtain multiple routes to different destinations from a single RouteRequest. The major difference between AODV and other on-demand routing protocols is that it uses a destination sequence number (DestSeqNum) to determine an up-to-date path to the destination. A node updates its path information only if the DestSeqNum of the current packet received is greater than the last DestSeqNum stored at the node.

A RouteRequest carries the source identifier (SrcID), the destination identifier (DestID), the source sequence number (SrcSeqNum), the destination sequence number (DestSeqNum), the broadcast identifier (BcastID), and the time to live (TTL) field. DestSeqNum indicates the freshness of the route that is accepted by the source. When an intermediate node receives a RouteRequest, it either forwards it or prepares a RouteReply if it has a valid route to the destination. The validity of a route at the intermediate node is determined by comparing the sequence number at the intermediate node with the destination sequence number in the RouteRequest packet. If a RouteRequest is received multiple times, which is indicated by the BcastID-SrcID pair, the duplicate copies are discarded. All intermediate nodes having valid routes to the destination, or the destination node itself, are allowed to send RouteReply packets to the source. Every intermediate node, while forwarding a RouteRequest, enters the previous node address and its BcastID. A timer is used to delete this entry in case a RouteReply is not received before the timer expires. This helps in storing an active path at the intermediate node as AODV does not employ source routing of data packets. When a node receives a RouteReply packet, information about the previous node from which the packet was received is also stored in order to forward the data packet to this next node as the next hop toward the destination.

#### Advantages

The main advantage of this protocol is that routes are established on demand and destination sequence numbers are used to find the latest route to the destination. The connection setup delay is less.

#### Disadvantages

One of the disadvantage of AODV protocol is that intermediate nodes can lead to inconsistent routes if the source sequence number is very old and the intermediate nodes have a higher but not the latest destination sequence number, thereby having stale entries and also multiple *RouteReply* packets in response to a single *RouteRequest* packet can lead to heavy control overhead.

### C. Location-aided routing protocol (LAR)

LAR reduces the control overhead by limiting the search area for finding a path. The efficient use of geographical position information, reduced control overhead, and increased utilization of bandwidth are the major advantages of this Location-aided routing protocol. The applicability of this protocol depends heavily on the availability of GPS infrastructure or similar sources of location information. Hence, this protocol cannot be used in situations where there is no access to such information.

## IV. LITERATURE SURVEY

### A. *A performance Comparison of Multi-hop Wireless ad hoc network Routing Protocols [1]*

Traditional routing algorithms prove to be inefficient in such a changing environment. Ad-hoc routing protocols such as dynamic source routing (DSR), ad-hoc on-demand distance vector routing (AODV) and destination-sequence distance vector (DSDV) have been proposed to solve the multi hop routing problem in ad-hoc networks. Performance studies of these routing protocols have assumed constant bit rate (CBR) traffic. Real-time multimedia traffic generated by video-on demand and teleconferencing services are mostly variable bit rate (VBR) traffic. Most of this multimedia traffic is encoded using the MPEG standard. When video traffic is transferred over MANETs a series of performance issues arise. This paper presents a performance comparison of three ad-hoc routing protocols - DSR, AODV and DSDV when streaming MPEG4 traffic. Simulation studies show that DSDV performs better than AODV and DSR. However all three protocols fail to provide good performance in large, highly mobile network environments.

### B. *A Survey on Position-based Routing in Mobile ad hoc networks [2]*

A survey on position-based routing in mobile ad hoc networks paper presents an overview of ad hoc routing protocols that make forwarding decisions based on the geographical position of a packet's destination. Other than the destination's position, each node need know only its own position and the position of its one-hop neighbors in order to forward packets. Since it is not necessary to maintain explicit routes, position-based routing does scale well even if the network is highly dynamic. The main prerequisite for position-based routing is that a sender can obtain the current position of the destination.

### C. *A framework for reliable routing in mobile ad hoc networks [21]*

Mobile ad hoc networks consist of nodes that are often vulnerable to failure. As such, it is important to provide redundancy in terms of providing multiple node-disjoint paths from a source to a destination. This papers propose a modified version of the popular AODV protocol that allows us to discover multiple node-disjoint paths from a source to a destination.

From that can conclude that it is necessary to place for call reliable nodes in the network for efficient operations. The proposed method of a deployment strategy that determines the positions and the trajectories of these reliable nodes such that achieve a framework for reliably routing information.

#### D. Survey on Opportunistic Routing in Multihop Wireless Networks [10]

The study of Opportunistic routing is based on the use of broadcast transmissions to expand the potential forwarders that can assist in the retransmission of the data packets. The receptors need to be coordinated in order to avoid duplicated transmissions. This is could be achieved by ordering the forwarding nodes and in the position-based packet forwarding strategies. This proposed Opportunistic routing protocols differ in the criterion to order the receptors and the way of the receptors coordinate. This paper presents a survey of the most significant opportunistic routing protocols for multihop wireless networks.

#### E. ExOR: Opportunistic MultiHop Routing for Wireless Networks [6]

This paper describes ExOR, an integrated routing and MAC protocol that increases the throughput of large unicast transfers in multi-hop wireless networks. ExOR chooses each hop of a packet's route after the transmission for that hop, so that the choice can react which intermediate nodes actually received the transmission. This deferred choice gives each transmission multiple opportunities to make progress. As a result ExOR can use long radio links with high loss rates, which would be avoided by traditional routing. ExOR increases a connection's throughput while using no more network capacity than traditional routing. ExOR's design faces the following challenges. The nodes that receive each packet must agree on their identities and choose one forwarder. The agreement protocol must have low overhead, but must also be robust enough that it rarely forwards a packet zero times or more than once. Finally, ExOR must choose the forwarder with the lowest remaining cost to the ultimate destination. For pairs between which traditional routing uses one or two hops, ExOR's robust acknowledgments prevent unnecessary retransmissions, increasing throughput by nearly 35%. For more distant pairs, ExOR takes advantage of the choice of forwarders to provide throughput gains of a factor of two to four.

### V. PROPOSED SYSTEM

The design process of Spot based Opportunistic Routing Protocol (SOR) is based on geographic routing and opportunistic forwarding. The nodes are assumed to be aware of their own locations and the positions of their direct neighbors. In SOR Routing Protocol choosing the candidates and assigning priority among nodes for the forwarding candidates plays a major role while designing this routing protocol. Actually candidate means the node assigned as next hop which is selected among all nodes in the direction of forwarding region. The forwarding area has been determined by sender and the next hop node.

#### A. SOR Protocol Design

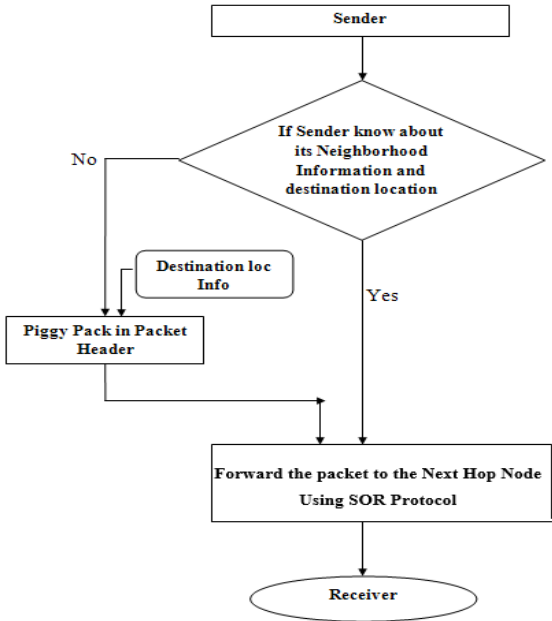


Figure 1 – SOR Protocol Designing Process flow

Spot based Opportunistic Routing Protocol (SOR) Designing process flow is shown in Figure 1. It takes the responsibility to solve the problems in the candidate selection and gives the priority to forwarding candidates. Also only the nodes located in the direction of forwarding region may get the chance to be backup nodes. Focused this concern an SOR Protocol has been designed in the following ways.

- Location Information adding in Packet header
- Candidate node selection
- Distance calculation between Each Node
- Priority assigning for the forwarding candidates
- Collects the Neighbor node List details
- Receiver details

From the above module details lists Candidate selection process has been primarily taken while designing SOR Protocol.

#### B. Candidate Selection Process in SOR protocol

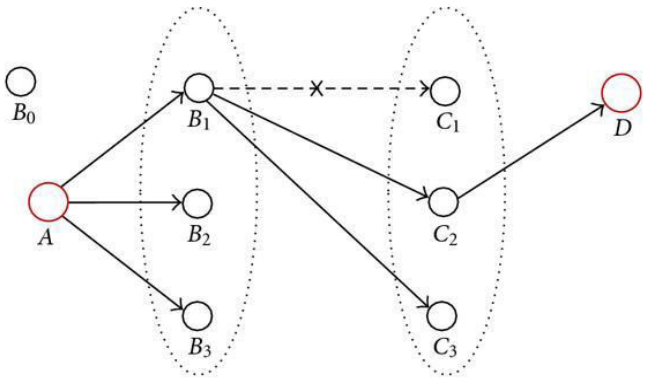


Figure 2: Candidate Selection process in SOR Protocol

Candidate selection process in Spot based Opportunistic Routing Protocol (SOR) is shown in the above Figure 2. Actual determination of the forwarding area is done by the sender and the next hop node. The node located in the forwarding area satisfies the following conditions. First one is that it might be present in the positive progress towards the destination. Next one important condition is that its distance to the next hop node should not exceed half of the transmission range of a wireless node. Based on that in above figure 2 the nodes from the A to D positive progress may take from the nodes B1, C1 and C2. Likewise the candidate process has been done by the SOR Routing Protocol.

## VI. CONCLUSION

This paper addresses the problem of reliable data delivery in highly dynamic mobile ad hoc networks. Constantly changing network topology makes conventional ad hoc routing protocols incapable of providing satisfactory performance. In the face of frequent link break due to node mobility, substantial data packets would either get lost, or experience long latency before restoration of connectivity. The efficacy of the involvement of forwarding candidates against node mobility, as well as the overhead due to opportunistic forwarding is analysed. Through implementation, confirm the effectiveness and efficiency of SOR high packet delivery ratio may achieve while the delay and duplication are the lowest.

## REFERENCES

- [1] S.J. Broch, D. A. Maltz, D. B. Johnson, Y.-C. Hu, and J. Jetcheva, "Performance comparison of multi-hop wireless ad hoc network routing protocols," in *MobiCom '98*. ACM, 1998, pp. 85–97.
- [2] M. Mauve, A. Widmer, and H. Hartenstein, "A survey on position-based routing in mobile ad hoc networks," *Network*, IEEE, vol. 15, no. 6, pp. 30–39, Nov/Dec 2001.
- [3] D. Chen and P. Varshney, "A survey of void handling techniques for geographic routing in wireless networks," *Communications Surveys & Tutorials*, IEEE, vol. 9, no. 1, pp. 50–67, Quarter 2007.
- [4] D. Son, A. Helmy, and B. Krishnamachari, "The effect of mobility induced location errors on geographic routing in mobile ad hoc sensor networks: analysis and improvement using mobility prediction," *Mobile Computing*, IEEE Transactions on, vol. 3, no. 3, pp. 233–245, July -Aug. 2004.
- [5] B. Karp and H. T. Kung, "Gpsr: greedy perimeter stateless routing for wireless networks," in *MobiCom '00*, 2000, pp. 243–254.
- [6] S. Biswas and R. Morris, "Exor: opportunistic multi-hop routing for wireless networks," in *SIGCOMM '05*, 2005, pp. 133–144.
- [7] S. Chachulski, M. Jennings, S. Katti, and D. Katabi, "Trading structure for randomness in wireless opportunistic routing," in *SIGCOMM '07*, 2007, pp. 169–180.
- [8] E. Rozner, J. Seshadri, Y. Mehta, and L. Qiu, "Soar: Simple opportunistic adaptive routing protocol for wireless mesh networks," *Mobile Computing*, IEEE Transactions on, vol. 8, no. 12, pp. 1622–1635, dec. 2009.
- [9] A. Balasubramanian, R. Mahajan, A. Venkataramani, B. N. Levine, and J. Zahorjan, "Interactive wifi connectivity for moving vehicles," in *SIGCOMM '08*, 2008, pp. 427–438.
- [10] K. Zeng, Z. Yang, and W. Lou, "Location-aided opportunistic forwarding in multirate and multihop wireless networks," *Vehicular Technology*, IEEE Transactions on, vol. 58, no. 6, pp. 3032–3040, July 2009.
- [11] A. Tsirigos and Z. Haas, "Analysis of multipath routing, part 2: mitigation of the effects of frequently changing network topologies," *Wireless Communications*, IEEE Transactions on, vol. 3, no. 2 pp. 500–511, March 2004.
- [12] X. Huang, H. Zhai, and Y. Fang, "Robust cooperative routing protocol in mobile wireless sensor networks," *Wireless Communications, IEEE Transactions on*, vol. 7, no. 12, pp. 5278–5285, December 2008.
- [13] A. Tsirigos and Z. Haas, "Analysis of multipath routing-part i: the effect on the packet delivery ratio," *Wireless Communications IEEE Transactions on*, vol. 3, no. 1, pp. 138–146, Jan. 2004.
- [14] B. Deb, S. Bhatnagar, and B. Nath, "Reinform: reliable information forwarding using multiple paths in sensor networks," in *Local Computer Networks*, 2003. LCN '03, Oct. 2003, pp. 406–415.
- [15] N. Arad and Y. Shavitt, "Minimizing recovery state in geographic ad hoc routing," *Mobile Computing*, IEEE Transactions on, vol. 8, no. 2, pp. 203–217, Feb. 2009.
- [16] Y. Han, R. La, A. Makowski, and S. Lee, "Distribution of path durations in mobile ad-hoc networks: palms theorem to the rescue," *Computer Networks*, vol. 50, no. 12, pp. 1887–1900, 2006.
- [17] R. Groenevelt, "Stochastic models for mobile ad hoc networks," Ph.D. dissertation, Universite de Nice, Sophia Antipolis, France, 2005.
- [18] W. Navidi and T. Camp, "Stationary distributions for the random waypoint mobility model," *Mobile Computing*, IEEE Transactions on, vol. 3, no. 1, pp. 99–108, Jan-Feb 2004.
- [19] "The network simulator ns-2," <http://www.isi.edu/nsnam/ns/>.
- [20] E. Felemban, C.-G. Lee, E. Ekici, R. Boder, and S. Vural, "Probabilistic qos guarantee in reliability and timeliness domains in Wireless sensor networks," in *INFOCOM 2005*, vol. 4, March 2005, pp. 2646–2657 vol. 4.
- [21] S. Mueller, R. Tsang, and D. Ghosal, "Multipath routing in mobile ad hoc networks: Issues and challenges," *Lecture Notes in Computer Science*, vol. 2965, pp. 209–234, 2004.
- [22] A. Tsirigos and Z. Haas, "Analysis of multipath routing, part 2: mitigation of the effects of frequently changing network topologies," *Wireless Communications*, IEEE Transactions on, vol. 3, no. 2 pp. 500–511, March 2004.

# DATA WEB FOR QUERY FORMULATION LANGUAGE

\*Dr. G. Silambarasan, \*\* Dr. V. Chandrasekar,

\*Assistant Professor, Dept. of Computer Science and Engineering,  
The Kavery College of Engineering, Salem, Tamilnadu, India,

\*\*Associate Professor, Dept. of Computer Science and Engineering,

Malla Reddy College of Engineering and Technology, Secunderabad, Telangana State, India,

\*\*[drchaTndru86@gmail.com](mailto:drchaTndru86@gmail.com), \* [gssilambarasan@gmail.com](mailto:gssilambarasan@gmail.com)

## ABSTRACT

This trend of structured data on the web (Data web) is shifting the focus of web technologies toward new paradigms of structured-data retrieval. Traditional search engines cannot serve such data as the results of a keyword-based. Query will not be precise or clean, because the query itself is still ambiguous although the underlying data are structured. To expose the massive amount of structured data on the web to its full potential, people should be able to query these data easily and effectively. We present a query formulation language (called MashQL) in order to easily query and fuse structured data on the web. The main novelty of MashQL is that it allows people with limited IT skills to explore and query one (or multiple) data sources without prior knowledge about the schema, structure, vocabulary, or any technical details of these sources.

**Keyword:** Web, RDF, SQL, SPARQL

## INTRODUCTION AND MOTIVATION

In this short article we propose a data mashup approach in a graphical and Yahoo Pipes' style. This research is still a work in progress, thus please refer to [13] for the latest findings.

In parallel to the continuous development of the hypertext web, we are witnessing a rapid emergence of the Data

Web. Not only is the amount of social metadata increasing, but also many

Companies (e.g., Google Base, Upcoming, Flickr, eBay, Amazon, and others) started to make their content freely accessible through APIs. Many others (see [linkeddata.org](http://linkeddata.org)) are also making their content directly accessible in RDF and in a linked manner [3]. We are also witnessing the launch of RDFa, which allows people to access and consume HTML pages as structured data sources

This trend of structured and linked data is shifting the focus of web technologies towards new paradigms of structured-data retrieval. Traditional search engines cannot serve such data because their core design is based on keyword-search over unstructured data. For example, imagine how would be the results when using Google to search a database of job vacancies, say “well-paid research-oriented job in Europe”. The results will not be precise or clean, because the query itself is still ambiguous although the underlying data is structured. People are demanding to not only retrieve job links but also want to know the starting date, salary, location, and may render the results on a map.

Web 2.0 mashups are a first step in this direction. A mashup is a web application that consumes data originated from third parties and retrieved via APIs. For example, one can build a mashup that retrieves only well-paid vacancies from Google Base and mix it with similar vacancies from LinkedIn. The problem is that building mashups is an art that is limited to skilled programmers. Although some mashup editors have been proposed by the Web 2.0 community to simplify this art (such as Google Mashups, Microsoft’s

Popfly, IBM’s sMash, and Yahoo Pipes), however, what can be achieved by these editors is limited. They only focus on providing encapsulated access to some APIs, and still require programming skills. In other words, these mashup methods are motivating for -rather than solving- the problem of structured-data retrieval. To expose the massive amount of structured data to its full potential, people should be able to query and mash up this data easily and effectively.

## RELATED WORK

Several approaches have been proposed by the DB community to query structured data sources, such as query-by-example [23] and conceptual queries [4, 6, 17]. However, none of these approaches was used by casual users. This is because they still assume knowledge about the relational/conceptual schema. Among these, we found ConQuer [4] has some nice features, especially the tree structure of queries, but it also assumes one to start from the schema. In the natural language processing community, it has been proposed to allow people to write queries as natural language sentences, and then translate these sentences into a formal language (SQL [15] or XQuery [16]). However, these

approaches are challenged with the language ambiguity and the “free mapping” between sentences and data schemes.

This topic started to receive a high importance within the Semantic Web community. Several approaches (GRQL [1], iSPARQL [11], NITELIGHT [19] and RDFAuthor [18]) are

Proposing to represent triple patterns graphically as ellipses connected with arrows. However, these approaches assume advanced knowledge of RDF and SPARQL. Other approaches use Visual Scripting Languages (e.g., SPARQLMotion [21] and DeriPipes [22]), by visualizing links between query modules; but a query module merely is a window containing a SPARQL script in a textual form. These approaches are inspired by some industrial mashup editors such as Popfly, sMash, and Yahoo Pipes. These industry editors provide a nice visualization of APIs’ interfaces and some operators between them. However, when a user needs to express a query over structured data, she needs to use the formal language of that editor, such as YQL for Yahoo Pipes. Although MashQL visualizes links between query modules, similar to Yahoo Pipes and other Mashup editors,

## THE MASHQL LANGUAGE

The main goal of MashQL is to allow people to mash up and fuse data sources easily. In the background MashQL queries are automatically translated into and executed as SPARQL queries. Without prior knowledge about a data source, one can navigate this source and fuse it with another source easily. To allow people to build on each other’s results MashQL supports query pipes as a built-in concept. The example below shows two web data sources and a SPARQL query to retrieve “the book titles authored by Lara and published after 2007”. The same query in MashQL. The first module specifies the query input, and the second module specifies the query body. The output can be piped into a third module (not shown here), which renders the results into a certain format (such as HTML, XML or CSV), or as RDF input to other queries. Notice that in this way, one can easily build a query to fuse the content of two sources in a linked manner [3].

## AMBIGUITY

The main problem is that this approach is fundamentally bounded with the language ambiguity multiple meanings of terms and the mapping between these terms and the elements of a data schema allows people

with IT-skills to explore and query one or multiple data sources with prior knowledge about the schema[7], structure, vocabulary, or any technical details of these sources. We do not assume that a data source should have -an offline or inline- schema. This poses several language-design and performance complexities that we fundamentally tackle. The rapid growth of structured data on the Web has created a high demand for making this content more reusable and consumable.

- There are several approaches to solve this problem and hence different solutions exist in the literature.
- We also observed that people are still not used with the Data Web paradigm (i.e., dealing with structured data and the difficulty of querying it).

## QUERY FORMULATION ALGORITHM

We present a novel query formulation algorithm, by which the complexity and the responsibility of understanding a data source (even if it is schema free) are moved from the user to the query editor. It allows end users to easily navigate and query [8] an

unknown data graph(s). . We addressed the challenge of achieving interactive performance during query formulation by introducing a new approach for indexing RDF data. We presented two different implementation scenarios of MashQL and evaluated our implementation on two large datasets. Furthermore, we plan to use our approach on keyword-search.

It saves the time and the user spending cost.

- It is allow the user to dynamically create a new file through the web.
- MashQL can be similarly used for querying relational databases and XML.
- MashQL can be used to query and mash up the Data Web as simple as filtering and piping web feeds.

## CONCLUSION AND FUTURE

We plan to extend this work in several directions. We will introduce a search box on top of MashQL to allow keyword search and then use MashQL to filter the retrieved results. To allow people use MashQL in a typical data integration scenario, several reasoning services will be supported, including Same As, Subtype, Sub property

## REFERENCES

- [1]. Athanasis N, Christophides V, Kotzinos D: Generating On theFly Queries for the Semantic Web. ISWC (2004)
- [2]. Abiteboul S, Duschkal O: Complexity of Answering QueriesUsing Materialized Views. ACM SIGACT-SIGMODSIGART.(1998)
- [3]. Bizer C, Heath T, Berners-Lee T:Linked Data: Principles andState of the Art. WWW (2008)
- [4]. Bloesch A, Halpin, T: Conceptual Queries using ConQuer–II.(1997)
- [5]. Chong E, Das S, Eadon G, Srinivasan J: An efficient SQLbasedRDF querying scheme. VLDB (2005)
- [6]. Czejdo B, and Elmasri R, and Rusinkiewicz M, and Embley D:An algebraic language for graphical query formulation usingan EER model. Computer Science conference. ACM. (1987)
- [7]. Deng Y, Hung E, Subrahmanian VS: Maintaining RDF views.Tech. Rep CS-TR-4612 University of Maryland. 2004
- [8]. Ennals R, Garofalakis M: MashMaker: mashups for themasses. SIGMOD Conference 2007:
- [9]. Goldman R, Widom J: DataGuides: Enabling QueryFormulation and Optimization in Semistructured Databases.VLDB (1997)
- [10]. Hofstede A, Proper H, and Weide T: Computer SupportedQuery Formulation in an Evolving Context. Australasian DB Conf. (1995)
- [11]. <http://demo.openlinksw.com/isparql> (Feb. 2009)
- [12]. Jarrar M, Dikaiakos: MashQL: A Query-by-Diagram ToppingSPARQL. Proceedings of ONISW'08 workshop. (2008).
- [13]. Jarrar M, Dikaiakos M: A query-by-diagram language(MashQL). Technical Article TAR200805. University of Cyprus, 2008.  
<http://www.cs.ucy.ac.cy/~mjarrar/JD08.pdf>
- [14]. Kaufmann E, Bernstein A: How Useful Are Natural LanguageInterfaces to the Semantic Web for Casual End-Users. ISWC (2007)
- [15]. Li Y, Yang H, Jagadish H: NaLIX: An interactive natural language interface for querying XML. SIGMOD (2005)
- [16]. Popescu A, Etzioni O, Kautz H: Towards a theory of naturallanguage interfaces to databases. 8th Con on Intelligent user interfaces. (2003)
- [17]. Parent C, Spaccapietra S: About Complex Entities, Complex Objects and Object-Oriented Data Models. Info. System Concepts(1989)

- [18].<http://rdfweb.org/people/damian/RDFA>  
author (Jan. 2009)
- [19]. Russell A, Smart R, Braines D,  
Shadbolt R.: NITELIGHT: A Graphical  
Tool for Semantic Query Construction. The  
Semantic Web User Interaction Workshop.  
(2008)
- [20].[http://esw.w3.org/topic/SPARQL/Exten](http://esw.w3.org/topic/SPARQL/Extensions?)  
sions? (Feb. 2009)
- [21].[http://www.topquadrant.com/sparqlmot](http://www.topquadrant.com/sparqlmotion)  
ion (Feb. 2009)
- [22]. Tummarello G, Polleres A, Morbidoni  
C: Who the FOAFknows Alice? A needed  
step toward Semantic Web Pipes.ISWC WS.  
(2007)
- [23]. Zloof M: Query-by-Example:a Data  
Base Language. IBMSystems Journal, 16(4).  
(1977)
- [24]. Zhuge Y, Garcia-Molina H, Hammer J,  
Widom J: ViewMaintenance in a  
Warehousing Environment. SIGMOD  
(1995)

# MOVING HUMAN ACTION RECOGNITION AND IMAGE CLASSIFICATION

<sup>1</sup>Dr. Manikandan.P  
mani.p.mk@gmail.com

<sup>2</sup>Dr.S.P. Anandaraj  
anandsofttech@gmail.com

<sup>1& 2</sup> Professor / CSE Department,  
Malla Reddy Engineering College for Women,

**ABSTRACT:** The proposed unified people localize, and label their activities, in complex short-duration video sequences. Image classification is one of classical problems of concern in image processing. There are various approaches for solving this problem. The aim of this paper is bring together areas in which are Artificial Neural Network (ANN) applying for image classification. we separate the image into many sub-images based on the features of images. Each sub-image is classified into the responsive class by an ANN. Finally, image has been compiled all the classify result of ANN. The experimental results show the feasibility of our proposal model.

Keyword: image classification; feature Extraction; artificial neural network.

## I. INTRODUCTION

A continuous video consists of two inter-related components: 1) tracks of the persons in the video and 2) localization and labels of the

activities of interest performed by these actors. Activity analysis of solving the recognition problems. The solution to one problem can help in finding the solution to the other. Similarly, information about the location and labels of activities in a scene can help in determining the movement of people in the scene. Therefore, we propose a method which performs the tasks in image classification and feature extraction. Image classification is one of classical problems of concern in image processing. The goal of image classification is to predict the categories of the input image using its features.

The ANN classifier, a conventional non-parametric, calculates the distance between the feature vector of the input image (unknown class image) and the feature vector of training image dataset. Artificial Neural Network (ANN), a brain-style computational model, has been used for many applications. Researchers have developed various ANN's structure in accordant with their

problem. After the network is trained, it can be used for image classification. Besides there are some integrated multi techniques model for classifying such as Multi Artificial Neural Network (MANN) applying for facial expression classification, and Multi Classifier Scheme applying for Adult image classification.

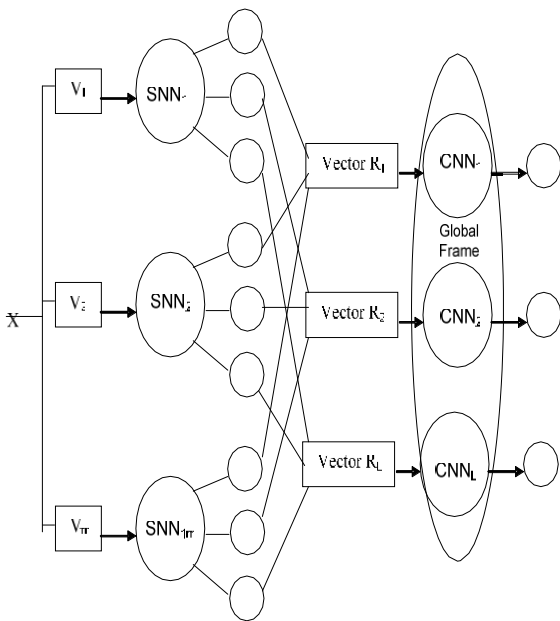


Fig. 1 Multi Artificial Neural Network model.

This model uses many Neural Networks so that the training phrase is complex and long. Besides multi classifier scheme has just been proposed for Adult image classification with low level feature. This experiment has showed that we need to choose the appropriate classifiers for the feature extraction to increase the precision of image classification. On the other hand, the

precision of classification system depends on the feature extraction and the classifier.

2. Background and Related Work  
2.1The stages of image classification

The main steps in the image classification process are shown in the following diagram:

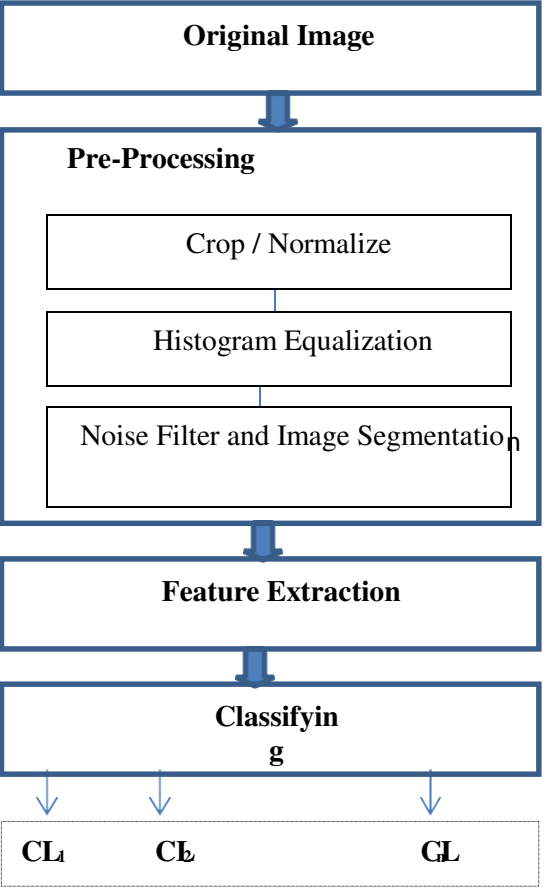


Fig. 2 Image Classification Process  
CL1, CL2, ..., CLn refers to the classes or categories that images are classified into. Step 1, pre-processing, is required before applying any image analysis methods. The images are normalized, performing histogram equalization,

applying the noise filter and segmenting. In the step 2, feature extraction, using the suitable transform to decompose an image the features of images are the input of our classification system. Finally, images are classified into the responsive classes by the suitable techniques.

## 2.2 Image Feature Extraction

The extraction of image features is the fundamental step for image classification. There are various types of features for image classification's aim as follow: color and shape features, statistical features of pixels, and transform coefficient features. In addition, some researchers have used algebraic feature for image recognition and image classification.

The output of image's feature extraction is often a vector or multi vectors. In this research, an image is extracted to k feature vectors based on k representing sub-space.

## 3. PROPOSED METHODOLOGY

Methodology is shown in flowchart. Step by step process is followed for pre-processing of image. MATLAB provides all image processing function and toolbox. MATLAB have large library functions and set of tools. Main features of MATLAB are following:

1. It provides advanced algorithm for high numerical computation.

2. Ability to define user define functions and large collection of mathematical functions.

3. For plotting and displaying data, two and three dimensional graphics are supported.

4. Online help is present which is very much helpful for new user.

5. Powerful, effective and efficient matrix and vector oriented high level programming language is provided by MATLAB.

6. Several toolboxes are provides for solving domain specific problems. Some of toolboxes are Image processing toolbox. Fuzzy logic, Digital signal processing toolbox, neural network toolbox etc.

## INPUT VIDEO

The input video format taken in this paper is AVI. AVI stands for audio video interleave. An AVI file actually stores audio video data under the format RIFF (Resource Interchange File Format). In AVI files, audio data and video data are stored next to each other, so that synchronous audio with-video playback can be allowed. Audio data is usually stored in AVI files in uncompressed pulse code modulation format with various parameters. Video data is usually stored in AVI files in compressed format with various parameters and codec.



Fig: 3 input image

**GABOR FILTER**

Its impulse response is defined by a sinusoidal wave multiplied by a Gaussian function. Gabor filter are directly related to Gabor wavelets, since they can be designed for a number of dilations and rotation. A set of Gabor filter with different frequencies and orientation may be helpful for extracting useful feature from an image. Gabor filter have been widely used in pattern analysis application.

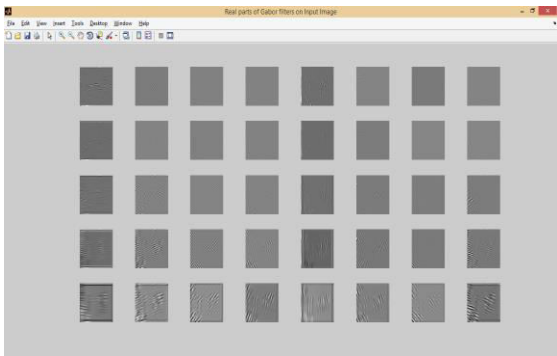


Fig: 4 real parts of Gabor filter on input image.

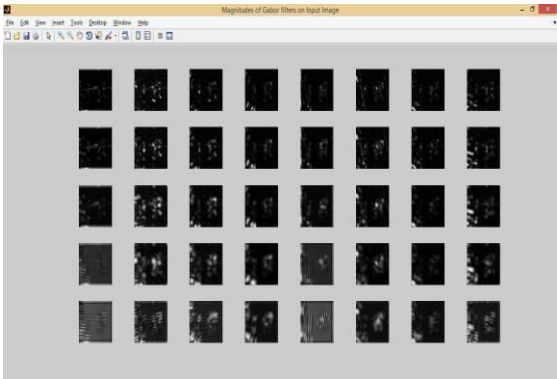


Fig: 5 magnitudes of Gabor filter on input image.

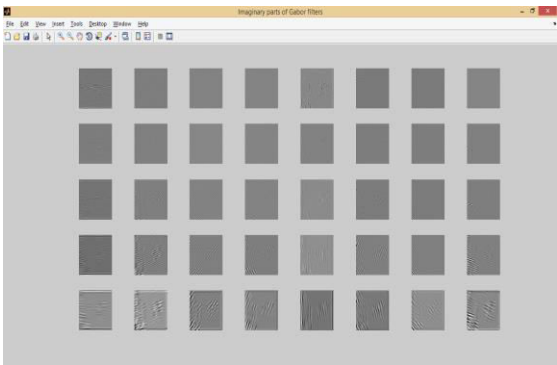


Fig: 6 imaginary parts of gabor filter on input image.

**FEATURE EXTRACTION**

In machine learning, pattern recognition and in image processing, feature extraction starts from an initial set of measured data and builds derived values (feature) intended to be informative and non-redundant, facilitating the subsequent learning and generalization steps, and in some cases leading to better human interpretations. Feature extraction is related to dimensionality reduction.

When the input data to an algorithm is too large to be processed and it is suspected to be redundant (e.g. the same

measurement in both feet and meters, or the repetitiveness of images presented as pixels), then it can be transformed into a reduced set of features (also named a features vector). This process is called *feature extraction*. The extracted features are expected to contain the relevant information from the input data, so that the desired task can be performed by using this reduced representation instead of the complete initial data.

### 3.3.1 CORNER DETECTION

Corner detection is an approach used with in computer vision system to extract certain kinds of feature and infer the content of an image. Corner detection is frequently used in motion detection, image registration, image tracking, image mosaicking, panorama stitching, 3D modeling, and object recognition.

Corner detection are not usually very robust and often require large redundancies introduced to prevent the efforts of individual error from domination recognition task.



Fig: 7 corner detection

### 3.4 ANN CLASSIFICATION

This feature vector is the input of ANN for image classification based on a sub-space. Every ANN has 3 layers: input, hidden and output. The number nodes of input layer are equal to the dimension of feature vector, called in. The number nodes of output are equal to n, the number of classes.

The simple integrating way is to calculate the mean value:

$$CL = \frac{1}{K} \sum_{i=1}^K CL\_SSi \quad (1)$$

Or weighted mean value:

$$CL = \frac{1}{K} \sum_{i=1}^K w_i CL\_SSi \quad (2)$$

Where  $w_i$  is the weight of classification result of subspace  $SS_i$ , and satisfies:

$$\sum_{i=1}^K w_i = 1 \quad (3)$$

Neural Network for identify the weights or importance of the local results. In this research, we suggest that the parameter of the hyper plans of ANN is instead of the weights  $w_i$ . Although ANN need to be trained first, the parameter of ANN is adjusted to suitable for the training data in the specific problem.



Fig: 8 ANN Classification.

## CONCLUSION:

In this research, we develop an integrated model of ANN for image classification. ANN is easy to design and deploy for the specific classification problem. The precision is high, but the performance of processing time need to improve, especially we apply for complex image classification such as facial image. Experiment on a variety of video with highly articulated object or complex background presented verified the effectiveness and robustness of our proposed method.

## REFERENCES:

- [1] K. K. Sung, "Learning and example selection for object and pattern recognition", Ph.D Thesis, MIT, Artificial Intelligence Laboratory and Center for Biological and Computational Learning, Cambridge, 1996.
- [2] D. Maltoni, D. Maio, A. K. Jain, and S. Prabhakar, Handbook of Fingerprint Recognition, 2nd ed. Berlin, Germany: Springer-Verlag, 2009.
- [3] FVC2006: The fourth international fingerprint verification competition. (2006). [Online]. Available: <http://bias.csr.unibo.it/fvc2006/>
- [4] V. N. Dvornychenko, and M. D. Garriss, "Summary of NIST latent fingerprint testing workshop," Nat. Inst. Standards Technol., Gaithersburg, MD, USA, Tech. Rep. NISTIR 7377, Nov. 2006.
- [5] Ming-Hsuan Yang, David J. Kriegman and Narendra Ahuja, "Detecting Faces in Images: A Survey", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 24, No. 1, pp. 34-58, 2002.
- [6] L. M. Wein and M. Baveja, "Using fingerprint image quality to improve the identification performance of the U.S. visitor and immigrant status indicator technology program," Proc. Nat. Acad. Sci. USA, vol. 102, no. 21, pp. 7772-7775, 2005.
- [7] S. Yoon, J. Feng, and A. K. Jain, "Face Detection Using Feed Forward Neural Network in Matlab," IEEE Trans. Pattern Anal. Mach. Intell., vol. 34, no. 3, pp. 451-464, Mar. 2012.
- [8] E. Tabassi, C. Wilson, and C. Watson, "Fingerprint image quality," Nat. Inst. Standards Technol., Gaithersburg, MD, USA, Tech. Rep. NISTIR 7151, Aug. 2004.
- [9] F. Alonso-Fernandez, J. Fiérrez-Aguilar, J. Ortega-Garcia, J. Gonzalez-Rodriguez, H. Fronthaler, K. Kollreider, and J. Bigün, "A comparative study of fingerprint image-quality estimation methods," IEEE Trans. Inf. Forensics Security, vol. 2, no. 4, pp. 734-743, Dec. 2007.
- [10] J. Fiérrez-Aguilar, Y. Chen, J. Ortega-Garcia, and A. K. Jain, "Incorporating image quality in multi-algorithm fingerprint verification," in Proc. Int. Conf. Biometrics, 2006, pp. 213-220.
- [11] L. Hong, Y. Wan, and A. K. Jain, "Face Detection using Neural Network & Gabor Wavelet Transform," IEEE Trans. Pattern Anal. Mach. Intell., vol. 20, no. 8, pp. 777-789, Aug. 1998.
- [12] S. Chikkerur, A. N. Cartwright, and V. Govindaraju, "Fingerprint enhancement using STFT analysis," Pattern Recognit., vol. 40, no. 1, pp. 198-211, 2007.
- [13] F. Turrone, R. Cappelli, and D. Maltoni, "Fingerprint enhancement using contextual iterative filtering," in Proc. Int. Conf. Biometrics, 2012, pp. 152-157.
- [14] Yang M. H., Ahuja N., and Kriegman D., "A survey on face detection methods", IEEE Transactions on Pattern Analysis and Machine Intelligence, to appear 2001.
- [15] J. Zhu, M. Vai and P. Mak, "Gabor wavelets Transform and extended nearest feature space classifier for face recognition", Proceedings of the Third International Conference on Image and Graphics, pp. 246-249, 2004.
- [16] Lawrence S. , Giles C., Tsoi A. and Back A., "Face Recognition: A Convolutional Neural Network Approach", IEEE Trans. on Neural Networks, vol. 8, pp. 98-113, 1997.
- [17] N. Ratha and R. Bolle, "Effect of controlled image acquisition on fingerprint matching," in Int. Conf. Pattern Recognit., 1998, vol. 2, pp. 1659-1661.

- [18] Y. Fujii, "Implementation of Artificial Neural Network for Face Recognition using," U.S. Patent No. 7 660 447, Feb. 9, 2010.
- [19] C. Dorai, N. K. Ratha, and R. M. Bolle, "Dynamic behavior analysis in compressed fingerprint videos," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 14, no. 1, pp. 58–73, Jan. 2004.
- [20] N. K. Ratha, K. Karu, S. Chen, and A. K. Jain, "A real-time matching system for large fingerprint databases," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 18, no. 8, pp. 799–813, Aug. 1996.
- [21] X. Chen, J. Tian, and X. Yang, "A new algorithm for distorted fingerprints matching based on normalized fuzzy similarity measure," *IEEE Trans. Image Process.*, vol. 15, no. 3, pp. 767–776, Mar. 2006.
- [22] A. M. Bazen, and S. H. Gerez, "Fingerprint matching by thin-plate spline modelling of elastic deformations," *Pattern Recognit.*, vol. 36, no. 8, pp. 1859–1867, Aug. 2003.
- [23] L. R. Thebaud, "Systems and methods with identity verification by comparison and interpretation of skin patterns such as fingerprints," U.S. Patent No. 5 909 501, Jun. 1, 1999.
- [24] Z. M. Kovacs-Vajna, "A fingerprint verification system based on triangular matching and dynamic time warping," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 11, pp. 1266–1276, Nov. 2000.
- [25] J. Feng, Z. Ouyang, and A. Cai, "Fingerprint matching using ridges," *Pattern Recognit.*, vol. 39, no. 11, pp. 2131–2140, 2006.
- [26] A. Ross, S. C. Dass, and A. K. Jain, "A deformable model for fingerprint matching," *Pattern Recognit.*, vol. 38, no. 1, pp. 95–103, 2005.
- [27] A. Ross, S. C. Dass, and A. K. Jain, "Fingerprint warping using ridge curve correspondences," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 1, pp. 19–30, Jan. 2006.
- [28] A. Senior, and R. Bolle, "Improved fingerprint matching by distortion removal," *IEICE Trans. Inf. Syst.*, vol. 84, no. 7, pp. 825–831, Jul. 2001.
- [29] D. Wan, and J. Zhou, "Fingerprint recognition using model-based density map," *IEEE Trans. Image Process.*, vol. 15, no. 6, pp. 1690–1696, Jun. 2006.
- [30] J. Feng, "Combining minutiae descriptors for fingerprint matching," *Pattern Recognit.*, vol. 41, no. 1, pp. 342–352, 2008.
- [31] A. M. Bazen, and S. H. Gerez, "Systematic methods for the computation of the directional fields and singular points of fingerprints," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 7, pp. 905–919, Jul. 2002.
- [32] C.-C. Chang, and C.-J. Lin, "LIBSVM: A Library for support vector machines," *ACM Trans. Intell. Syst. Technol.*, vol. 2, pp. 27:1–27:27, 2011.
- [33] F. L. Bookstein, "Principal warps: Thin plate splines and the decomposition of deformations," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 11, no. 6, pp. 567–585, Jun. 1989.
- [34] S. Novikov and O. Ushmaev, "Principal deformations of fingerprints," in *Audio and Video Based Biometric Person Authentication*. Berlin, Germany: Springer-Verlag, 2005, pp. 250–259.
- [35] D. Rueckert, A. F. Frangi, and J. A. Schnabel, "Automatic construction of 3-D statistical deformation models of the brain using nonrigid registration," *IEEE Trans. Med. Imag.*, vol. 22, no. 8, pp. 1014–1025, Aug. 2003.
- [36] S. Tang, Y. Fan, G. Wu, M. Kim, and D. Shen, "RABBIT: Rapid alignment of brains by building intermediate templates," *Neuroimage*, vol. 47, no. 4, pp. 1277–1287, 2009.
- [37] D. H. Ballard, "Generalizing the Hough transform to detect arbitrary shapes," *Pattern Recognit.*, vol. 13, no. 2, pp. 111–122, 1981.
- [38] J. Dai, J. Feng, and J. Zhou, "Robust and efficient ridge-based palmprint matching," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 8, pp. 1618–1632, Aug. 2012.



# Digital image encryption based on fuzzy logic

<sup>1</sup>S R Mahipal

*Computer Science and Engineering*  
*Malla Reddy college of Engineering*  
Hyderabad, India  
email:srmahipal@gmail.com

<sup>2</sup>V V Ramanjaneyulu

*Computer Science and Engineering*  
*Malla Reddy college of Engineering*  
Hyderabad, India  
email:ramu5b4@gmail.com

<sup>3</sup>Dr. Sunil Tekale

*Computer Science and Engineering*  
*Malla Reddy college of Engineering*  
Hyderabad, India  
email:sunil.tekale2010@gmail.com

**Abstract**—Today it is a big challenge to see that the data that is transferred from one place to another over network is highly secured, we need to have a secured network where there is no chance for the hackers to play their role, as many organizations are dependent on technology and it is the major source of the growth, we need to provide a highly secured and reliable data transfer. In this paper we have designed an algorithm which helps in transmission of data on network in a secured manner. The data in this is encrypted using matrix manipulation concept and the data after encryption is transferred on the network. Encryption is process of hiding the information, when the information is transferred through a network and decryption is the process of extracting the information from an encrypted information. For this encryption and decryption, we need some encryption and decryption algorithm which are proposed as a part of this paper.

**Keywords**—Encryption, Decryption, Image, Fuzzy Logic.

## I. INTRODUCTION

The world has recently witnessed major development in the information and communications technology and the digital world. The computer science is used in all areas of life, including sending and receiving digital images as the importance of which are tremendously increasing. The images are sent and treated automatically, and this requires careful secret storage of data to be sent as there are many reasons to protect the image of the breach. Cryptography is the science which deals with ways that help us to protect and store information and transfer in a wide range and these methods depend on a secret key that is used to encrypt data. (1).

Security is the main problem in the modern digital world. There are a lot of cyber-crimes have arisen with the development of technology. [3] As solutions for these security risks users can shut down unused services, keep patches updated, reduce permissions and access rights of applications and users.

Another solution for this problem can be provided by using cryptography. [4]Cryptography consists of cryptology and crypto analysis. Encryption comes under cryptology. It is the process of converting a readable message into an unreadable format. [5] A set of rules is using for that process. It is called an encryption algorithm. Most of the nowadays existing encryption algorithms only concern on security. [6] However,

users who have connections with low bandwidths need an encryption algorithm, which uses a low processing power. High security algorithms tend to take little more processing power than the low security algorithms. Nevertheless, newly implemented encryption algorithm, which has the facility to control both desired security level and the processing level, will be a great improvement for current real world applications. Various algorithms have been proposed to implement encryption in digital images. They can be categorized into three major clusters (i) value transformation [2], (ii) pixel position permutation [7, 8] (iii) chaotic systems [1517].

Fuzzy logic is a problem-solving control system methodology that presents itself to implementation in systems ranging from simple, small, embedded micro-controllers to large, networked, multi-channel PC or workstation-based data acquisition and control systems. [18][19] Fuzzy logic provides a simple way to arrive at a definite conclusion based upon vague, ambiguous, imprecise, noisy, or missing input information. [20] In fuzzy logic rules and membership sets are used to make a decision. [21] To achieve security and low processing, the algorithm uses variable keys. 0<sup>th</sup> position gives a fully low processing algorithm, and 1<sup>st</sup> position gives fully secured algorithm. The fuzzification changes depending on the key size and the number of mapping tables of the encryption algorithm. Users can input the desired key. One character will be 8-bit long. The main algorithm structure defines different key sizes up to 128bit. User can enter desired key- (application defines as the password) and also depending on the number of mapping tables' algorithm will allocate weight dynamically. Allocation of the weights will differ from 0.0 to 1.0 range; and the number of security levels will be vary from 1-16. The number of rounds will be determined by pre-defined mapping tables and the users initial input. Mapping tables are predefined in the algorithm and consists of mathematically defined values, and then those values will dynamically choose the relevant algorithm procedure once the user input the key to encryption.

### A. Fuzzification

Fuzzification is the operation of making a crisp quantity Fuzzy. It is simply done by recognizing that many of the quantities that are regarded as crisp deterministic are actually not deterministic at all; they carry considerable doubt. If the

form of doubt happens to arise, because of imprecision, opacity, or ambiguity, then the variable is probably Fuzzy and can be represented by a membership function.

### B. Defuzzification

For a given input, several IF/THEN rules could be begun at the same time. Each rule will have a different strength because a given input may belong to more than one Fuzzy set, but with different membership values.

## II. THE AIMS OF RESEARCH

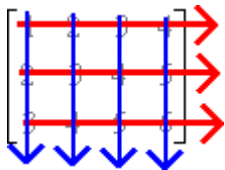
The main aim of the research is to build encryption system based on fuzzy logic to secure confidential trading images. Here, the principles of information technology are applied to encrypt images and decrypt them. Besides, it allows the sender to make sure that the images will reach just the people to whom the images are sent to, and the right way that no one can decode it except the receiver person.

## III. PROPOSED ALGORITHM

In this algorithm the data to be encrypted is considered at the lowest level and the data then is taken into account byte by byte, which is then transformed into matrix form and the same is considered after interchange of rows and columns. Then the resultant data is made to undergo swapping of 4 bits on both sides which results in final version of the data. The same is decrypted on the other end by using the decryption keys, which helps the user to get the data in original form. The format of this will be

0	1	1	0	0	1	1	1
---	---	---	---	---	---	---	---

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \text{ After interchanging rows and cols}$$



0	0	1	1	1	1	0	1
---	---	---	---	---	---	---	---

Of each matrix data will again undergo transformation nibbles in reverse order so as to get

1	1	0	0	1	0	1	1
---	---	---	---	---	---	---	---

## IV. FEATURES

The proposed algorithm has been proved to provide high protection to the images data from illegal intrusions. It is fast in

the process of encryption and decryption. The decryption process does not induce any loss of image data, and it can deal with different format of images as will.

## V. EXPERIMENTAL DETAILS AND RESULT

The proposed encryption algorithm can be classified into multiple criteria such as lossless, maximum distortion, maximum performance and maximum speed. In this section, the proposed algorithm is applied on different sizes and types of images.

The test images employed show a positive result. The encryption and decryption algorithm are implemented in JAVA and core2duo of 2.66 GHz machine. The decryption algorithm takes between 76 and 100 Milliseconds to get executed. Calculating the lossless by this formula.

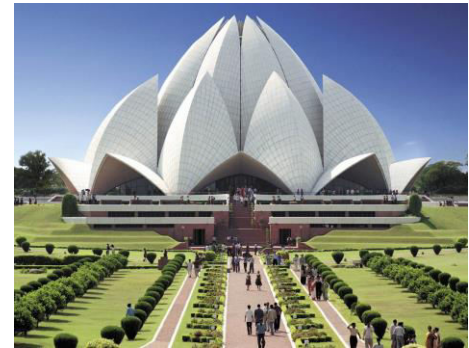


Fig. 1. Original Image

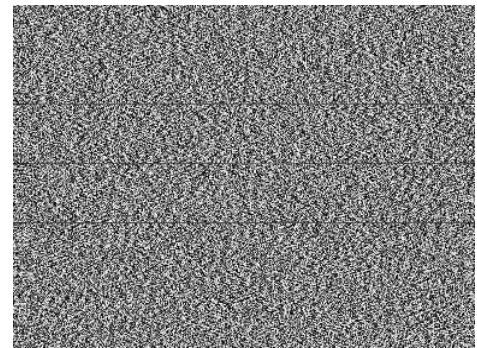


Fig. 2. Encrypted Image

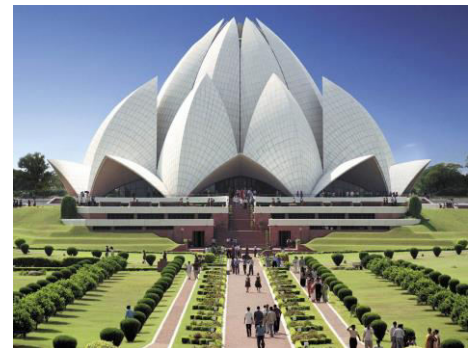


Fig. 3. Decrypted Image

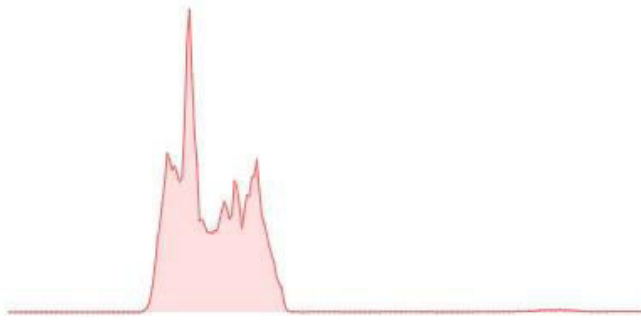


Fig. 4. Histogram of Original Image

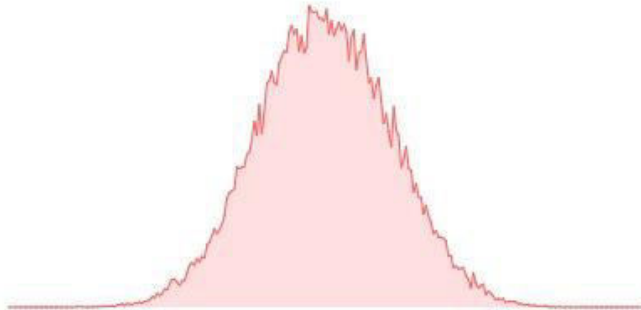


Fig. 5. Histogram of Encrypted Image



Fig. 6. Histogram of Decrypted Image

#### A. Security analysis

Security analysis The strength is the most essential feature that a good quality encryption algorithm should possess. If the encryption algorithm is unable to prevent all types of attack including statistical and brute force attacks, it will not be sufficient for protecting the data. Many experiments are done for defining the competency of the proposed technique. In this part, the proposed technique is applied on images which have different formats and sizes.

#### B. Statistical Analysis

The encrypted images should hold certain random properties to prevent statistical attacks. A statistical analysis has been done by calculating the histograms, the entropy, the correlations and differential analysis for the plain image and the encrypted image for proving the strength of the proposed algorithm. After various images are tested, it appears that the intensity values are good.

#### C. Histogram Analysis

An image histogram is a commonly used method of analysis in image processing and data mining applications. One of the various benefits of the histogram is that it shows the shape of the distribution for a large set of data. Therefore, an image histogram provides a clear illustration of how the pixels in an image are distributed by graphing the number of pixels at each color intensity level. It is essential to make sure that the encrypted and original images possess different statistics. The histogram analysis shows the ways that pixels in an image are distributed by plotting the number of pixels at each intensity level. The Fig. 3. shows the results of the experiment on the plain image, its corresponding cipher image and their histograms. The histogram of each plain image explains how the pixels are distributed by graphing the number of pixels at every grey level [26]. The results show that the histogram of the encrypted image is uniformly distributed and significantly different from the respective histograms of the original images.

TABLE I. IMAGE PROPERTIES

S. No	Image name	Dimensions	Size before encryption	Size after encryption
1	Lotus temple	256*256	22.9KB	21.5KB
2	Hi tech city	259*194	21.6KB	19.17KB
3	Parliament of India	350*350	32KB	29.4KB

#### CONCLUSION

We have designed an algorithm which helps in transmission of data on network in a secured manner. The data in this is encrypted using matrix manipulation concept and the data after encryption is transferred on the network. Encryption is process of hiding the information, when the information is transferred through a network and decryption is the process of extracting the information from an encrypted information, here the interchange of data is done by taking the byte of data in matrix form and then the same data is reversed by using the 2 nibbles.

#### REFERENCES

- [1] Gamil R.S. Qaid , Sanjay N. Talbar, "Encryption and Decryption of Digital Image Using Color Signal" IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 2, No 2, PP. 588 -592, March 2012.
- [2] Komal D. Patel, Sonal Belani "Image Encryption Using Different Techniques: A Review", International Journal of Emerging Technology and Advanced Engineering, Volume 1, Issue 1, PP. 30 -34, 2011.
- [3] Borko Furht, Edin Muharemagic, Daniel Socek "Multimedia Encryption and Watermarking", Springer, USA, 2005.

- [4] Aloha Sinha, Kehar Singh, "A technique for image encryption using digital signature", Optics Communications, Vol-218 (2203), PP.229-234.
- [5] S.S.Maniccam, N.G. Bourbakis, "Lossless image compression and encryption using SCAN", Pattern Recognition, 34 (2001), PP.1229- 1245.
- [6] Ahmed Bashir Abugharsa, Abd Samad Bin Hasan Basari, Hamida Almangush, "A Novel Image Encryption using an Integration Technique of Blocks Rotation based on the Magic cube and the AES Algorithm", IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 4, No 1, PP.41-47 July 2012.
- [7] Liu, Z., et al., "Image encryption scheme by using iterative random phase encoding in gyrator transform domains", Optics and Lasers in Engineering, 2011. 49(4): PP. 542-546.
- [8] Guo, Q., Z. Liu, and S. Liu, "Color image encryption by using Arnold and discrete fractional random transforms in HIS space", Optics and Lasers in Engineering, 2010. 48(12): PP. 1174-1181.
- [9] Liu, Z., et al., "Image encryption by using gyrator transform and Arnold transform", Journal of Electronic Imaging, 2011. 20: PP. 013020.
- [10] R. Tao, X. Y. Meng, and Y.Wang, "Image encryption with multi orders of fractional Fourier transforms In Information Forensics and Security", 2010: IEEE Transactions on Image Processing.
- [11] Zunino, R., "Fractal circuit layout for spatial décor relation of images", Electronics Letters, 1998. 34(20): PP. 1929-1930.
- [12] Zhang, G. and Q. Liu, "A novel image encryption method based on total shuffling scheme". Optics Communications, 2011.
- [13] Zhao, X. and C. Gang, "Ergodic matrix in image encryption", 2002.
- [14] Zhu, Z., et al., "A chaos-based symmetric image encryption scheme using a bit-level permutation" Information Sciences, 2011. 181(6): PP. 1171-1186.
- [15] Huang, C. and H. Nien, "Multi chaotic systems based pixel shuffle for image encryption", Optics Communications, 2009. 282(11): PP. 2123-2127.
- [16] Wang, K., et al., "On the security of 3D Cat map based symmetric image encryption scheme. Physics Letters A, 2005. 343(6): PP. 432-439.
- [17] Wang, X.Y., et al., "A chaotic image encryption algorithm based on perceptron model. Nonlinear Dynamics, 2010. 62(3): PP. 615-621.
- [18] Fuzzy Logic: An Introduction [online] <http://www.seattlerobotics.org>
- [19] "Europe Gets into Fuzzy Logic", Electronics Engineering Times, 1991
- [20] "Fuzzy Sets and Applications: Selected Papers by L.A. Zadeh", ed. R.R. Yager et al. (John Wiley, New York, 1987).
- [21] "U.S. Loses Focus on Fuzzy Logic" (Machine Design, June 21, 1990).
- [22] Wang, Y., et al., "A new chaos-based fast image encryption algorithm. Applied Soft Computing, 2011. 11(1): p. 514-522.
- [23] Monisha Sharma, Shri Shankarcharya, Manoj Kumar Kowar, "Image Encryption Techniques Using Chaotic Schemes: A Review" International Journal of Engineering Science and Technology. Vol. 2(6), PP. 2359-2363. 2010.
- [24] Nawal El-Fishawy, Osama M. Abu Zaid, "Quality of Encryption Measurement of Bitmap Images with RC6, MRC6, and Rijndael Block Cipher Algorithms" International Journal of network Security, Vol.5, No.3, PP.241-251, Nov.2007.
- [25] Sara Tedmori, Nijad Al-Najdawi "Lossless Image Cryptography Algorithm Based on Discrete Cosine Transform" IAJIT First Online Publication vol.3, 2011.
- [26] Nassiba Wafa Abderrahim, Fatima Zohra Benmansour, Omar Seddiki "Integration of chaotic sequences uniformly distributed in a new image encryption algorithm" IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 2, No 3, March 2012.

## RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF GRAPHS

M. S. MALCHIJAH RAJ  
*Research Scholar,  
 Research and Development Centre,  
 Bharathiar University,  
 Coimbatore - 641046, India.  
 malchijahraj@gmail.com*

**ABSTRACT.** Let  $G = (V, E)$  be a connected graph of order  $p$ . A set of vertices  $S$  in a graph  $G$  is a *restrained geodetic set* if  $S$  is a geodetic set and the subgraph  $G[V - S]$  induced by  $V - S$  has no isolated vertices. The minimum cardinality of a restrained geodetic set, denoted by  $g_r(G)$ , is called the *restrained geodetic number* of  $G$ . A  $g_r(G)$  - set is a restrained geodetic set of cardinality  $g_r(G)$ . Correspondingly, a set  $W$  of vertices of a graph  $G$  is a *restrained Steiner set* if  $W$  is a Steiner set, and if either  $W = V$  or the subgraph  $G[V - W]$  induced by  $V - W$  has no isolated vertices. The minimum cardinality of a restrained Steiner set of  $G$  is the *restrained Steiner number* of  $G$ , and is denoted by  $s_r(G)$ . In this paper we study the restrained geodetic number and the restrained Steiner number of some standard graphs.

**Key Words:** geodetic set, restrained geodetic number, Steiner number, restrained Steiner number.

**AMS Subject Classification:** 05C12

### 1. INTRODUCTION

By a graph  $G = (V, E)$ , we mean a finite undirected graph without loops or multiple edges. The order and size of  $G$  are denoted by  $p$  and  $q$  respectively. The distance  $d(u, v)$  between two vertices  $u$  and  $v$  in a connected graph  $G$  is the length of a shortest  $u - v$  path in  $G$ . An  $u - v$  path of length  $d(u, v)$  is called an  $u - v$  geodesic. It is

known that this distance is a metric on the vertex set  $V(G)$ . For a vertex  $v$  of  $G$ , the *eccentricity*  $e(v)$  is the distance between  $v$  and a vertex farthest from  $v$ . The minimum eccentricity among the vertices of  $G$  is the radius,  $radG$  and the maximum eccentricity is its diameter,  $diamG$  of  $G$ . Two vertices  $x$  and  $y$  are *antipodal* if  $d(x, y) = diamG$ . A vertex  $v$  is an *extreme vertex* or a *simplicial vertex* of a graph  $G$  if the

## 2RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF GRAPHS

subgraph induced by its neighbors is complete. Note that every end vertex is simplicial. If  $e = \{u, v\}$  is an edge of a graph  $G$  with  $d(u) = 1$  and  $d(v) > 1$ , then we call  $e$  a pendant edge,  $u$  a leaf or end vertex and  $v$  a support.

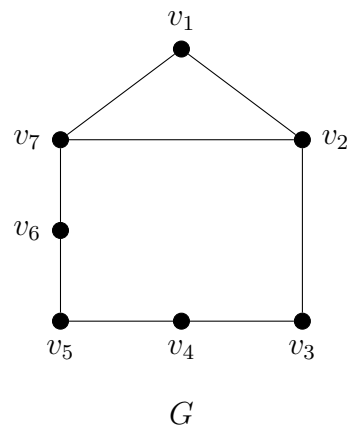
A *geodesic set* of  $G$  is a set  $S \subseteq V(G)$  such that every vertex of  $G$  is contained in a geodesic joining some pair of vertices in  $S$ . The *geodesic number*  $g(G)$  of  $G$  is the minimum cardinality of its geodesic sets and any geodesic set of cardinality  $g(G)$  is a *minimum geodesic set*. A set of vertices  $S$  in a graph  $G$  is a *restrained geodesic set* if  $S$  is a geodesic set and the subgraph  $G[V - S]$  induced by  $V - S$  has no isolated vertices. The minimum cardinality of a restrained geodesic set, denoted by  $g_r(G)$ , is called the *restrained geodesic number* of  $G$ . A  $g_r(G)$  - set is a restrained geodesic set of cardinality  $g_r(G)$ .

For a nonempty set  $W$  of vertices in a connected graph  $G$ , the *Steiner distance*  $d(W)$  of  $W$  is the minimum size of a connected subgraph of  $G$  containing  $W$ . Necessarily, each subgraph is a tree and is called a *Steiner tree* with respect to  $W$  or a *Steiner  $W$ -tree*.  $S(W)$  denotes the set of all vertices that lie on Steiner  $W$ -trees. A set  $W \subseteq V(G)$  is called a *Steiner set* of  $G$  if every vertex of  $G$  lies on some Steiner  $W$ -tree or if  $S(W) = V(G)$ . A Steiner set of minimum cardinality

is a *minimum Steiner set* or simply a  *$s$ -set* and this cardinality is the *Steiner number*  $s(G)$  of  $G$ . A set  $W$  of vertices of a graph  $G$  is a *restrained Steiner set* if  $W$  is a Steiner set, and if either  $W = V$  or the subgraph  $G[V - W]$  induced by  $V - W$  has no isolated vertices. The minimum cardinality of a restrained Steiner set of  $G$  is the *restrained Steiner number* of  $G$ , and is denoted by  $s_r(G)$ .

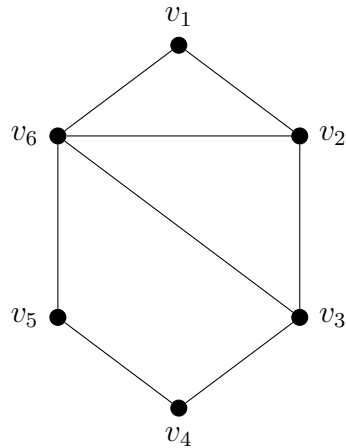
We have  $2 \leq g(G) \leq s(G) \leq p$  and also that, Every Steiner set in a connected graph is a geodesic set.

As an example, consider the graph  $G$  given in Figure 1. The set  $S = \{v_1, v_4, v_6\}$  is a  $g$ -set of  $G$ . Hence  $g(G) = 3$ . Since the Steiner minimum tree for  $S$  has Steiner distance 4 without covering the vertices  $v_2$  and  $v_3$ ,  $S$  is not a Steiner set for  $G$ . The sets  $W_1 = \{v_1, v_2, v_5\}$ ,  $W_2 = \{v_1, v_3, v_6\}$ ,  $W_3 = \{v_1, v_4, v_5\}$  and  $W_4 = \{v_1, v_4, v_7\}$  are all  $s$ -sets of  $G$  so that  $s(G) = 3$ . The sets  $W_i, i = 1, 2, 3, 4$  are also geodesic sets or geodesic covers of  $G$ .



# RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF GRAPHS3

Figure 1



G

Figure 2

Consider the graph  $G$  in Figure 2. The set  $S = \{v_1, v_4\}$  is a  $g$  - set as well as a  $s$  - set of  $G$  so that  $g(G) = s(G) = 2$ . Also the set  $S_1 = \{v_1, v_3, v_5\}$  is also a geodetic as well as a Steiner set of  $G$ . Since the subgraph  $G[V - S_1]$  induced by  $V - S_1$  has an isolated vertex  $v_4$ ,  $S_1$  is not a restrained geodetic or restrained Steiner set of  $G$ . Therefore the set  $S_2 = \{v_1, v_3, v_4, v_5\}$  is a  $g_r(G)$  - set as well as  $s_r(G)$  - set so that  $g_r(G) = s_r(G) = 4$ . We have  $2 \leq g(G) \leq g_r(G) \leq p$ . Similarly  $2 \leq s(G) \leq s_r(G) \leq p$  holds for every connected graph  $G$ .

We make use of the following results in this paper for characterizing different classes of graphs with respect to their restrained geodetic number or restrained Steiner number. Since each restrained geodetic set is a geodetic set, and since the

complement of each restrained geodetic set has cardinality different from 1, we have  $g_r(G) \neq p - 1$ . The similar argument is true for the restrained Steiner number also. So we have that there is no graph of order  $p$  with  $s_r(G) = p - 1$ .

**Theorem 1.1.** *If  $G$  is a connected nontrivial graph, then every simplicial vertex belongs to every geodetic set of  $G$  as well as every restrained geodetic set of  $G$ .*

**Theorem 1.2.** *Each extreme vertex of a graph  $G$  belongs to every Steiner set of  $G$  as well as every restrained Steiner set of  $G$ .*

**Theorem 1.3.** *For the complete bipartite graph  $G = K_{m,n}$ ,*

$$s(G) = \begin{cases} 2 & \text{if } m = n = 1 \\ n & \text{if } n \geq 2, m = 1 \\ \min\{m, n\} & \text{if } m, n \geq 2 \end{cases}$$

## 2. RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF VARIOUS GRAPHS

In this section the restrained geodetic and the restrained Steiner number of various graphs are described using results from literature.

**Class 2.1.** *Complete graphs  $K_p$*

Since each vertex in a complete graph is an extreme vertex, the set of all vertices of the complete graph  $K_p$  is the unique restrained geodetic set as well as the unique restrained Steiner set of  $K_p$ . So we have  $g_r(G) = s_r(G) = p$ .

#### 4RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF GRAPHS

##### **Class 2.2.** *Cycles $C_p$*

Let  $G = C_p$ . It is clear that  $g_r(C_p) = s_r(C_p) = p$  for  $p \in \{3, 4, 5\}$ . Let  $p \geq 6$ . Let  $p$  be even. Then any pair of antipodal vertices of  $C_p$  forms the restrained geodetic set and the restrained Steiner set of  $C_p$ , and so  $g_r(C_{2p}) = s_r(C_{2p}) = 2$ . Let  $p$  be odd. It is clearly verified that no two element subset of  $C_p$  is a geodetic set or a Steiner set of  $G$  and so  $g_r(G) = s_r(G) \geq 3$ . For any vertex  $u$ , let  $v$  and  $w$  be the two antipodal vertices. Let  $W_1 = \{u, v, w\}$ . Then  $W_1$  is a restrained geodetic set and a restrained Steiner set so that  $g_r(C_{2p+1}) = s_r(C_{2p+1}) = 3$ .

##### **Class 2.3.** *Trees $T$*

Let  $W$  denote the set of all end vertices of  $T$ . If  $T$  is not a Star, then  $W$  is the unique restrained geodetic set as well as the restrained Steiner set of  $T$  and so  $g_r(T) = s_r(T) = |W|$ . If  $T$  is a star  $K_{1,p-1}$  then  $g_r(T) = s_r(T) = p$ .

##### **Class 2.4.** *Complete bipartite graphs $K_{m,n}$*

Let  $G = K_{m,n}$  ( $2 \leq m \leq n$ ). Let  $X$  and  $Y$  be the bipartite sets of  $G$ , where  $X = \{x_1, x_2, \dots, x_m\}$  and  $Y = \{y_1, y_2, \dots, y_n\}$ . Let  $2 \leq m \leq n$ . If  $m = 2$ , then obviously  $g_r(K_{2,n}) = 2 + n$ . If  $m \geq 3$ , then it is clear that  $\{x_1, x_2, y_1, y_2\}$  is a restrained geodetic set of  $G$  and so  $g_r(K_{m,n}) = 4$ . By Theorem 1.3,  $W = X$  is a  $s$ -set of

$G$ . Since the subgraph  $G[V - W]$  has isolated vertices,  $W$  is not a restrained Steiner set of  $G$  and so  $s_r(G) \geq m + 1$ . Let  $W'$  be a restrained Steiner set of  $G$  with  $|W'| \geq m + 1$ . Then either  $W' \subseteq Y$  or  $W' \subsetneq X \cup Y$ . If  $W' \subseteq Y$ , then  $W' = Y$  is the only Steiner set of  $G$ . Since the subgraph  $G[V - W']$  has isolated vertices,  $W'$  is not a restrained Steiner set of  $G$ , which is a contradiction. If  $W' \subsetneq X \cup Y$ , then  $\langle W' \rangle$  is connected. Then the Steiner  $W'$ -tree contains elements of  $W'$  only. Therefore  $W'$  is not a restrained Steiner set of  $G$ . Hence  $W' = X \cup Y$ . This implies  $s_r(G) = m + n$ .

##### **Class 2.5.** *Wheel $W_p = K_1 + C_{p-1}$ , ( $p \geq 5$ )*

Let  $G = W_p$ . Let  $v$  be the vertex of  $K_1$  and let  $v_1, v_2, \dots, v_{p-1}, v_1$  be the cycle  $C_{p-1}$ . Then  $S = \{v_i / i \text{ is odd}\}$  is a restrained geodetic set so that  $g_r(W_p) = \lfloor \frac{p}{2} \rfloor$ . Now for the restrained Steiner number, we have the following result. For  $p = 5$ , let  $W = \{v_1, v_3\}$ . Then  $W$  is a Steiner set of  $G$ . Since the subgraph  $G[V - W]$  has no isolated vertices,  $W$  is a restrained Steiner set of  $G$  so that  $s_r(W_p) = 2 = p - 3$ . Let  $p \geq 6$ . Let  $W$  be a Steiner set of  $G$ . If  $v \in W$ , then  $\langle W \rangle$  is connected. Then the Steiner  $W$ -tree contains elements of  $W$  only. Therefore  $v \notin W$ . Hence  $W \subseteq V(C_{p-1})$ . Let  $W = \{v_1, v_3, v_4, \dots, v_{p-2}\}$ . Then  $W$  is a restrained Steiner set of  $G$  and so  $s_r(G) \leq p - 3$ . We prove that

# RESTRAINED GEODETIC AND RESTRAINED STEINER NUMBER OF GRAPHS

$s_r(G) = p - 3$ . If not let  $W'$  be a restrained Steiner set of  $G$  with  $|W'| \leq p - 4$ . Then  $v \notin W'$ . Therefore there exists at least one  $v_i (1 \leq i \leq p - 1)$  such that  $v_i \notin W'$ . Then  $v_i (1 \leq i \leq p - 1)$  does not lie on any Steiner  $W'$  - tree of  $G$  and so  $W'$  is not a restrained Steiner set of  $G$ , which is a contradiction. Hence  $s_r(G) = p - 3$ .

**Class 2.6.** *Hyper cube*  $Q_n (n \geq 3)$

$Q_n$  has  $2^n$  vertices, which may be labeled  $(a_1 a_2 a_3 \dots a_n)$ , where each  $a_i (1 \leq i \leq n)$  is either 0 or 1. It is easily seen that  $\{(0, 0, 0, \dots, 0), (1, 1, 1, \dots, 1)\}$  is a  $g_r(Q_n)$  set as well as  $s_r(Q_n)$  set for  $n \geq 3$ . Hence  $g_r(Q_n) = s_r(Q_n) = 2$ .

## 3. CONCLUSION

The study of relationship between the restrained concept of geodetic and Steiner number of a graph can be done for various other concepts such as upper restrained edge geodetic and upper restrained edge Steiner, restrained edge geodetic and restrained edge Steiner, restrained forcing etc.

## REFERENCES

- [1] H. Abdollahzadeh Ahangar, V. Samodivkin, S. M. Sheikholeslami

- and Abdollah Khodkar, The restrained geodetic number of a graph, *The Bulletin of the Malaysian Mathematical Society Series*, **2** (2015), 1143 – 1155.
- [2] F. Buckley and F. Harary, Distance in Graphs, *Addition-Wesley, Redwood City, CA*, 1990.
- [3] G. Chartrand, F. Harary and P. Zhang, On the geodetic number of a graph, *Networks* **39** (2002), 1 – 6.
- [4] G. Chartrand, F. Harary and P. Zhang, The Steiner Number of a Graph, *Discrete Math.* **242** (2002), 41 – 54.
- [5] D. P. Day, O. R. Oellermann, H. C. Swart, Steiner distance-hereditary graphs. *SIAMJ. Discrete Math.* **7** (1994), 437 – 442.
- [6] F. Harary, Graph Theory, *Addition-Wesley*, 1969.
- [7] F. Harary, E. Loukakis and C. Tsouros, The geodetic number of a graph, *Math. Comput. Modelling* **17** (1993), 89 – 95.
- [8] J. John, M. S. Malchijah Raj, The Restrained Steiner Number of a Graph, (Communicated).
- [9] O. R. Oellermann, Tian, Steiner centers in graphs, *Graph Theory* **14** (5), (1990), 585 – 597.
- [10] I. M. Pelayo, Comment on "The Steiner number of a graph" by G. Chartrand and P. Zhang, *Discrete Math.* **242** (2002), 41 – 54.
- [11] M. Raines, P. Zhang, The Steiner distance dimension of graphs, *Australasian J. Combin.* **20** (1999), 133 – 143.

# Time Slicing Approach for Resource Allocation in Cloud Computing

M.AHARONU  
Assistant Professor  
Department of CSE  
MRCE,Hyderabad,Telangana.  
[aharon.mattakoyya@gmail.com](mailto:aharon.mattakoyya@gmail.com)

V .DeviPriya  
Assistant Professor  
Department of CSE  
MRCE,Hyderabad,Telangana.  
[vaddi.devipriya@gmail.com](mailto:vaddi.devipriya@gmail.com)

Dr.Sunil Tekale  
Professor,  
Department of CSE  
MRCE, Hyderabad, Telangana  
[sunil.tekale2010@gmail.com](mailto:sunil.tekale2010@gmail.com)

## Abstract:

Cloud computing is an on-demand service resource which includes applications to data centers on a pay-per-use basis. In order to allocate these resources properly and satisfy users' demands, an efficient and flexible resource allocation mechanism is needed. Due to increasing user demand, the resource allocating process has become more challenging and difficult. One of the main focuses of research scholars is how to develop optimal solutions for this process. In this paper, an algorithm is proposed which will help to perform the allocation of resources in a better and optimized manner[1].

**Keywords** - Cloud Computing resource s allocation time sharing optimization Network

## 1. Introduction

Resource management is a major task in cloud computing and in any other computing environments[1][2]. Cloud computing attempts to provide cheap and easy access to computational resources, which include servers, networks, storage, and, possibly, services. Cloud providers have to efficiently manage, provide, and allocate these resources to provide services to cloud consumers based on service level agreements (SLAs) which both sides agree to prior to the consumer using the services[1][3]. Therefore, providers must maintain a reliable allocating mechanism in order to satisfy the cloud users' requirements, while stabilizing an appropriate profit margin for themselves. Due to the increasingly high use of the Internet, and thereby cloud services, the typically static allocation and management of resources have become impractical, and the development of dynamic mechanisms have become more appropriate and worth studying[1].

Nonetheless, even these dynamic mechanisms present issues and challenges to be overcome and solutions to be found. Many researchers have tried, and are still trying, to provide the optimal solutions for the resource allocation and management problem in cloud computing environments[1].

The four deployment models present in cloud computing are: <sup>2</sup>

**1. Public cloud:** In the public cloud, the cloud provider provides resources for free to the public. Any user can make use of the resources; it is unrestricted. The public cloud is connected to the public internet for anyone to leverage[4].

**2. Private cloud:** In a private cloud, the planning and provisioning of the cloud are operated and owned by the organization or the third party. Here the hosted services are provided to a restricted number of people or group of individuals.

**3. Community cloud:** These type of cloud infrastructures exist for special use by a group of users. These are a group of users who share a common mission or have specific regulatory requirements, and it may be managed by the third party or organizations.

**4. Hybrid Cloud:** Hybrid Cloud provides the best of above worlds. It is created by combining the benefit of different types of cloud (private cloud & public cloud). In these clouds, some of the resources are provided and managed by public cloud and others as a private cloud.

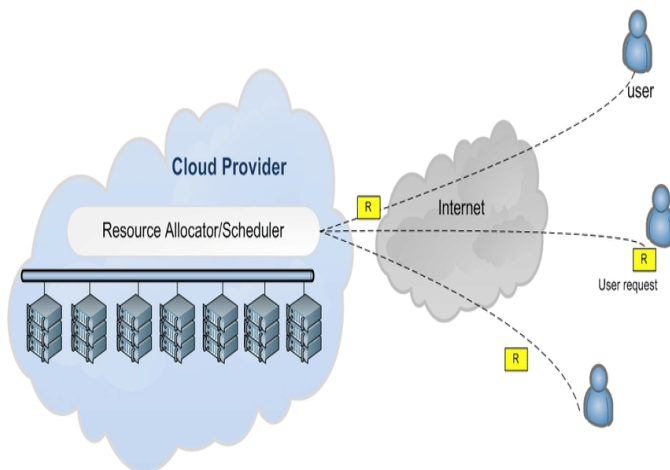
The three different service models present in cloud computing are:

**1. Infrastructure as a Service (IaaS):** IaaS model provides just the hardware and the network. It allows users to develop and install their operating system, software and run any application as per their needs on cloud hardware of their own choice[5][6].

**1. Platform as a Service (PaaS):** In PaaS model, an operating system, hardware, and network are provided to the user. It enables users to build their applications on cloud making use of supplier specific tools and languages

**2. Software as a Service (SaaS):** In SaaS model, a pre-built application together with any needed software, hardware, operating system and the network is provided to the user.

**Fig shows Resource allocator with users demand.**



## RESEARCH CHALLENGES

The research on RA in cloud systems is still at an early stage. Several existing issues have not been fully addressed while new challenges keep emerging. Some of the challenging research issues are given as follows:

1) Migration of VM: This migration problem occurs due to the need of the user to switch to another provider in order to get better data storage[1][8].

2) Control: There often is a lack of control mechanism over the resources as they are rented by the users from the remote server[7].

3) Energy Efficiency: Due to the emergence of huge data centers that have various computing operations, there is a need for energy efficient allocation. These centers lead to the release of large quantities of carbon emission.

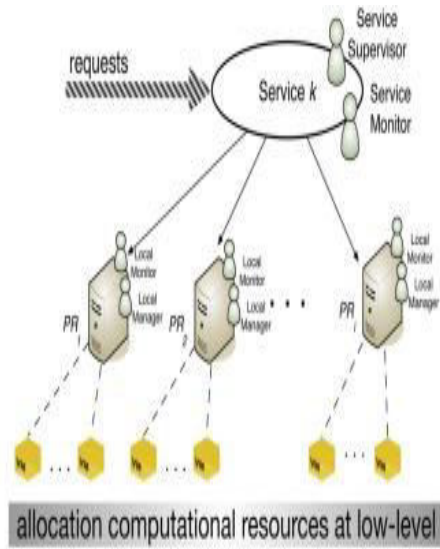
4) The Scheduling of Parallel Jobs: Parallel jobs in the field of computing increase the job that is serve. There are two types of jobs: dependent and independent. The first type must be done very carefully. These jobs include communication issues. Independent jobs can be performed using several VMs at the same time.

5) Reduction of Cost and Maximizing of Resources: It is important to handle the constraints that must be met in the allocation of resources in terms of cloud operating costs and to maximize the use of all resources. In other words, the service provider must provide users with low-cost services.

6) Maintaining High Availability: The availability of resources in the cloud must be guaranteed in case there is a job with long running computations that can take many hours. Thus, there is a need for some techniques to automatically handle any interruption or unavailability in resources and switch the jobs to an available resource. these techniques should support the transparency property by which the user cannot observe the unavailability or any failure problem[4][5][6][7]

7) Elasticity: In the cloud, elasticity refers to what extent resource requirements can be handled dynamically. Demand for resources may increase over time, and the cloud should automatically detect the size of these demands to be met and the necessary resources required to meet them.

**Fig shows allocation of Resources.**



### Proposed Algorithm :

Resource allocation is one of the most important aspect of cloud computing. As the number of resources will be less and demand for resources will be high, there needs to be a perfect policy for allocation of these resources. Resources allocation becomes difficult when allocation of resources needs to be done in peak time, which is going to be a big challenge. An algorithm is proposed to handle the issues of allocation of resources at peak time and on demand basis. The resources will be initially allocated to task which comes on first come first basis and then the same resources are to be allocated to new incoming task ,which means there will be multiple tasks handles by the same resource by using the concept of time slicing internally. These makes allocation of resources to the jobs and also helps in making them to wait for a minimum amount of time. A set of  $n$  jobs is to be processed on a single machine. The machine is available at time zero, and pre-emption of jobs is not allowed.

The  $i$ -th job is characterized by its processing time,  $P_i$ , and a "weight"  $w_i$  determining its priority. In the general case the precedence relations between the jobs are specified. The problem is to find an order  $(\pi_1, \dots, \pi_n)$  of jobs minimizing the total weighted completion time:

$$\sum_{i=1}^n w_i \sum_{j=1}^i P_{\pi_j}.$$

In the special case when the jobs have no precedence relations specified (actually, this is the case considered by Smith [a2]), the following Smith rule solves the problem: any sequence putting the jobs in order of non-decreasing ratios  $P_j / w_j$  is optimal, [2]. Adding arbitrary precedence constraints between jobs results in  $\mathcal{NP}$ -hardness, and the

same happens as soon as arbitrary arrival dates or deadlines are added to the model [a10].

**Scheduling Problems** Suppose that  $m$  machines  $M_j$  ( $j = 1, \dots, m$ ) have to process  $n$  jobs  $J_i$  ( $i = 1, \dots, n$ ). A schedule is for each job an allocation of one or more time intervals to one or more machines. Schedules may be represented by Gantt charts as shown in Figure 1.1. Gantt charts may be machine-oriented (Figure 1.1(a)) or job-oriented (Figure 1.1(b)). The corresponding scheduling problem is to find a schedule satisfying certain restrictions.

The general resource allocation model is below. When the parameters are given specific numerical values the result is an instance of the general model.

#### Parameters.

$n$ : Number of activities. Activities are indexed by  $J=1..n$

$m$ : Number of Resources .Resources are indexed by  $1..m$

$P_j$  : Profit for activity  $J$

$b_i$  : Amount available Resource  $i$

$a_{ij}$ : Amount of Resource  $i$  used by a unit of activity  $j$

Variables:

$X_j$  : Amount of activity  $j$  selected

Model

$$\text{Maximize profit} = \sum_{C_j X_j}^n$$

Subject to

$$\sum_{j=1}^n a_{ij} x_j \leq b_i \text{ for } i = 1..m$$

$$X_j \geq 0 \text{ for } J=1..n$$

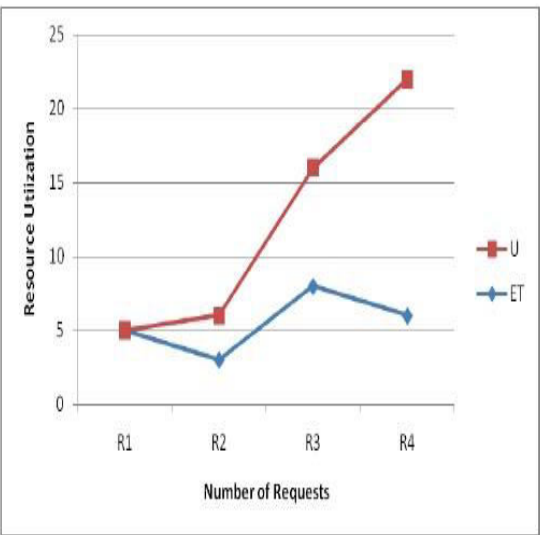
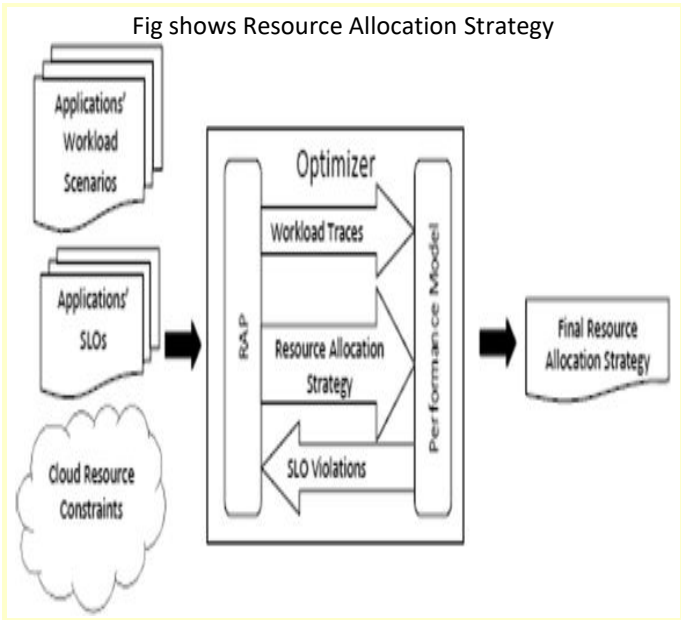


Fig: 3Shows Resource Utilization

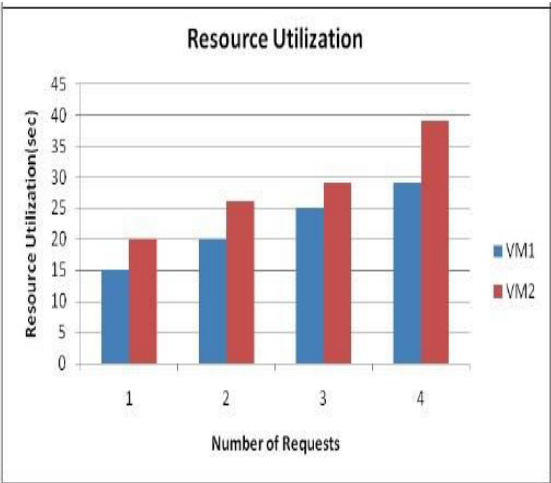


Fig shows no. of Resources Vs Utilization

### Conclusion:

Resource Allocation is the major issue in cloud computing which can be easily handled by the proposed algorithm, where we can assign multiple task to the same resource there by performing time slicing operation internally. This algorithm should work in a best possible manner. This helps to utilize the resources in a optimized way and the performance of the system will improve as all the incoming task are easily assigned to resources which are minimum in number. Proposed algorithm will help to solve the major issues which arises in allocation of Resources when the demand for the resources is high.

### References

- [1] R. Buyya, C. S. Yeo, S. Venugopal, J. Broberg and I. Brandic, "Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility", Future generation computer systems, vol. 25, no. 6, pp. 599–616, June 2009.
- [2] Kousik Dasgupta, Brototi Mandal, Paramartha Dutta, Jyotsna Kumar Mondal, Santanu Dam, "A Genetic Algorithm (GA) based Load Balancing Strategy for Cloud Computing", Elsevier, International Conference on Computational Intelligence: Modeling Techniques and Applications (CIMTA) 2013, 340-347

- [3] Kun Li, Gaochao Xu, Guangyu Zhao, Yushuang Dong, Dan Wang, "Cloud Task scheduling based on Load Balancing Ant Colony Optimization", IEEE, 2011, 978-0-7695- 4472-4/11
- [4] LI Kun – lun, Wang Jun, Song Jian, Song Jian, Dong Qing – yun, "Improved GEP Algorithm for Task Scheduling in Cloud Computing",
- [9] E. Gelenbe and C. Morfopoulou, A Framework for Energy Aware Routing in Packet Networks. accepted for publication in The Computer Journal.
- [10] E. Gelenbe and T. Mahmoodi, Energy-Aware Routing in the Cognitive Packet Network In International Conference on Smart Grids, Proceeding of Energy 2011 conference, 2011.

Second International Conference on Advanced Cloud and Big Data, IEEE, 2014, 978-1-4799- 8085-7/14 : 2231-2307, Volume-2, Issue-4, September 2012

[5] Rajveer Kaur, Supriya Kinger, "Enhanced Genetic Algorithm based Task Scheduling in Cloud Computing", International Journal of Computer Applications, Volume 101– No.14, September 2014, (0975 – 8887)

[6] Chun-Yan LIU, Cheng-Ming ZOU, Pei WU," A task scheduling algorithm based on genetic algorithm and ant colony optimization in cloud computing", 13th International Symposium on Distributed Computing and Applications to Business, Engineering and Science, IEEE, 2014, 978-1-4799-4169-8/14

[7] Liping Zhang, Weiqin Tong, Shengpeng Lu, "Task scheduling of cloud computing based on Improved CHC algorithm", ICALIP, IEEE, 2014, 978-1-4799-3903-9/14 Zhen Xiao, "Dynamic Resource Allocation Using Virtual Machines for Cloud Computing Environment", IEEE Transactions on Parallel and Distributed Systems, Vol. 24, No. 6, JUNE 2013.

[8] A. Berl, E. Gelenbe, M. di Girolamo, G. Giuliani, H. de Meer, M.- Q. Dang, and K. Pentikousis. Energy-Efficient Cloud Computing. The Computer Journal, 53(7), September 2010, doi:10.1093/comjnl/bxp080.

# A Novel Weighted Scan-Based Test Pattern For Built-In Self-Test

D.Gaspin Beautly

PG Student,ME VLSI Design,Marthandam College Of Engineering and Technology,Tamilnadu—629177

Mail.id : gaspinbeautly@gmail.com

**Abstract-** In this paper a new LP BIST method has been proposed using weighted test-enable signal- based pseudorandom test pattern generation and LP deterministic BIST and reseeding. In the existing systems, more power is consumed since all the scan chains are active in both the phases. To overcome this drawback and to design a low power BIST this system is proposed. This new method consists of two separate phases namely, LP weighted pseudorandom pattern generation and LP deterministic BIST with reseeding. The first phase selects weights for test-enable signals of the scan chains in the activated sub circuits. A new procedure has been proposed to select the primitive polynomial and the number of extra inputs injected at the LFSR. A new LP reseeding scheme, which guarantees LP operations for all clock cycles, has been proposed to further reduce test data kept on-chip.

**Index Terms**—Low-power (LP) built-in self-test (BIST), reseeding, scan-based BIST, weighted test-enable signals.

## I. INTRODUCTION

The gap between functional and test power consumption is growing bigger and bigger, with the latter reaching 2X to 5X of the former due to the ever-shrinking functional power and ever-increasing test power. Problems, such as excessive heat that may reduce circuit reliability, formation of hot spots, difficulty in performance verification, reduction of the product yield and lifetime, and so on, have become severe. More details on how to provide more accurate power model can be found from previous paper. A fast simulation approach was proposed for low-power (LP) off-chip interconnect design in [1]. An important through silicon via (TSV) modeling/simulation technique for LP 3-D stacked IC design was presented in [12]. Furthermore, the power dissipation of scan-based built-in self-test (BIST) is much higher than power dissipation in deterministic scan testing due to excessive switching

important. Weighted pseudorandom testing schemes can effectively improve fault coverage. However, these approaches usually result in much more power consumption due to more frequent transitions at the scan flip flops in many cases. Therefore, we intend to propose an LP scan-based pseudorandom pattern generator (PRPG). This is one of the major motivations of this paper.

Most of the previous deterministic BIST approaches did not include LP concerns. We intend to present a new method that effectively combines an efficient LP PRPG and LP deterministic BIST. In order to reduce test power in deterministic BIST, we will propose a new LP reseeding scheme, since there is no other effective solution in this field. This is another motivation of this paper.

In this paper, we propose a new LP scan-based BIST architecture, which supports LP pseudorandom testing, LP deterministic BIST and LP reseeding. We present the major contributions of this paper in the following.

- 1) A new LP weighted pseudorandom test pattern generator using weighted test-enable signals is proposed using a new clock disabling scheme. The design - for - testability (DFT) architecture to implement the LP BIST scheme is presented. Our method generates a series of degraded subcircuits. The new LP BIST scheme selects weights for the test-enable signals of all scan chains in each of the degraded subcircuits, which are activated to maximize the testability.
- 2) A new LP deterministic BIST scheme is proposed to encode the deterministic test patterns for random pattern- resistant faults. Only a part of flip flops are activated in each cycle of the whole process of deterministic BIST. A new procedure is proposed to select a primitive polynomial and the number of extra variables injected into the linear-feedback shift register (LFSR) that encode all deterministic patterns. The new LP reseeding scheme can cover a number of vectors with fewer care bits.

which allows a small part of flip flops to be activated in any clock cycle.

The rest of this paper is organized as follows. The related work is presented in Section II. The new LP weighted pseudorandom test generation approach is described in Section III. The new LP deterministic BIST method with reseeding is presented in Section IV. Experimental results are shown in Section V. This paper is concluded in Section VI.

## II. RELATED WORK

Scan flip flops, especially, the ones close to the scan-in pins, are not observable in most of shift cycles. A novel BIST scheme that inserts multiple capture cycles after scan shift cycles during a test cycle. Thus, the fault coverage of the scan-based BIST can be greatly improved. An improved method of the earlier work, presented in [2], selects different numbers of capture cycles after the shift cycles. In this paper, a new LP scan-based BIST technique is proposed based on weighted pseudorandom test pattern generation and reseeding. A new LP scan architecture is proposed, which supports both pseudorandom testing and deterministic BIST.

Weighted pseudorandom testing schemes can effectively improve fault coverage. A weighted test-enable signal-based pseudorandom test pattern generation scheme was proposed for scan-based BIST in [6], according to which the number of shift cycles and the number of capture cycles in a single test cycle are not fixed. A reconfigurable scan architecture was used for the deterministic BIST scheme in using the weighted test-enable signal-based pseudorandom test generation scheme. The proposed a new scan segmentation approach for more effective BIST.

LP BIST approaches were proposed early in a distributed BIST control scheme in order to simplify the BIST execution of complex ICs. The average power was reduced and the temperature was reduced. The methods reduced switching activity during scan shifts by adding extra logic. A new random single-input change test generation scheme in generates LP test patterns that provide a high level of defect coverage during LP BIST of digital circuits. An LP BIST scheme was proposed based on circuit partitioning.

New pseudorandom test generators were proposed to reduce power consumption during testing. A new encoding scheme is proposed in [3], which can be used in conjunction with any LFSR-

reseeding scheme to significantly reduce test power and even further reduce test data volume. A new LP PRPG for scan-based BIST using a restricted scan chain reordering method to recover the fault coverage loss. A low-transition test pattern generator in was proposed to reduce the average and peak power of a circuit during test by reducing the transitions among patterns. Transitions are reduced in two dimensions: 1) between consecutive patterns and 2) between consecutive bits. The [1] proposed a PRPG to generate test vectors for test-per-scan BISTs in order to reduce the switching activity while shifting test vectors into the scan chain. Furthermore, a novel algorithm for scan-chain ordering has been presented. A new adaptive low shift power pseudorandom test pattern generator was presented to improve the tradeoff between test coverage loss and shift power reduction in logic BIST. This is achieved by applying the information derived from test responses to dynamically adjust the correlation among adjacent test stimulus bits. The proposed LP programmable generators capable of producing pseudorandom test patterns with desired toggling levels.

A new LP BIST technology that reduces shift power by eliminating the specified high frequency parts of vectors and also reduces capture power. A novel approach to reduce peak power and power droop during capture cycles in scan based logic BIST. An efficient BIST architecture was recently presented in [3] for targeting defects in dies and in the interposer interconnects.

A novel low-power BIST technology was proposed in [4] that reduces shift power by eliminating the specified high frequency parts of vectors and also reduces capture power. Multi cycle tests support test compaction by allowing each test to detect more target faults. The ability of multi cycle broadside tests to provide test compaction depends on the ability of primary input sequences to take the circuit between pairs of states that are useful for detecting target faults. This ability can be enhanced by adding DFT logic that allows states to be complemented in [3].

A new DFT scheme for launch-on-shift testing was proposed which ensures that the combinational logic remains undisturbed between the interleaved capture phases, providing computer-aided-design tools with extra search space for minimizing launch-to-capture switching activity through test pattern ordering.

Complete fault coverage can be obtained [9] when the pseudorandom test generator is modified. A combination of a pseudorandom test generator and a

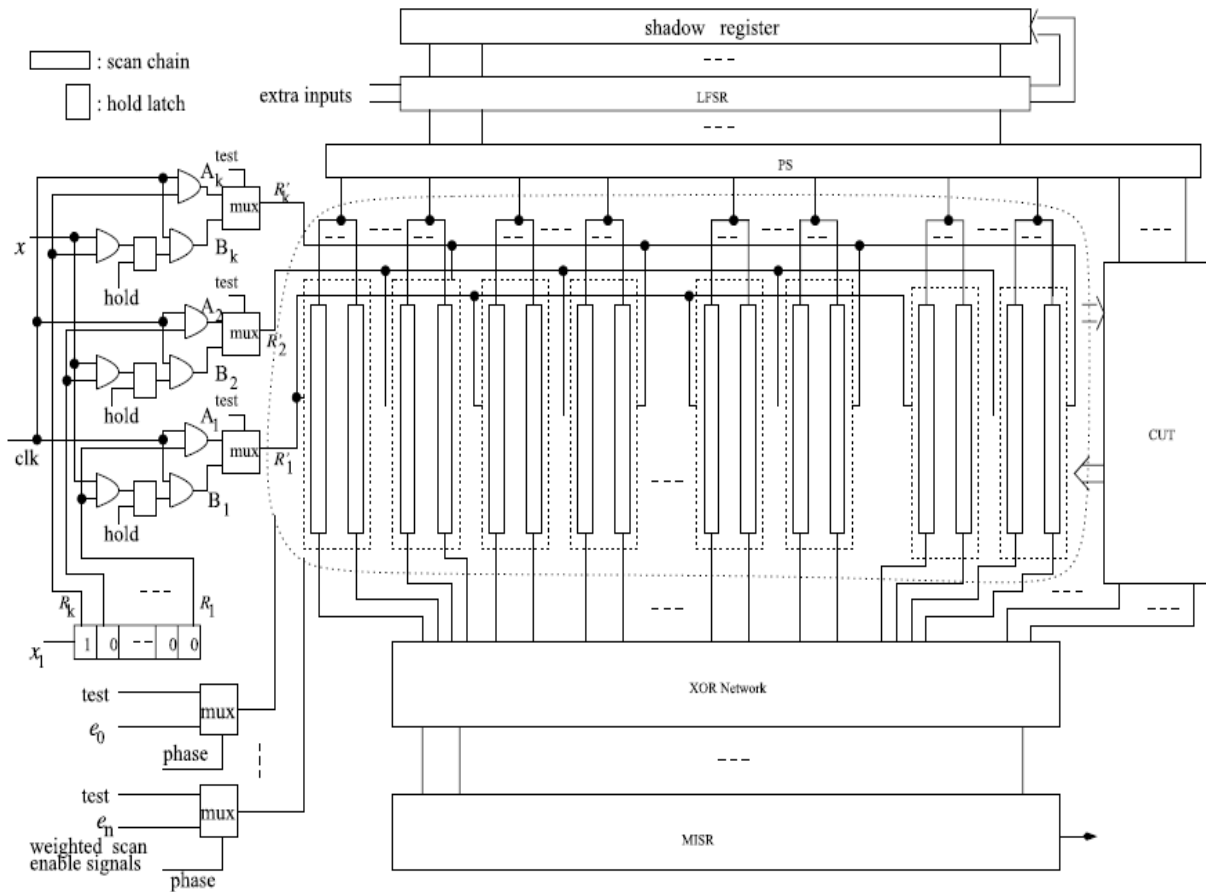


Fig.1. General DFT architecture for LP scan-based BIST.

combinational mapping logic was constructed by produce a given target pattern set of the hard-to-detect faults.

Deterministic vectors can be encoded into LFSR seeds which was proposed the seminal work, which encoded deterministic vectors into seeds. The requirement on the average size of the LFSR can be reduced by using multiple primitive polynomials .Deterministic vectors were encoded by using a folding counter and compressed by a tree architecture in a reconfigurable scan architecture for effective deterministic BIST. LP design was implemented in the new methodology in [2] to increase the encoding efficiency by combining reseeding and bit fixing.

The work is about LP delay testing, whose scan architecture and test application scheme are completely different from the new method. Our method is about scan-based BIST for single stuck-at faults based on a new weighted pseudorandom test generator and an LP deterministic BIST approach. The scan architecture is quite similar to the method in

[2]. Both methods do not require test response shift-out operations, which do not cause zero aliasing.

The proposed novel X-filling method by assigning 0 and 1 s to unspecified (X) bits in a test cube obtained during ATPG. This method reduces the circuit switching activity in capture mode and can be easily incorporated into any test generation flow to achieve capture power reduction without any area, timing, or fault coverage impact. A new scan shifting method based on the clock gating of multiple groups was proposed in [45] by reducing the toggle rate of the internal combinational logic. This method prevents cumulative transitions caused by shifting operations of the scan cells, because all scan flip flops are connected to the XOR network for test response compaction.

It is possible to implement LP scan testing in a test compression environment without any increase on test application cost which was proposed a new scan architecture to compress test data and compact test responses for delay testing. An important TSV

modeling/simulation technique for LP 3-D stacked IC design was presented in [4]. The connectivity of TSVs in many important circuits also needs to be tested in an efficient way.

### III. NEW LOW-POWER WEIGHTED PSEUDORANDOM PATTERN TEST GENERATOR

We present the DFT architecture to implement the LP BIST method in Section III-A. The process to implement LP pseudorandom pattern generation is presented in Section III-B.

#### A. DFT Architecture

As shown in Fig. 1, the scan-forest architecture [57] is used for pseudorandom testing in the first phase. Each stage of the phase shifter (PS) drives multiple scan chains, where all scan chains in the same scan tree are driven by the same stage of the PS. Unlike the multiple scan-chain architecture used in the previous methods the scan-forest architecture is adopted to compress test data and reduce the deterministic test data volume. Separate weighted signals  $e_0, e_1, \dots, e_n$  are assigned to all scan chains in the weighted pseudo-random testing phase (phase = 0), as shown in Fig. 1, which is replaced by the regular *test* in the deterministic BIST phase (phase = 1). Each scan-in signal drives multiple scan chains, as shown in Fig. 1, where different scan chains are assigned different weights. This technique can also significantly reduce the size of the PS compared with the multiple scan-chain architecture where each stage of the PS drives one scan chain. The compactor connected to the combinational part of the circuit is to reduce the size of the MISR. The shadow register is used for LP deterministic and reseeding, more details of which are described in Section IV-B.

The size of the LFSR needed for deterministic BIST depends on the maximum number of care bits of all deterministic test vectors for most of the previous deterministic BIST methods. In some cases, the size of the LFSR can be very large because of a few vectors with a large number of care bits even when a well-designed PS is adopted. This may significantly increase the test data volume in order to keep the seeds. This problem can be solved by adding a small number of extra variables to the LFSR or ring generator [10] without keeping a big seed for each vector.

We propose a new weighted PRPG for the new LP BIST approach. The new design is significantly different from the ones in [5] and [7]. This is mainly because the proposed LP design uses

the gating technique to disable most of the scan chains, where the pseudo primary inputs (PPIs) of the disabled scan chains are set to constant values. As shown in Fig. 1, all scan chains in the same scan tree are selected into the same subset of scan chains, which are driven by the same clock signal. Our method selects weights for each scan chain in the degraded subcircuits. Let the scan chains be partitioned into  $k$  subsets, where only one subset of scan chains is activated in any clock cycle. Our method selects optimal weights for all scan chains in the subset of scan chains in each round. It requires  $k$  separate rounds to determine optimal weights for all scan chains.

#### B. Weighted Pseudorandom Test Pattern Generation

Our method generates the degraded subcircuits for all subsets of scan chains in the following way. All PPIs related to the disabled scan chains are randomly assigned specified values (1 and 0). Note that all scan flip flops at the same level of the same scan tree share the same PPI. For any gate, the gate is removed if its output is specified; the input can be removed from a NAND, NOR, AND, and OR gates if the input is assigned a noncontrolling value and it has at least three inputs. For a two-input AND or OR gate, the gate is removed if one of its inputs is assigned a noncontrolling value. For a NOR or NAND gate, the gate degrades to an inverter if one of its inputs is assigned a noncontrolling value.

For an XOR or NXOR gate with more than three inputs, the input is simply removed from the circuit if one of its inputs is assigned value 0; the input is removed if it is assigned value 1, an XOR gate changes to an NXOR gate, and an NXOR gate changes to an XOR gate. For an XOR gate with two inputs, and one of its inputs is assigned value 0, the gate is deleted from the circuit. For a two-input NXOR gate, the gate degrades to an inverter. If one of its inputs is assigned value 1, a two-input XOR gate degrades to an inverter. If one of its inputs is assigned value 1, a two-input NXOR gate can be removed from the circuit. We first propose a new procedure to generate the weights of the test-enable signals for all scan chains in the LP DFT circuit after the degraded subcircuits for each subset of scan chains, which are driven by a single clock signal, have been produced. The  $i$ -controllability  $C_i(l)$  ( $i \in \{0, 1\}$ ) of a node  $l$  is defined as the probability that a randomly selected input vector sets

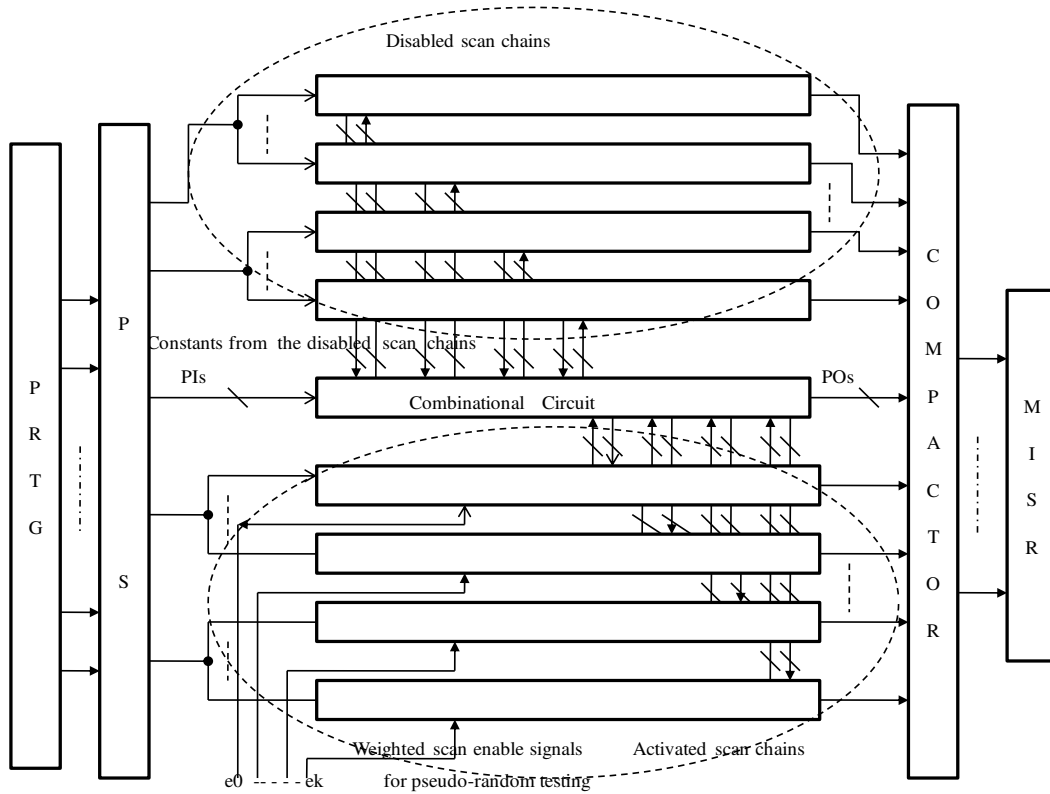


Fig.2. Weighted pseudorandom test generator for scan-tree-based LP BIST.

$l$  to the value  $i$ . The observability  $O_{-}(l)$  is defined as the probability that a randomly selected input vector propagates the value of  $l$  to a primary output. The *signal probability* of a node is defined in the same manner as its 1-controllability measure.

In the scan-based BIST architecture, as shown in Fig. 2, different weights  $e_0, e_1, \dots$ , and  $e_k$  are assigned to the test-enable signals of the scan chains  $SC_0, SC_1, \dots$ , and  $SC_k$ , respectively, where  $e_0, e_2, \dots, e_k \in \{0.5, 0.625, 0.75, 0.875\}$ . Scan flip flops in all disabled scan chains are set to constant values. Our method randomly assigns constant values to all scan flip flops in the disabled scan chains. The circuit is degraded into a smaller sub circuit. All

weights on the test enable signals are selected in the degraded subcircuit.

The gating logic is presented in Fig. 1. We do not assign weights less than 0.5 to the test-enable signals, because we do not want to insert more capture cycles than scan shift cycles. We have developed an efficient method to select weights for the test-enable signals of the scan chains. The selection of weights for the test-enable signals is determined by the following testability gain function:

$$G = \sum_{l/i \in F} \frac{|c'_1 - c'_0(l)|}{o'(l)} \quad (1)$$

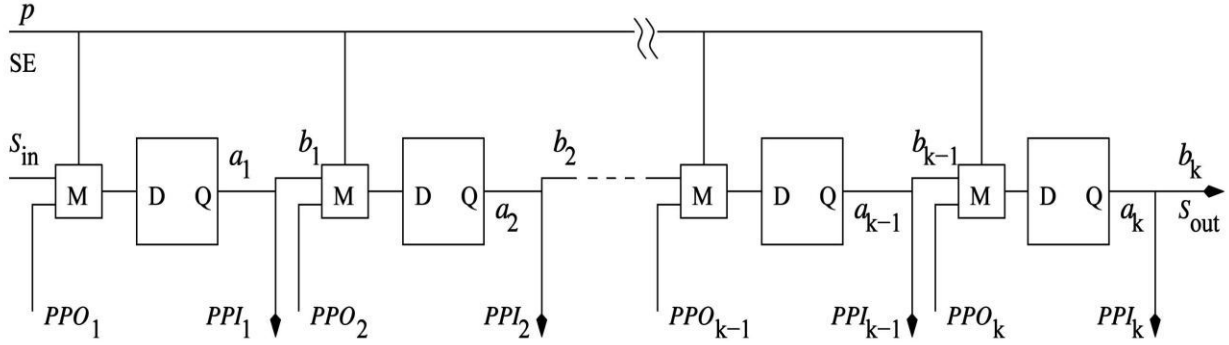


Fig. 3. Scan chain with a weighted test-enable signal.

where  $l/i$  represents the stuck-at  $i$  ( $i \in \{0, 1\}$ ) fault at line  $l$ . In (1),  $F$  is the random-pattern-resistant fault set, defined as the set of faults whose detection probability is no more than ten times that of the hardest fault [6]. We attempt to minimize the testability gain function as given in (1).

Fig. 3 presents a scan chain with a weighted test-enable signal, where  $S_{in}$ ,  $S_{out}$ , and test are the scan-in signal, scanout signal, and the test-enable signal of the scan chain, respectively. Initially, all PPIs are assigned signal probability 0.5, and the observability of the pseudoprimary outputs (PPOs) is set to  $1/n$ . Let  $p$  be the selected weight of the test-enable signal, as shown in Fig. 3. Then

$$C'_1(PPI_i) = p \cdot C'_1(a_{i-1}) + (1-p) \cdot C'_1(PPO_i). \quad (2)$$

The observability of  $PPO_i$  (PPO) can be estimated as follows:

$$O'(PPI_i) = (1-p) \cdot O'(a_i) \quad (3)$$

$$O'(a_i) = (1-(1-O'(b_i)) \cdot (1-O'(PPI_i))) \quad (4)$$

$$O'(b_{i-1}) = p \cdot O'(a_i). \quad (5)$$

The observability of the scan-out signal is set to 1. Even though the output of a scan chain is connected to the test response compactor, we can achieve zero aliasing by carefully connecting the scan chains to the XOR gates based on very simple structural analysis. Testability measures of the internal nodes, PPIs, and PPOs can be calculated iteratively using the controllability/observability program (COP) measures [11] and (2)–(5). We find that the testability measures for all nodes in the benchmark circuits converge within a few iterations.

The controllability of the PPI of the  $i$ th scan flip flop in a scan chain is set to 0.5, and the observability of the PPO of the  $i$ th scan flip flop is

set to  $1/d$ , where  $d$  is the length of the scan chains. Iterative testability estimation is adopted for all nodes based on (2)–(5) and the COP measure. It is found that testability measures for all nodes become stable after a quite few rounds of testability calculation.

Separate weights from the set  $\{0.5, 0.625, 0.75, 0.875\}$  are assigned to the test-enable signals of the scan chains. Algorithm 1 presents the details to select weights for all scan chains that are in the same scan chain subset. The inputs of Algorithm 1 are the scan chain set  $SC$ , which are partitioned into subsets  $\{SC_0, SC_1, \dots, SC_{k-1}\}$ . Our method generates  $k$  degraded subcircuits for each scan chain subset  $SC_i \in \{SC_0, SC_1, \dots, SC_{k-1}\}$ .

#### Algorithm 1 Select-Weights-for-Test-Enables()

- ```

{
1) Assign the same values to the scan enable
   signals as the regular test-per-scan BIST
   scheme to all scan segments.
2) While the scan chain set  $s \neq 0$ , do
3) Select a scan chain  $SC$  from the scan chain
   set  $S = S - \{SC\}$ 
4) For each weight in
    $\{0.5, 0.625, 0.75, 0.875\}$ , testability
   estimation is adopted to evaluate the cost
   function as presented in equation (1)
5) Select the best weight
    $W \in \{0.5, 0.625, 0.75, 0.875\}$  that makes the
   cost function as presented in equation (1)
6) For each scan chain, if no weight can be
   selected just leave its scan enable signal as
   the one in the conventional test-per-scan
   test scheme
}
```

Our method selects a weight for the first scan chain testenable signal to minimize the gain function. After the best weight has been selected for the first scan chain, a weight for the test-enable signal of the second scan chain is selected to minimize the cost function in (1). If no weight can be selected for any scan chain, our method sets its test-enable signal to the same value as the one in the conventional *test-per-scan* BIST scheme (the number of shift cycles is equal to the length of the scan chains, and a capture cycle follows). Continue the above process until appropriate weights have been chosen for all test-enable signals of the scan chains in  $SC_i$ .

The proposed DFT architecture, as shown in Fig. 1, has an implicit advantage over other BIST architectures [1]. Each stage of the PS drives a scan tree [5] instead of a single scan chain, while each stage of the PS requires a few number of XOR gates. In any case, flip flops of all disabled scan chains are assigned with specified values. Therefore, no unknown signals are produced to corrupt the compacted test responses kept in the MISR.

Fig. 2 presents a degraded subcircuit based on the proposed LP BIST method. PPIs corresponding to scan flip flops in all disabled scan chains are assigned with randomly selected constant values in the period of weighted pseudorandom test pattern application for the current subset of scan chains. The proposed LP weighted pseudorandom test pattern generation process is as follows. The first subset of scan chains is activated when all the remaining scan chains are disabled. The generated weighted pseudorandom pattern is applied to the degraded subcircuits if a scan chain is set to the capture cycle; otherwise, the scan chain is set to the scan shift mode. Our method turns to the next phase when the second subset of scan chains is activated after the given number of clock cycles. This process continues until all subsets of scan chains have been processed. The first subset of scan chains is again activated, and the above process is executed again. The process continues until the whole given number of clock cycles has run over.

The proposed LP weighted pseudorandom test generator is shown to be able to improve fault coverage compared with the conventional test-per-scan BIST approaches according to the experimental results presented in the experimental result section. The amount of test data to be stored on-chip is also significantly reduced.

#### IV. LOW-POWER DETERMINISTIC BIST

We use the same LFSR for both pseudorandom pattern generation and deterministic phases. First, we propose a new algorithm to select a proper primitive

polynomial; after that the LP deterministic BIST and LP reseeding schemes are presented.

##### A. Selecting a Primitive Polynomial and the Extra Variable Number

Some extra variables are injected just like EDT [42]. We propose a new scheme to select the size of the LFSR and the number of extra variables simultaneously in order to minimize the amount of deterministic test data. Usually, a small LFSR constructed by a primitive polynomial is sufficient when a well-designed PS is adopted in the pseudorandom testing phase. In our method, a combination of a small LFSR and the PS from [4] is used to generate test patterns in the pseudorandom testing phase. The weighted test-enable signal-based pseudorandom test generator generates weighted pseudorandom test patterns. The size of the LFSR is not determined by the maximum number of care bits for any deterministic test vector. That is, the same LFSR is used for both phases.

For any degree less than 128, it is computationally feasible to generate enough primitive polynomials in reasonable time, out of which one (whose degree is equal to the maximum number of care bits in the deterministic vectors) can be selected to encode all deterministic test vectors. The tool that we used to generate primitive polynomials can only handle polynomials up to degree 128 of the word-length limit of the computer. However, only very small LFSRs are used for all circuits according to all experimental results (no more than 30). This is mainly because we inject some extra variables to the LFSR. To encode a few deterministic test vectors with a large number of care bits, the injected extra variables and the seed kept in the LFSR are combined just like the EDT tool. Therefore, it is not necessary to provide an LFSR whose size is at least the maximum number of care bits by injecting some extra variables.

A well-designed LFSR is needed in order to encode all deterministic vectors after the pseudorandom testing phase. A new procedure is proposed to select a primitive polynomial with the minimum degree that can encode all deterministic test vectors for the hard faults. An efficient algorithm is used to generate primitive polynomials of any desired degree. For any  $i \leq 30$ , assume that all primitive polynomials are kept in  $Q_i$ . As for  $i > 30$ , only a number of primitive polynomials are provided in  $Q_i$ . The following procedure returns a primitive polynomial with the minimum degree that encodes all deterministic vectors for the random pattern-resistant (hard) faults. Usually, the numbers of care bits of all deterministic test vectors is quite different. Therefore,

it is recommended to use an LFSR whose size is more than the maximum number of care bits  $S_{\max}$  of all deterministic vectors. Unlike the method in the new method selects a primitive polynomial of relatively low degree when some extra variables are injected into the LFSR. The commercial tool EDT used similar technique to reduce the amount of test data stored in the on-chip ROM or automatic test equipment (ATE).

## B. Low-Power Deterministic BIST and Reseeding

An effective seed encoding scheme is used here to reduce the storage requirements for the deterministic test patterns of the random-pattern-resistant faults. The encoded seed is shifted into the LFSR first. A deterministic test vector is shifted into the scan trees that are activated by the gating logic, where each scan-in signal drives a number of scan trees, and only one of the scan trees driven by the same scan-in signal is activated. The extra variables are injected into the LFSR when the seed is shifted into the activated scan trees. The gating logic, as shown in Fig. 1, partitions scan trees into multiple groups.

The first group of scan trees is disabled after they have received the test data. The second group of scan trees is activated simultaneously, and all other scan trees are disabled. The seed can be stored in an extra shadow register, which is reloaded to the LFSR in a single clock cycle. The scan shift operations are repeated when the extra variables are injected into the LFSR. This process continues until all scan trees have received test data.

Let us describe the details about constructing the scan forest. Assume that the number of scan flip flops at each level in the same scan tree is  $l$  and the depth of the scan forest is  $d$ . For a given scan-in pin,  $l$  scan flip flops are selected among all scan flip flops for the first level of the scan tree. The routing overhead is minimized when constructing the scan trees, which can be easily estimated using tools, such as *Astro* from synopsys [4]. Experimental results reported in this paper were obtained using the *Astro* tool. All scan flip flops at the same level in the same scan tree meet the following condition. Each pair of scan flip flops has no combinational successor in the circuit. Each scan flip flop  $p$  at the first level of the scan tree is connected to a scan flip flop  $f$  at the second level that has the minimum distance from  $p$  among all scan flip flops that can be placed at the second level of the scan tree, where all scan flip flops at the second level of the same scan tree have no common combinational successor. Repeat the above process until the scan trees have been constructed.

We propose an LP deterministic BIST scheme with reseeding. The deterministic test vectors for the random-pattern resistant faults are ordered according to the number of care bits. Our method partitions all scan chains into multiple subsets, while only one subset of scan trees is activated at any clock cycle. The gating logic controls the whole test application process. The first deterministic test vector is shifted into all scan trees as follows. The seed is first shifted into the LFSR. The extra variables with calculated values are injected into the LFSR when the seed is applied to the first subset of activated scan trees. The same values on the extra inputs are delivered after the same seed is loaded to the LFSR again for the second subset of activated scan trees. This process continues until all scan trees have received the test vector.

Our method turns to the reseeding process. The final values in the LFSR remain unchanged. The activated subset of scan trees performs  $d$  shift cycles when the extra variables with the same values are injected. The second subset of activated scan trees performs  $d$  shift cycles when the same values of the extra variables are injected. This process continues until the values of the extra variables have been shifted into all scan trees. Our method begins to check the values of the scan trees to see whether they are compatible with any remaining deterministic test vector. If so, the test vector is deleted from the ordered test sequence, and another LP capture period is applied as stated earlier from this state.

If the values kept in the scan chains are compatible with a deterministic vector, our method continues the responses capturing process. Assume that the initial values kept in the LFSR are stored in the shadow register. The first subset of scan trees is activated, which captures the test responses. The values kept in the shadow register are reloaded to the LFSR. The values of the extra variables are injected again when activated scan trees are filled. The above process continues until all scan trees have captured test responses.

If the values kept in the scan flip flops are incompatible with any other deterministic test vector, our method starts another LP shift-in period when injecting the extra variables that are stated earlier. The reseeding process continues until the given number of reseeding processes has been completed. In each round of the reseeding processes, the states of the scan trees are checked to see whether they are compatible with any deterministic test vector. If so, the deterministic vector is deleted. Our method copes with the second deterministic vector after the reseeding processes have been completed.

TABLE- 1

FAULT COVERAGE COMPARISON OF THE LP WEIGHTED PSEUDORANDOM TEST GENERATOR

| -          | The Proposed Method |        |        |        | [10]   |        |        |        |
|------------|---------------------|--------|--------|--------|--------|--------|--------|--------|
|            | FC                  | FC(10) | FC(20) | FC(30) | FC     | FC(10) | FC(20) | FC(30) |
| circuits   |                     |        |        |        |        |        |        |        |
| s38417     | 99.165              | 99.053 | 99.077 | 99.107 | 97.879 | 97.365 | 97.561 | 97.613 |
| b19        | 84.832              | 84.332 | 84.645 | 84.679 | 83.237 | 82.859 | 82.883 | 82.994 |
| wb_conmax  | 93.527              | 93.266 | 93.437 | 93.471 | 91.793 | 91.412 | 91.486 | 91.506 |
| usb_funct  | 92.811              | 92.742 | 92.787 | 92.798 | 92.016 | 91.621 | 91.795 | 91.814 |
| pci_bridge | 95.447              | 95.003 | 95.224 | 95.287 | 94.841 | 94.597 | 94.772 | 94.768 |
| des_perf   | 96.901              | 96.887 | 96.889 | 96.892 | 95.396 | 95.013 | 95.175 | 95.223 |
| ethernet   | 96.318              | 96.089 | 96.117 | 96.226 | 95.904 | 95.457 | 95.682 | 95.726 |
| vga_lcd    | 92.031              | 91.683 | 91.778 | 91.859 | 91.303 | 90.768 | 90.917 | 91.162 |
| netcard    | 95.194              | 93.982 | 94.337 | 94.756 | 94.546 | 93.349 | 93.651 | 94.173 |

V. EXPERIMENTAL RESULTS

The proposed method has been implemented and evaluated on a Dell Precision 7810 workstation. The pseudorandom testing phase was used with the scan-forest scan architecture, and separate weighted test-enable signals were assigned to the scan chains. A very small number of scan-in pins were used, making the size of the PS very small. That is, the area overhead can be reduced significantly.

Performance comparison for the proposed LP BIST scheme and the one in [10] is presented in Table I on the fault coverage of the pseudorandom test generators. The column FC shows the fault coverage of the original weighted pseudorandom test generator. The columns FC(10), FC(20), and FC(30) present the fault coverages of the proposed LP BIST method after 500k clock cycles, where the number given in the bracket shows the percentages of the activated scan flip flops for the proposed LP BIST method and the one presented in [10].

As for the circuit netcard, both LP BIST methods reach 93.98% and 93.35% fault coverage when only 10% scan chains are activated. The numbers of the deterministic test vectors for both methods are 173 and 192, respectively, and the final amount of on-chip data for the seeds is reduced approximately 26.8 times. It is shown that the number of maximum care bits of the deterministic vectors for both methods are 379 and 9873, respectively, which makes the amount of seeds to be kept on-chip completely different.

Fig. 3 presents the performance of the proposed LP PRPG for circuits netcard and vga when different percentages (10%, 20%, 30%, and 100%) of scan chains are activated. It is shown that the fault coverage is less when fewer scan chains are activated. Finally, the fault coverages for all four cases are quite close after 500 000 clock cycles. In a few cases, the fault coverage with 30% activated scan

chains is slightly more than that with 100% activated scan chains for the circuit netcard. This anomaly also occurs for the circuit vga, as shown in Fig. 3.

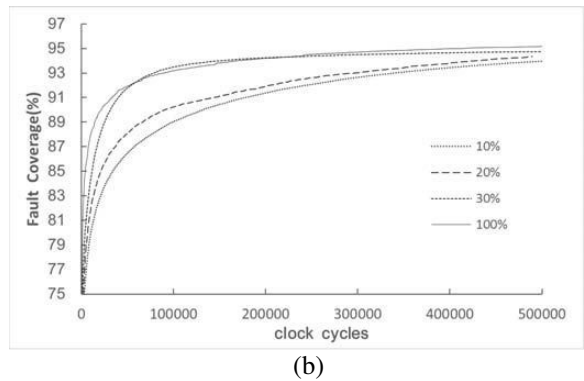
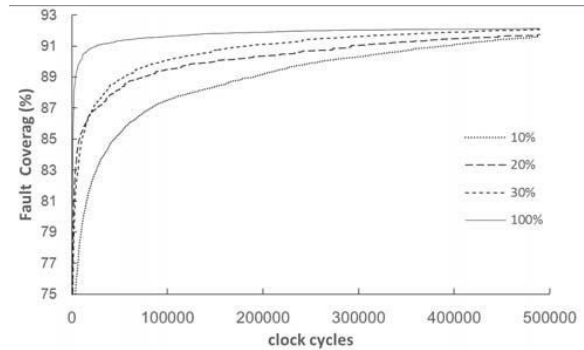


Fig. 5. Fault coverage with different toggle rates. (a) vga\_lcd. (b) Netcard

TABLE II

POWER REDUCTION

| -          | low-power deterministic BIST |                    |         | low-power PRPG   |                    |         |
|------------|------------------------------|--------------------|---------|------------------|--------------------|---------|
|            | peak(mW, before)             | lp peak(mW, after) | rate(%) | peak(mW, before) | lp peak(mW, after) | rate(%) |
| circuits   |                              |                    |         |                  |                    |         |
| s38417     | 6724                         | 643                | 9.6     | 7695             | 704                | 9.1     |
| b19        | 41933                        | 5660               | 13.4    | 45334            | 5779               | 12.7    |
| wb_conmax  | 10299                        | 1111               | 10.8    | 11197            | 1255               | 11.2    |
| usb_funct  | 6261                         | 610                | 9.7     | 6342             | 639                | 10.1    |
| pci_bridge | 8119                         | 984                | 12.1    | 8783             | 1014               | 11.5    |
| des_perf   | 22771                        | 2181               | 9.5     | 25552            | 2398               | 9.4     |
| ethernet   | 28865                        | 3746               | 12.9    | 29331            | 3823               | 13      |
| vga_lcd    | 46372                        | 4775               | 10.2    | 47432            | 4396               | 9.3     |
| netcard    | 95165                        | 9426               | 9.9     | 103676           | 9970               | 9.6     |

Table II presents the performance of the proposed LP deterministic BIST scheme on peak power (milli-Watt, mW) reduction when 10% scan chains are activated. The supply voltage and frequency are set to 1.5 V and 200 MHz, respectively. The column's peak (mW, before) and lp peak (mW, after) show the peak power for the original deterministic BIST and weighted test-enable-based PRPG, and the proposed LP BIST method. The column rate(%) shows the percentage of peak power for the proposed method compared with the one without the LP design for both the weighted pseudorandom test generation phase and the deterministic BIST phase. Experimental results in Table II show that the proposed LP PRPG phase reduces the peak power to less than 13% for all circuits, and the LP deterministic BIST scheme reduces the peak power to less than 14% in all cases. Experimental results show that the peak power for the PRPG phase is a little more than that for the deterministic BIST phase for the all circuits except s38417 before the LP design is included. This is mainly because only 10% flip flops are activated in any case during the LP weighted pseudorandom testing and the LP deterministic BIST phases, as shown in Fig. 1.

## VI. CONCLUSION

A new low-power (LP) scan-based built-in self-test (BIST) technique is proposed based on the weighted pseudorandom test pattern generation and reseeding. A new LP scan architecture is proposed in this paper, which supports both the pseudorandom testing and deterministic BIST. During pseudorandom testing phase, an LP weighted random test pattern generation scheme is proposed by disabling a part of scan chains. During the deterministic BIST phase, the design-for testability architecture is modified slightly while the linear-feedback shift register is kept short. In both the cases, only a small number of scan chains are activated in a single cycle thus low power scan based BIST technique is designed.

## REFERENCES

1. Abu-Issa A. S. and Quigley S. F. (2009), "Bit-swapping LFSR and scan-chain ordering: A novel technique for peak- and average-power reduction in scan-based BIST," *IEEE Trans. Comput.-Aided Des. Integr. Circuits Syst.*, vol. 28, no. 5, pp. 755–759.
1. Agrawal V. D., Kime C. R., and Saluja K. K. (1993), "A tutorial on built-in self-test.I. Principles," *IEEE Des. Test Comput.*, vol. 10, no. 1, pp. 73–82.
2. Al-Yamani A., Devta-Prasanna N., Chmelar E., Grinchuk M., and Gunda A.(2007), "Scan test cost and power reduction through systematic scan reconfiguration," *IEEE Trans. Comput.-Aided Des. Integr. Circuits Syst.*,vol. 26, no. 5, pp. 907–918.
3. Banerjee S., Chowdhury D. R, and Bhattacharya B. B. (2007), "An efficient scan tree design for compact test pattern set," *IEEE Trans. Comput.- Aided Des. Integr. Circuits Syst.*, vol. 26, no. 7, pp. 1331–1339.
4. Bardell P. H., Mc Anney W. H., and Savir J. (1987), *Built in Test for VLSI: Pseudorandom Techniques*. New York, NY, USA: Wiley.
5. Basturkmen N. Z., Reddy S. M., and Pomeranz I. (2003), "A low power pseudorandom BIST technique," *J. Electron. Test., Theory Appl.*, vol. 19, no. 6, pp. 637–644.
6. Bushnell M. L. and Agrawal V. D. (2000), *Essentials of Electronic Testing*. Norwell, MA, USA: Kluwer.
7. Chatterjee M. and Pradhan D.K. (2003), "A BIST pattern generator design for near-perfect fault coverage," *IEEE Trans. Comput.*, vol. 52, no. 12, pp. 1543–1558.
8. Filipek M. et al. (2015), "Low-power programmable PRPG with test compression capabilities," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 23, no. 6, pp. 1063–1076.
9. Gerstendörfer S. and Wunderlich H. J. (1999), "Minimized power consumption for scan-based BIST," *J. Electron. Test.*, vol. 16, no. 3, pp. 203–212.
10. Girard P., Landrault C., Pravossoudovitch S., Virazel A., and Wunderlich H. J. (2002), "High defect coverage with low-power test sequences in a BIST environment," *IEEE Des. Test Comput.*, vol. 21, no. 5, pp. 44–52.
11. Girard P., Guiller L., Landrault C., and Pravossoudovitch S. (2000), "Low power BIST design by hyper graph partitioning: Methodology and architectures," in *Proc. Int. Test Conf.*, pp. 652–661.
12. Hellebrand S., Rajski J., Tarnick S., Venkataraman S., and Courtois B. (1995), "Built-in test for circuits with scan based on reseeding of multiple polynomial linear feedback shift registers," *IEEE Trans. Comput.*, vol. 44, no. 2, pp. 223–233.

## ENERGY-EFFICIENT SECURE DATA AGGREGATION FRAMEWORK (ESDAF) PROTOCOL IN HETEROGENEOUS WIRELESS SENSOR NETWORKS

Dr. G. Silambarasan<sup>1</sup>, , Dr. V.Bhoopathy<sup>2</sup> Dr. V. Chandrasekar<sup>3</sup>

<sup>1</sup>Assistant Professor, Dept. of Computer Science and Engineering,  
The Kavery College of Engineering, Salem, Tamilnadu, India,

<sup>2</sup>Professor, Dept. of Computer Science and Engineering,  
Malla Reddy College of Engineering Secunderabad, Telangana State

<sup>3</sup>Associate Professor, Dept. of Information Technology,  
Malla Reddy College of Engineering and Technology, Secunderabad, Telangana State,  
[gssilambarasan@gmail.com](mailto:gssilambarasan@gmail.com), [v.bhoopathy@gmail.com](mailto:v.bhoopathy@gmail.com), [drchandru86@gmail.com](mailto:drchandru86@gmail.com)

**ABSTRACT** - Wireless Sensor Networks (WSNs) are constrained in terms of memory, computation, communication, and energy. In the existing secure data aggregation techniques, reduction in the energy consumption is not much discussed and combined solution for both integrity and authentication is not addressed. Data aggregation is a very important technique, but it gives extra opportunity to the adversary to attack the network, inject false messages into the network and trick the base station to accept false aggregation results. This paper presents an energy-efficient secure data aggregation framework (ESDAF) protocol WSN. The goal of the framework is to ensure data integrity and data confidentiality. ESDAF uses two types of keys. Base station shares a unique key with each sensor node that is used for integrity and the aggregator shares a unique key with each sensor node (within that cluster) that is used for data confidentiality. Sensor nodes calculate a message authentication code (MAC) of the sensed data using shared key with base station, which verifies the MAC for message integrity. Sensor nodes encrypt the sensed data using shared key with aggregator, which ensures data confidentiality. Proposed framework has low communication overhead as the redundant packets are dropped at the aggregators.

**Keyword:** Wireless Sensor Network (WSN), Message Authentication Code (MAC), Energy-Efficient Secure Data Aggregation Framework (ESDAF).

### 1. Introduction

#### 1.1. Wireless Sensor Networks

Wireless sensor networks comprises of the upcoming technology that has attained noteworthy consideration from the research community. Sensor networks comprise of many small, low cost devices and are naturally self organizing ad hoc systems. The function of the sensor network is monitoring the physical environment,

collect and transmit the information to other sink nodes. In general the range of the radio transmission for the sensor networks are in the orders of the magnitude which is smaller than the geographical extent of the intact network. Hence, the data has to be transmitted hop-by-hop towards the sink in a multi-hop manner. The consumption of energy in the network can be reduced if the amount of data to be relayed is reduced. [1].

Wireless sensor network comprises of a great number of minute electromechanical sensor devices which possess the sensing, computing and communication abilities. These devices can be utilized for gathering sensory information, like measurement of temperature from an extended geographical area [2].

Many of the features of the wireless sensor networks give rise to challenging problems [3]. The most important three characteristics are:

- Sensor nodes are the ones which are prone to maximum failures.
- Sensor nodes make use of the broadcast communication pattern and have severe bandwidth restraint.
- Sensor nodes have limited amount of resources.

### 1.2. Data Aggregation

Data aggregation is considered as one of the fundamental distributed data processing procedures for saving the energy and minimizing the medium access layer contention in wireless sensor networks [4]. Data aggregation is presented as an important pattern for routing in the wireless sensor networks. The basic idea is to merge the data from various sources, reroute it with the elimination of the redundancy and thus reducing the number of transmissions and saving the energy [5]. The inbuilt redundancy in the raw data gathered from various sensors can be prevented by the in-network data aggregation. Additionally, these operations use raw materials for obtaining application specific information. To preserve the energy in the system for

maintaining longer lifetime in the network, it is important for the network to maintain high incidence of the in-network data aggregation [6].

### 1.3. Secure Data Aggregation

The issues related to the security in the data aggregation of WSN are as follows [7]:

- **Data Confidentiality:** In particular, the basic security issue is the data confidentiality which safeguards the transmitted data that is sensitive from passive attacks like eavesdropping. The importance of the data confidentiality is in the hostile environment, where the wireless channel is more susceptible to eavesdropping. Even though cryptography has provided plenty of methods, the operation related to complicated encryption and decryption, like modular multiplication of large numbers in public key based cryptosystems, uses the sensor's power quickly.
- **Data Integrity:** It prevents the alteration of the final aggregation value by the compromised source nodes or aggregator nodes. Sensor nodes can be easily compromised due to the lacking of the expensive tampering-resistant hardware. The otherwise used hardware may not be reliable at times. A compromised message is capable of modifying, forging and discarding the messages.

In general, for secure data aggregation in wireless sensor networks, two methods can be used. They are hop by hop encrypted data aggregation and end to end encrypted data aggregation [7].

- **Hop-by-Hop encrypted data aggregation:** In this technique, the encryption of the data is performed by

the sensing nodes and decryption by the aggregator nodes. The aggregator nodes aggregate the data and again encrypt the aggregation result. At the end, the sink node on obtaining the final encrypted aggregation result decrypts it.

- End to End encrypted data aggregation: In this technique, the aggregator nodes in between have no decryption keys and can only perform aggregation on the encrypted data.

## 2. Related Work

Yingpeng Sang et al [7] have classified the security issues, data confidentiality and integrity in data aggregation into two cases: hop-by-hop encrypted data aggregation and end-to-end encrypted data aggregation. They have also proposed two general frameworks for these two cases respectively. The framework for end-to-end encrypted data aggregation has higher computation cost on the sensor nodes, but achieves stronger security, in comparison with the framework for hop-by-hop encrypted data aggregation.

Prakash G.L et al [8] have presented privacy-preserving data aggregation scheme for additive aggregation functions. The goal of their work is to bridge the gap between collaborative data collection by wireless sensor networks and data privacy. They have presented simulation results of their schemes and compared their performance to a typical data aggregation scheme TAG, where no data privacy protection is provided. Results show the efficacy and efficiency of their schemes. But, due to the algebraic properties of the polynomials, the communication overhead increases and becomes more complex.

Tamer AbuHmed et al [9] have presented a dynamic and secure scheme for data aggregation in WSN. Their proposal scheme includes level-based key derivation, data aggregation, and a new node join phases. Furthermore, they have done a security analysis for a related Level-based Key Management (LBKM) scheme proposed by Kim et al. Their analysis shows that LBKM is insecure for one node compromising and neighbor nodes misbehavior. To this end, they proposed different levelbased key management scheme for secure data aggregation. Their scheme is secure and more efficient than LBKM scheme in term of communication overhead and security. However, the proposed work is operated only in the tree based structure. Moreover, the overhead is greater in the case of the threshold cryptography.

Wenbo He et al [10] have presented two privacy-preserving data aggregation schemes for additive aggregation functions. Their first scheme is Energy Efficient Secure Data Aggregation (EESDA) which leverages the clustering protocol and algebraic properties of polynomials. Their second scheme is Slice-Mix-AggRegaTe (SMART) which builds on slicing techniques and the associative property of addition. The goal of their work is to bridge the gap between collaborative data collection by wireless sensor networks and data privacy. They assessed the two schemes by privacy-preservation efficacy, communication overhead, and data aggregation accuracy. Their Simulation results show the efficacy and efficiency of our schemes. But the bandwidth

consumption is increased in the case of their proposed SMART technique.

Shih-I Huang et al [11] have proposed a Secure Encrypted-data Aggregation (SEA) scheme in mobile wireless sensor networks (MWSN) environment. Their design for data aggregation eliminates redundant sensor readings without using encryption and maintains data secrecy and privacy during transmission. In contrast to conventional schemes, their proposed scheme provides security and privacy, and duplicate instances of original readings will be aggregated into a single packet; therefore, more energy can be saved. But integrity is not discussed in their proposed SEA scheme.

### 3. Secure Data Aggregation

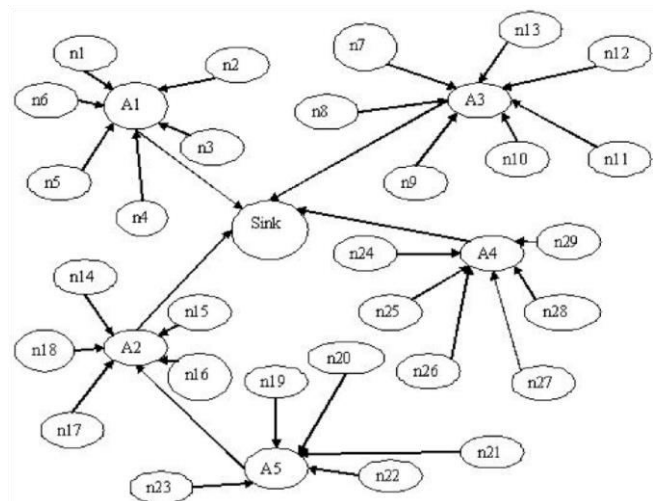
#### 3.1. System Overview

In a clustered WSN, the network is grouped into clusters. In each cluster, there is an aggregator which consists of a very powerful wireless transceiver that is capable of transmitting the data directly to the backend server. In our work, our assumption is that each sensor performs the transmission of the data only to the aggregator. As a result, each sensor will be able to reduce the overhead in transmitting the data packets. We assume that the sensor nodes do not have any mobility, i.e., the sensor nodes are all attached to a position and are cannot move.

The verification information is built by the source using the shared key. Verification information is included with data packet during the transmission. On reception of the packet, the source is

verified by the aggregator using the shared key. In case of failure in the verification, the packet will be discarded, otherwise it will be forwarded. On reception of the data packet by the sink, the source will be tested again for its validity. If the validity of the source fails then it will be discarded. A MAC based authentication code is used in order to maintain the integrity of the data packet. The sink can detect any changes performed by the aggregator including the verification information, by checking of the MAC value using its shared key. If the data packet is found to be modified, then it will be discarded.

The power consumption is reduced in our proposed data aggregation method, along with the maintenance of the secrecy and privacy. In case of the secrecy, encryption is performed by each sensor node and this encrypted data is then transmitted to the aggregator. Hence, it will not be possible for the adversaries to read the data packet.



**Figure 1:** System Architecture

### 3.2 Encryption and Decryption

After the selection of the aggregator, each sensor nodes communicate with the aggregator, *aggr* using a symmetric key  $K_{ch,i}$ . The sensor nodes send the encrypted data using this key to the *aggr*. Then the *aggr* receives the encrypted data and decrypts the data using the same key  $K_{ch,i}$ . Now the *aggr* identifies the malicious or compromised nodes, and filter out their data in the networks based on *MAC* function.

Each *aggr* determines a *MAC* value for the aggregated data and finally all the aggregated data are encrypted and transmitted to the sink. This data is encrypted using a symmetric key  $K_{ch,s}$ . The sink decrypts the received data using the same key  $K_{ch,s}$ .

### 3.3 Algorithm for the Aggregator

1. The sensors send its data to the nearest aggregator, *aggr* since each sensor node has a *aggr* to ensure its connectivity.
2. Each sensor node encrypts the data using the symmetric key  $K_{ch,i}$  and sends it to its *aggr*.
3. When *aggr* receives the data packet from any node  $S$ , it decrypts the data using the symmetric key  $K_{ch,i}$ .
4. The *aggr* then calculates the *MAC* using the hash functions  $MAC(aggr)$ .
5. By calculating the *MAC*, the *aggr* ensures that the sensor sending the data is valid and authenticates the sensor, else the sensor is considered to be invalid and it is deauthenticated.
6. *aggr* again encrypts the data along with the *MAC* by the symmetric key  $K_{ch,s}$  and transmits it to the sink

7. When all the aggregated data from *aggr* reaches the sink, it decrypts the data using symmetric key  $K_{ch,s}$ .
8. The sink checks if the aggregated data is valid without any change in its content by checking its *MAC*.
9. If the *MAC* is not valid, the *aggr* is prohibited from further transmissions.

## 4. Simulation Results

### 4.1. Simulation Setup

The performance of Energy-efficient secure data aggregation framework (**ESDAF**) protocol is evaluated through NS2 simulation [12]. A random network deployed in an area of 50 X 50 m is considered. We vary the number of Attackers as 1, 2,..5. Initially the nodes are placed randomly in the specified area. The base station is assumed to be situated 100 meters away from the above specified area. The initial energy of all the attackers assumed as 3.1 joules. The IEEE 802.15.4 MAC layer is used for a reliable and single hop communication among the devices, providing access to the physical channel for all types of transmissions and appropriate security mechanisms. The IEEE 802.15.4 specification supports two PHY options based on direct sequence spread spectrum (DSSS), which allows the use of low-cost digital IC realizations. The PHY adopts the same basic frame structure for low-duty-cycle low-power operation, except that the two PHYs adopt different frequency bands: low-band (868/915 MHz) and high band (2.4 GHz). The PHY layer uses a common frame structure, containing a 32-bit preamble, a frame length.

The simulated traffic is FTP with TCP source and sink. The number of sources is varied from 1 to 4. Table 1 summarizes the simulation parameters used

**Table 1:**Simulation Parameters

| No. of Attackers   | 1, 2, 3,...5  |
|--------------------|---------------|
| Area Size          | 100 X 100     |
| Mac                | IEEE 802.15.4 |
| Simulation Time    | 50 sec        |
| Transmission Range | 40m           |
| Routing Protocol   | CBQR          |
| Traffic Source     | FTP           |
| Packet Size        | 100           |
| Transmit Power     | 0.660 w       |
| Receiving Power    | 0.395 w       |
| Idle Power         | 0.335 w       |
| Initial Energy     | 3.1 J         |

## 4.2. Performance Metrics

The performance of ESDAF is compared with the ESDAF [10] protocol. The performance is evaluated mainly, according to the following metrics.

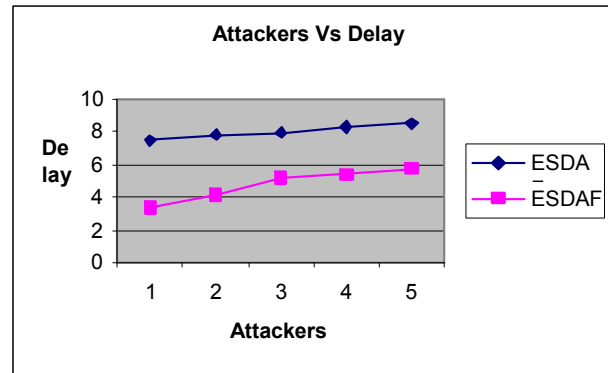
- **Average end-to-end Delay:** The end-to-end-delay is averaged over all surviving data packets from the sources to the destinations.
- **Average Packet Delivery Ratio:** It is the ratio of the number.of packets received successfully and the total number of packets transmitted.
- **Energy Consumption:** It is the average energy consumption of all nodes in sending, receiving and forward operations

The simulation results are presented in the next section.

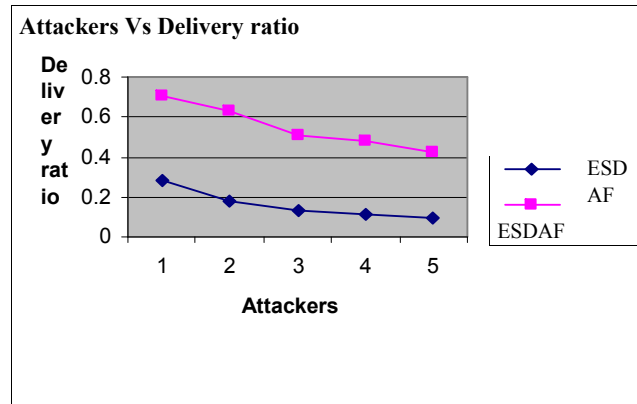
### Based on Attackers

In our initial experiment, we vary the number of Attackers as 1,2,3,4 and 5.

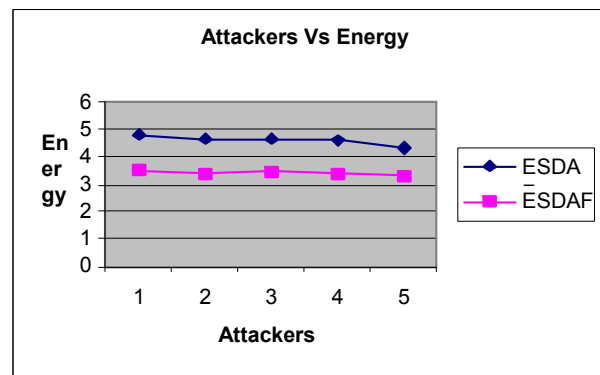
**Figure 2:** Attackers Vs Delay



**Figure 3:** Attackers Vs Delivery ratio



**Figure 4:** Attackers Vs Energy



**Figure 5:** Attackers Vs Overhead

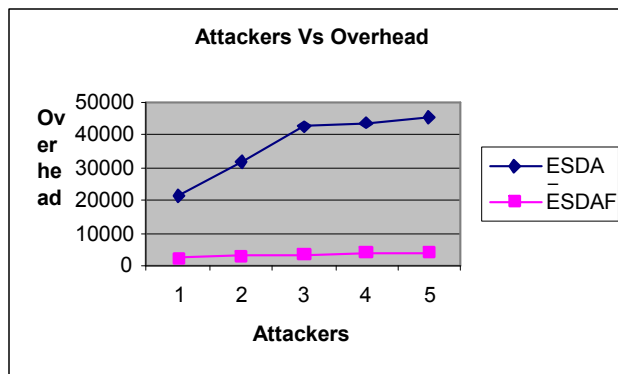


Figure 2 gives the average end-to-end delay for both the protocols when the number of nodes is increased. From the figure, it can be seen that the average end-to-end delay of the proposed ESDAF protocol is less when compared with ESDAF protocols.

Figure 3 presents the packet delivery ratio of both the protocols. Since reliability is achieved using the link stability, ESDAF achieves good delivery ratio, compared to ESDAF protocol.

Figure 4 shows the results of energy consumption for both the protocols. From the results, we can see that ESDAF protocol has less energy consumption than ESDAF protocol, since it has the energy efficient path.

Figure 5 shows the results of Overhead for both the protocols. From the results, we can see that ESDAF protocol has less Overhead than ESDAF protocol.

## 5. Conclusion

In this paper, we have developed a secure data aggregation protocol for wireless sensor networks which maintains energy efficiency. For data aggregation, the system is grouped such that each group is headed by an aggregator. This aggregator acts as a

link between the sensor nodes and the sink. During the transmission of the data, first encryption is performed by the sensor nodes when transferring data to the aggregator. The aggregator on reception of the data decrypts it using the key and reads it. The aggregator then determines the MAC value using hash function to check the validity of the source sensor. If the estimated MAC value is valid then the source is authenticated. Second encryption is performed by the aggregator when transferring data along with the MAC value to the sink. Hence integrity of the system is maintained. Due to the double encryption of the data during data aggregation, adversaries cannot affect the system. Hence the system remains secure even in the wireless environment. Simulation results show that our proposed protocol has reduced energy consumption while attaining good packet delivery ratio.

## References

- [1] Dorottya Vass, Attila Vidacs, "Distributed Data Aggregation with Geographical Routing in Wireless Sensor Networks", *Pervasive Services, IEEE International Conference* on July 2007.
- [2] Jukka Kohonen, "Data Gathering in Sensor Networks", Helsinki Institute for Information Technology, Finland. Nov 2004.
- [3] Gregory Hartl, Baochun Li, "Loss Inference in Wireless Sensor Networks Based on Data Aggregation", *IPSN 2004*.
- [4] Zhenzhen Ye, Alhussein A. Abouzeid and Jing Ai, "Optimal Policies for Distributed Data Aggregation in Wireless Sensor Networks", *Draft Infocom2007 Paper*.

- [5] Bhaskar Krishnamachari, Deborah Estrin and Stephen Wicker, "The Impact of Data Aggregation in Wireless Sensor Networks", *Proceedings of the 22nd International Conference on Distributed Computing Systems*, 2002.
- [6] Kai-Wei Fan, Sha Liu, and Prasun Sinha, "Structure-free Data Aggregation in Sensor Networks", *IEEE Transactions on Mobile Computing*, 2007.
- [7] Yingpeng Sang, Hong Shen, Yasushi Inoguchi, Yasuo Tan and Naixue Xiong, "Secure Data Aggregation in Wireless Sensor Networks: A Survey", *Seventh International Conference on Parallel and Distributed Computing, Applications and Technologies*, 2006.
- [8] Prakash G L, S H Manjula, K R Venugopal and L M Patnaik, "Secure Data Aggregation Using Clusters in Sensor Networks", *International Journal of Wireless Networks and Communications* Volume 1, Number 1 (2009), pp. 93–101.
- [9] Tamer AbuHmed and DaeHun Nyang, "A Dynamic Level-based Secure Data Aggregation in Wireless Sensor Network", Information Security Research Laboratory Graduate School of IT & Telecommunication InHa University.
- [10] Wenbo He, Xue Liu, Hoang Nguyen, Klara Nahrstedt and Tarek Abdelzaher, "PDA: Privacy-preserving Data Aggregation in Wireless Sensor Networks", 26th IEEE International Conference on Computer Communications. IEEE INFOCOM 2007.
- [11] Shih-I Huang and Shihpyng Shieh, "SEA: Secure Encrypted-Data Aggregation in Mobile Wireless Sensor Networks", *International Conference on Computational Intelligence and Security* 2007.
- [12] Bhoopathy, V. and Parvathi, R.M.S. "Energy Efficient Secure Data Aggregation Protocol for Wireless Sensor Networks", *European Journal of Scientific Research*, Vol. 50, Issue 1, pp.48-58, 2011.
- [13] Bhoopathy, V. and Parvathi, R.M.S. "Secure Authentication Technique for Data Aggregation in Wireless Sensor Networks" *Journal of Computer Science*, Vol. 8, Issue 2, pp 232-238, 2012.
- [14] Bhoopathy, V. and Parvathi, R.M.S. "Energy Constrained Secure Hierarchical Data Aggregation in Wireless Sensor Networks" *American Journal of Applied Sciences*, Vol. 9, Issue 6, pp. 858-864, 2012.
- [15] Bhoopathy, V. and Parvathi, R.M.S. "Securing Node Capture Attacks for Hierarchical Data Aggregation in Wireless Sensor Networks" *International Journal of Engineering Research and Applications*, Vol. 2, Issue 2, pp. 458-466, 2012.

# Security and Privacy Issues of Healthcare Application and Implication of Predictive Analytics in Big Data

A.S. Gousia Banu, Research Scholar, Department of Computer Science  
D. Saritha, Research Scholar of JNTU Kakinada, Department of Computer Science  
K Narasimhulu, Research Scholar, Department of Computer Science

**Abstract**— Every day extensive amount of data is being generated every day from Big data. This paper investigates the impact of Big data in Healthcare information. Big data analytics uses the Hadoop Tools which plays a key role in the real-time analysis to organize the unstructured and heterogeneous data which comes in immense volume. This paper describes the big data analytics in Healthcare, its profits and about the implication of predictive analytics in healthcare and also enlightens the data security and privacy problem in health care application.

**Keywords**— Big Data, Analytics, Healthcare, Hadoop, Predictive Analytics, privacy, security.

## 1. INTRODUCTION

The healthcare application historically has generated large amounts of data, driven by record keeping, compliance & regulatory requirements, and patient care. Healthcare is generating extensive amount of data which are both structured and also unstructured to keep the number of patient records. Driven by mandatory requirements and the potential to improve the quality of healthcare delivery meanwhile reducing the costs, these massive quantities of data (known as 'big data') hold the assurance of supporting a wide range of medical and healthcare services, including among others clinical decision support, disease surveillance, and population health management. This paper enlightens the impact of Healthcare along with its security and privacy problem in Big data analytics.

## 2. RELATED WORK

In [19] the authors have mentioned about the various types of challenges on Big data in healthcare application. They have also enlightened about the uses of Big data in Healthcare. In [20] profits of using the Big data in healthcare organization for the individuals' and for the patients, for the Government etc. Even some extensive assurance of big data in

healthcare like the quality of data, missing data are some challenges that need to be taken up.

In [21] the authors mentioned that by the massive set of data in the health care application provides the possibilities to do the predictive analysis and find the solutions to many problems. It is also mentioned that Predictive Analysis where we use various statistical method, machine learning techniques and, data mining approaches to process, analyze data and predict the outcome for the unknown bag of data. Healthcare domain is still in early stage to take up the new possibilities, which can be offered by big data solution and use it to do effective decision-making [22].

In [23] the author has designed a prediction algorithm which collects the data, aggregates, apply the case attribute capture the performance which can be implemented in healthcare organizations.

In [24] the authors have mentioned how the predictive analysis is being accomplished in Healthcare and what are the tools and techniques that are being used.

## 3. IMPACT OF BIG DATA IN HEALTHCARE INFORMATION

What exactly is big data? A report delivered to the U.S. Congress in August 2012 defines big data

as “large volumes of high velocity, complex, and variable data that require advanced techniques and technologies to enable the capture, storage, distribution, management and analysis of the information”.

By digitizing, bringing together and effectively using big data, healthcare organizations ranging from single-physician offices and multi-provider groups to large hospital networks and accountable care organizations stand to realize significant profits. There are some areas in which the big data and the Analytics can enhance the health care organization:

**3.1 Targeting the right people:** Identifying the people who are at risk and who could profit from additional treatments and screenings for contrary consequences. Data analytics can also optimize the health status. When dealing with large populations, it becomes even more important to know who can potentially profit from interference as a way to improve health and lower costs. During the management of the people, we can understand who can get the profit from the interference and also it is one of the way to enhance the patient’s health and cost can be reduced.

**3.2 Right intervention:** The ability to deliver the right intervention at the right time will improve as people begin to understand their own risks, monitor their health and share pertinent information with their care providers. The big data analytics has the capability to produce the right intrusion makes the people to understand the risks and monitor the health.

**3.3 Perceiving the spreading diseases in advance:** Predicting the viral diseases earlier before spreading based on the live analysis. This can be identified by analyzing the social logs of the patients suffering from a disease in a particular geo-location [8]. This helps the experts who are working in health care organizations to take some precautionary actions.

**3.4 Identifying the cost for treatments:** Pinpointing the treatments for the diseases which is very expensive and by using some analytics and with some effective analytics we can replace with lower costs.

#### 4. ARCHITECTURAL FRAMEWORK

We can easily understand the concept of frame work for Analytics in Big data project which is similar to conventional method of health informatics. The key difference between the health informatics and Big data Analytics is the way in which the processing take place. The analysis are performed

with the help of the business Intelligence tools are used in the health information processing. Hadoop/Map Reduce are the two open source platforms applicable on the cloud which stimulates the big data Analytics in the department of Healthcare. The interfaces of analytical tools which are used in Conventional Method and Big data are totally different from one another but the algorithms and other models are same. The tools in the conventional method are apparent and user friendly. The tools which are used in the Big data Analytics are very complex to understand and requires exhaustive programming knowledge and skills to apply those tools.

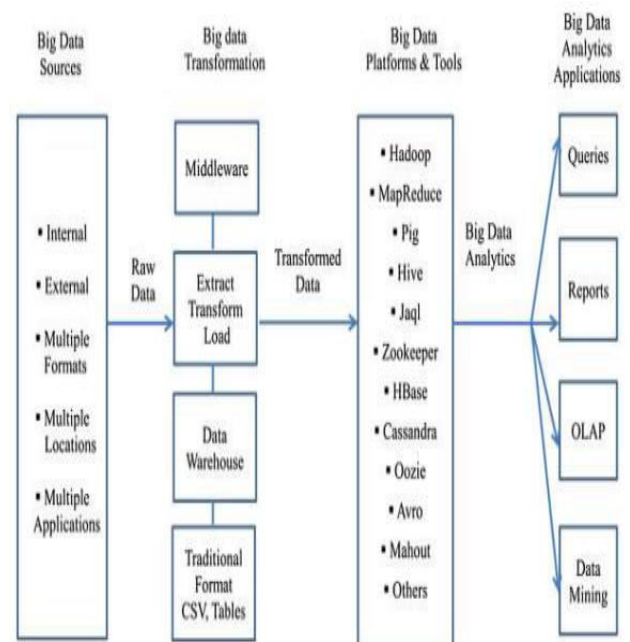


Fig: 1 Architecture of big data analytics

#### 5. PREDICTIVE ANALYTICS IN HEALTH CARE

Predictive analytics (PA) uses technology and statistical methods to search through massive amounts of information, analyzing it to predict outcomes for individual patients. That information can include data from past treatment outcomes as well as the latest medical research published in peer-reviewed journals and databases. [18] Predictive algorithms are used to diagnose the individual patient’s accurately by the physicians. When the patient’s come with the complaint of chest pain to the healthcare organization it is difficult to predict whether that patient should be hospitalized or not. If the details about the patient’s and their health status are entered properly in a system then the Predictive

algorithm will evaluate the patient's condition that he should be hospitalized or could be sent home.

Over the last few years, electronic health records (EHRs) have been widely implemented in the United States, and health care systems now have access to vast amounts of data. While they are beginning to apply "big data" techniques to predict individual outcomes like post-operative complications and diabetes risk, big data remains largely a buzzword, not a reality, in the normal delivery of health care [13] Predictive analytics encompasses a variety of statistical techniques from predictive modeling, machine learning, and data mining that analyze current and historical facts to make predictions about future or otherwise unknown events.[10][11] Predictive analytics and machine learning in healthcare are rapidly becoming some of the most-discussed, perhaps most-hyped topics in healthcare analytics. [9]

Some measurements to be taken for patient's care:

*To increase the way of diagnosing the diseases:* By predictive algorithms the physicians can make accurate diagnose of the patient's diseases. Health care and insurance costs. It is predicted that in 10 years' people in every world region will suffer more death and disability from non-communicable diseases than infectious diseases [12].

*Increasing the care for patients:* Health care application is taking much effort for anticipatory methods for treating the patients. For Instance, the patients can be classified as the one who will develop prolonged conditions and the other who will immediately respond to some sorts of treatments.

*Resource Optimization :* Based on the estimation of the admissions of the patients, the utilization of bed and the duration of the stay can also be analyzed for future database. .

## 6. BIG DATA SECURITY AND PRIVACY PROBLEM

*Information about Health and Exchanges and Electronic Health Records:* Kaiser Permanente (one of the healthcare providers in US) notified its 49,000 patients that their health information had been compromised due to theft of an unencrypted USB flash drive containing patient records [14]. This kind of healthcare offers Health Information exchanges (HIEs) and Electronic Health Records (EHRs) to share with the patients. This lead to the thieves to steal the data when amassing quantities of data are being shared between multiple providers within the

network. A study on patient privacy and data security showed that 94% of hospitals had somewhat one security breach in] the past two years. In most cases, the attacks were from an insider rather than external [15]

*Data Standardization and Structure problem:* The data which are applicable in healthcare application are in an unstructured format which will be in the form of graphs, charts, tables, images etc. Structured form of data will be heterogeneous. Understanding those kinds of data is a challenging issue.

*Imminent of Big data in Health care:* The usage of Big data and its acceptance throughout the world has provided different scopes to different observation in real-time activities. The following are some of the aspects that should be deliberated in future for the of Healthcare application.

Big Data has become an emerging force for the growth of IOT. Gartner estimates 26 billion IOT devices will be functional by 2020 and the amount of traffic generated by such devices will be large enough to place it in the category of big data [17].

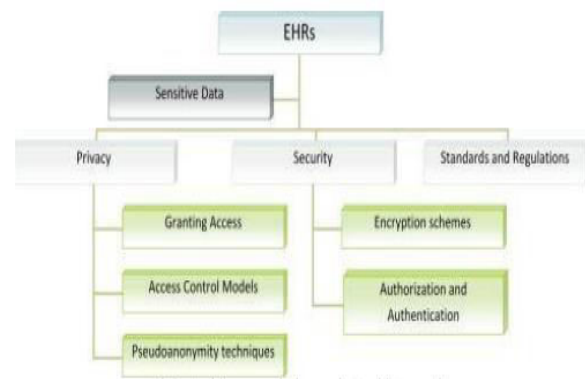


Fig2: Architecture of Electronic Health Record.

## 7. RESULT ANALYSIS

Predictive analytics has the potential to use the power of big data to improve health of patients and lower cost of health care

- Multiple models exist that will benefit from use of predictive analytics
- Several challenges need to be overcome to take advantage of full use of predictive analytics

### Predictive Analytics

Predictive analytics uses technology and statistical methods to search through massive

amounts of information, analyzing it to predict outcomes for individual patients.

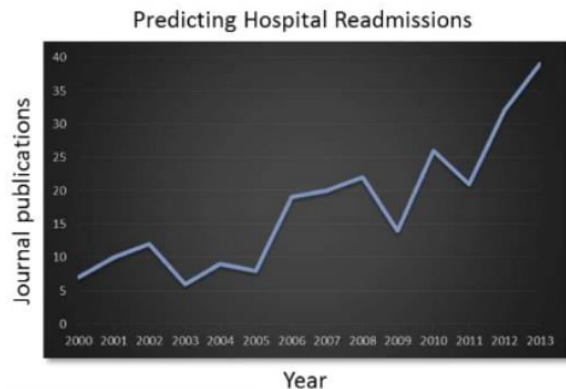


Fig:- Predicting Hospital Readmissions

Readmissions:-

- Use an algorithm to predict which patients are likely to be readmitted to the hospital.
- Tailor intervention to individual patient.
- Ensure patients actually receive the precise interventions intended for them.
- Monitor specific patients after discharge to find out if having problems before they decompensate .
- Ensure low ratio of patients flagged for intervention to patients who experience a readmission.

## 8. CONCLUSION

This paper enlightens the overview of big data in healthcare, the significant impact of big data in the healthcare organization. There are many opportunities and as well as many challenges are there in the healthcare application. Predictive analytics and algorithms are used in healthcare for saving the lives of the patients. It is also accepted the predictive analytics plays a key role in big data especially in healthcare organization. Even though many technologies are applicable some uncertainty of data or information from the electronic health record could lead to mortality or the patient's diseases cannot be diagnosed properly. Healthcare predictions are mostly profited by the emergency care to build tools with more potential variables as input which is applicable in Electronic Health Record.

## 9. REFERENCES

1. Raghupathi W: Data Mining in Health Care. Social insurance Informatics: Improving Efficiency and Productivity. Altered by: Kudyba S. , Taylor and Francis, 211-223,2010
2. Burghard C: Big Data and Analytics Key to Accountable Care Success. , IDC Health Insights ,2012
3. Dembosky An: "Information Prescription for Better Healthcare." Financial Times, December 12, 2012, p.
4. Feldman B, Martin EM, Skotnes T: "Enormous Data in Healthcare Hype and Hope." October 2012. Dr. Bonnie 360. 2012,<http://www.westinfo.eu/records/huge-information-in-healthcare.pdf>, Google Scholar
5. Fernandes L, O'Connor M, Weaver V: Big information, greater results. J AHIMA. , 38-42,2012
6. Burghard C: Big Data and Analytics Key to Accountable Care Success. , IDC Health Insights,2012
7. IHTT: Transforming Health Care through Big Data Strategies for utilizing huge information in the medicinal services application. 2013, IHTT: Transforming Health Care through Big Data Strategies for utilizing huge information in the medicinal services application. 2013,<http://ihealthtran.com/wordpress/2013/03/ih%20discharges-enormous-information-inquire-about-report-download-today>
8. Yanglin Ren, Monitoring patients by means of a protected and portable social insurance framework, IEEE Symposium on remote communication,2011
- 9.<https://www.healthcatalyst.com/prescient-investigation>
10. Nyce, Charles , Predictive Analytics White Paper(PDF), American Institute for Chartered Property Casualty Underwriters/Insurance Institute of America, p. 1,2007
11. Eckerson, Wayne , Extending the Value of Your Data Warehousing Investment, The Data Warehouse Institute.,2007

12. WHO, (2011). Worldwide Health and Aging, National Institute on Aging and National Institute of Health, U.S. Branch of Health and Human Services.

13. Ravi B. Parikh, Ziad Obermeyer, David Westfall Bates, <http://hbr.org/2016/04/making-prescient-examination-a-standard-piece-of-quiet-care>

14. E. McCann, "Kaiser reports second fall information rupture," Healthcare IT News, 2013.

15. P. Organization, "Third Annual Benchmark Study on Patient Privacy and Data Security," Pokémon Institute LLC, 2012.

16. P. Forests, B. Kayyali, D. Knott and S. V. Kuiken, "The 'enormous information' unrest in human services," McKinsey and Company, 2013.

17. P. Middleton, P. Kjeldsen and J. Tully, "Figure: The Internet of Things, Worldwide," Gartner, 2013

18. Linda A. Winters-Miner, PhD Posted on 6 October 2014 <https://www.elsevier.com/interface/seven-ways-prescient-examination-can-improve-human-services>

19. Priyanka Prof Nagarathna Kulennavar,. "A Survey On Big Data Analytics In Health Care International Journal of Computer Science and Information Technologies, Vol. 5 (4) , pp.5865-5896,2014

20. Jaslene Kaur Bains," Big Data Analytics in Healthcare-Its Benefits, Phases and Challenges", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 6, Issue 4,pp, 430-435 April 2016

21. Raol Priyanka Ajaysinh, Hinal Somani, "A Survey on Machine learning helped Big Data Analysis for Health Care Domain" IJEDR |Volume 4, Issue 4,pp-550-554,2016

22. Dharavath Ramesh, Pranshu Suraj, Lokendra Saini, "Enormous information Analytics in social insurance: A Survey Approach", IEEE, pp. 1 – 6, DOI: 10.1109/MicroCom.2016.7522520,2016

23. R. Sathiyavathi," A Survey: Big Data Analytics on Healthcare System", Contemporary Engineering Sciences, Vol. 8, no. 3,pp-121 – 125,2015

24. Ravi B. Parikh, Ziad Obermeyer, David Westfall Bates, <http://hbr.org/2016/04/making-prescient-examination-a-normal-piece-of-quiet-care>

#### First A. Author



A. S. Gousia Banu  
Research Scholar  
Department of Computer  
Science  
[gbanuzia@gmail.com](mailto:gbanuzia@gmail.com)

#### Second B. Author



D. Saritha  
Research Scholar of JNTU  
Kakinada  
Department of Computer  
Science  
[Ziauddinb17@gmail.com](mailto:Ziauddinb17@gmail.com)

#### Third C. Author



K. Narasimhulu  
Research Scholar,  
Department of Computer Science  
[narasimha.konduru@gmail.com](mailto:narasimha.konduru@gmail.com)

## ***Student Learning Experience by Data Mining & Social Media***

***S RASHEEDUDDIN***

***ASST.PROF CSE DEPT***

***MALLA REDDY COLLEGE OF ENGINEERING***

***rasheeduddin.phd@gmail.com***

***N.ANANTH RAM REDDY***

***ASST.PROF CSE DEPT***

***MALLA REDDY COLLEGE OF ENGINEERING***

***ananthramreddy\_cse@mrce.in***

### **Abstract:**

Many issues like depression, suicide, anger, anxiety are increasing among students. These issues are necessary to seek out and analyze, but students never discuss their issues with anyone. Today Social media is very popular medium where individuals share their feeling and opinion. Students also terribly active on social sites like Facebook and Twitter. Their unceremonious discussion on social media (e.g. Twitter, Facebook) illuminates light on their educational experiences—vote, sentiment, opinions, feelings, and concerns about the learning process. Data from such environments can supply valuable information which is helpful knowledge to understand student learning experiences. Analyzing such data can be challenging. The augmenting scale of data demands automatic data analysis techniques. This paper depicts a workflow to integrate both qualitative analysis and large-scale data mining techniques. This Paper emphasized on student's twitter posts to learn problems in student life as well as positive things occurred in their educational life. First conducted a qualitative analysis on sample tweets related to student's college life. Students face issues such as heavy work load of study, lack of social engagement, and sleep deprivation, employment issue, etc. In this paper "positive things" happen in student's life is also taken in to consideration. To classify tweets reflecting student's problem multilabel classification algorithms is implemented. Naïve Bayes and Linear Support Vector Machine Learning algorithms are used. The performance of these algorithms is compared in terms of accuracy, precision, recall and F1-Measure. Support Vector Machine learning algorithm have more accuracy than Naïve Bayes Algorithm.

**Keywords:** Education, computers and education, social networking, web text analysis, Twitter Multi-Label Classifier.

## INTRODUCTION

Data mining research provides several techniques, tools, and algorithms for enormous amounts of data to answer real-world issues. Social media plays powerful role in today's era. As social media is generally used for various purposes, vast amounts of user created data can be made available for data mining. Main objectives of the data mining procedure are to communally handle large-scale data, extract useful patterns, and gain required knowledge. Social media sites such as Twitter, Facebook, and YouTube provides stage to share happiness, struggle, sentiment, stress and acquire social support. On various social media sites, students discuss about their daily encounters in a comfortable and informal manner. They share their happiness and sorrows related to studies on social media in the form of judgmental comments, tweets, posts etc. Student's digital information provides large amount of implicit useful and reliable information for educational researchers to understand student's experiences outside the closed environment of classroom. This understanding can enhance education standard, and thus improve student employment, preservation, and accomplishment [1][2]. The Massive quantity of information on social sites gives prospective to recognize student's problem, however conjointly promotes some methodological complexities in use of social media data for educational reasons. The complexities such as assortment of Internet slangs, absolute data volumes and moment of students posting on the web. physical analysis cannot contract with the ever growing scale of data, while pure automatic algorithms cannot capture in-depth significance inside the data [3]. One important reason why social media will be relayed on is that the comments and posts are ad hoc emotions and feelings of students. This study will be terribly helpful and may prove revolutionary for an educational institute as crucial changes will be made in educational nature of the institute. The research goal of this learning are:-

- a) To make the enormous amount of data useful for educational purpose, as well as to combine both qualitative analysis and large-scale data mining techniques.
- b) To examine student's informal tweets on twitter in order to investigate the problems and issues faced by students in their life.

### I. LITRATURE SURVEY

The theoretical foundation for worth of informal data on the web is drawn from Goffman's theory of social performance. Although developed to elucidate face-to-face interactions, Goffman's theory of social performance is widely used to explain mediated interactions on the web today. One amongst the foremost basic aspects of this theory is that notion of front-stage and back-stage of people's social performances. Compared with the front-stage, the relaxing atmosphere of backstage typically encourages more spontaneous actions. Whether a social setting is front-stage or back-stage could be a relative matter. For students, compared with formal classroom settings, social media is a relatively informal and relaxing. When students post content on

social media sites, they usually post what they think and feel at that moment. In this sense, the data collected from on-line

conversations may be more authentic and unfiltered than responses to formal research prompts. These conversations act as a zeitgeist for students' experiences [1]. Many studies show that social media users may purposefully manage their on-line identity to "look better" than in real life. Human identity is complex and multifaceted. Humans acquire identity through social interaction and enact different roles, depending on context and social teams. For example, one may be a commanding, determined business executive at work, caring mother at home, and funny friend at a social gathering [2]. The facet of identity one enacts at a given point in time depends upon context and the particular social group (i.e. family, coworkers, friends) present in that context. Other studies show that there's an absence of awareness regarding managing online identity among college students, and that young people typically regard social media as their personal space to hang out with peers outside the sight of parents and teachers. Students' online conversations reveal aspects of their experiences that are not simply seen in formal classroom settings, thus are usually not documented in academic literature. Researchers from diverse fields have analyzed twitter content to generate specific knowledge for their respective subject domains. For example, Gaffney analyzes tweets with hashtag #iranElection using histograms, user networks, as well as frequencies of top keywords to quantify online activism. Similar studies are conducted in different fields including healthcare, marketing, and athletics, just to name a few. Analysis methods used in these studies usually include qualitative content analysis, linguistic analysis, network analysis, and few simplistic methods such as word clouds and histograms. This model was then applied and validated on a brand new data set. Therefore, emphasize not only the insights gained from one data set, but also the application of the classification algorithm to other data sets for detecting student issues. The human effort is thus augmented with large-scale data analysis. Researchers have analyzed twitter data and place efforts for introducing methods for classifying twitter data For example, Alec Go introduced a completely unique approach for automatically classifying the sentiment of Twitter messages [3]. These messages are classified as either positive or negative with relevance to a query term. This is useful for consumers who want to re- search the sentiment of products before purchase, or companies that want to monitor the public sentiment of their brands. There is no previous research on classifying sentiment of messages on microblogging services like Twitter. algorithms for classifying the sentiment of Twitter messages using distant supervisions. This is extremely useful because it allows feedback to be aggregated without manual intervention.

### II. PROPOSED SYSTEM

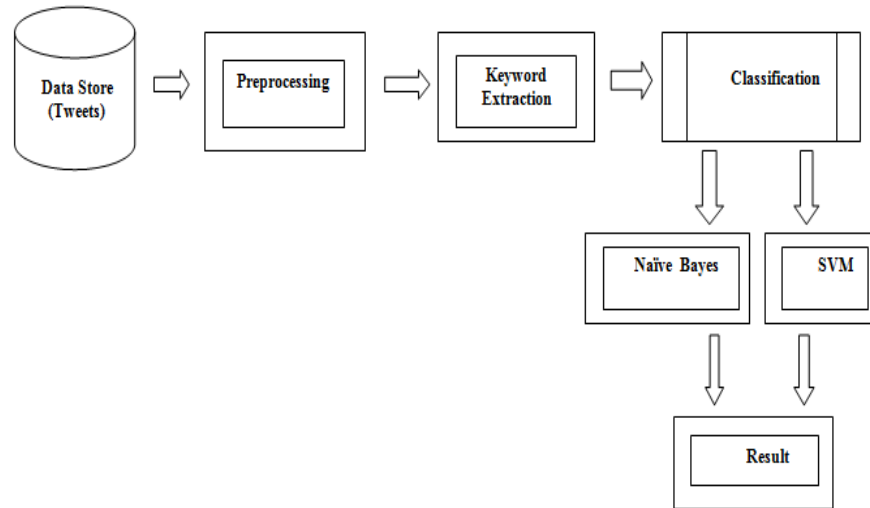
This approach uses totally different machine learning classifiers and has feature extractors. The proposed system works on student's tweets. These tweets are related to student's educational experiences. This system defines seven labels which are: heavy study load, lack of social engagement, negative emotions, sleep problem, diversity issues, positive things and struggle. The objective is to explore student's informal conversations on twitter in order to understand issues and problems of students encounter in their learning experiences. The tweets are loaded and

Proceeding of International Conference on Emerging Technologies in Computer Science, ISBN: 978-93-85100-13-0

processed by standard text mining procedure called Pre-processing. This system is used to understand student learning experiences. The existing system uses classification algorithms finds only negative emotions of students learning, whereas the proposed system will use to discover positive as well as negative emotions. Naïve bayes and support machine algorithm is used to discover the positive as well as negative emotion of student about their learning experiences. The comparison of the results of these algorithms is done using parameters accuracy, precision, recall and F1- Measure. The figure 1 shows the architecture of proposed system. First step is to collect the data for processing. This data is nothing but student's tweets which have positive and negative expression about their academic experiences. In next step collected data is explored and define the categories into which tweets can be differentiated. The tweets are preprocessed i.e. stemming, stop word cleaning. Stemming reduces inflected words to their stem, base or root form. In stop word cleaning, there are list of stop words which are removed by preprocessing from text documents. However, on tokenization stream of text is break into words, phrases or symbols. The model is

trained using multilabel classifier. The multi-label Support Vector Machine and Multi- label Naive Bayes classifier are implemented and compared. The procedure of the proposed system is as follow:

- 1) In the step one data collection is done from twitter.
- 2) Inductive Content analysis procedure is performed and categories are identified.
- 3) The preprocessing and tfidf is calculated in the step 2.
- 4) Naïve Bayes classifier, Linear SVM is applied on dataset in order to demonstrate its application in detecting student's issues is step 4.



**Figure.1. System Architecture**

### III. IMPLEMENTATION

This project implements using following processes Inductive content analysis is used to identify the relevant and irrelevant tweets. In this categories are identified in which tweets are going to classify. Naïve bayes and support vector machine algorithm are used to classify the tweets of student's informal discussions.

#### 4.1 DATA COLLECTION

It is challenging to collect social media data related to student's experiences because of the irregularity and diversity of the language used. Data searched on Twitter. The Twitter APIs can configure to accomplish this task. The search process was exploratory. Data searching is based on different Boolean combinations of possible keywords such as engineer, students, campus, class, homework, professor, and lab. Twitter API downloaded from twitter using tool My TwitterScraper. User should have twitter account and have to make an application for downloading tweets. These tweets can be download using hashtags. Twitter API can also be found location wise and particular person's account. MyTwitterScraper is java based tool which is used to download the Twitter API. These tweets are saved in excel format. But this data is unclean and contain many errors as well as bugs. Similarly all tweets are not relevant to this system. System required tweets which contain positive as well as negative emotion or experiences of student in their educational life. Using MyTweetScraper Tool 25000 tweets are downloaded among these tweets only 1700 tweets are useful for this system.

#### 4.1 Inductive content Analysis

Social media content like tweets contain a bulky amount of informal language, sarcasm, acronyms, and misspellings, meaning is often ambiguous and subject to human interpretation. Rost et. al argue that in large scale social media data analysis, faulty assumptions are likely to arise if automatic algorithms are used without taking a qualitative look at the data. According to study there is no appropriate unsupervised algorithms could reveal in-depth meanings in our data. For example, LDA (Latent Dirichlet Allocation) is a popular topic modeling algorithm that can detect general topics from very large scale data. LDA has only produced meaningless word groups from our data with a lot of overlapping words across different topics. There were no pre-defined categories of the data, so it is necessary to explore what students were saying in the tweets. Thus, first step is to perform an inductive content analysis on the dataset. Inductive content analysis is one popular qualitative research method for manually analyzing text content.

#### 4.2 Categories of Data

As a consequence proposed system first conducted an inductive content analysis on the n dataset. This paper classified the student tweets in to seven categories. Existing system have five prominent themes and proposed system consist of seven prominent categories:

##### i. Heavy Study Load

Analyses show that, classes, homework, exams, and labs control the student's life. Libraries, labs, and the college building are their most frequently visited places. Some

illustrative tweets are "Study for 30 hours..", "Doing homework since morning still incomplete.", and "OS project due Thursday.", and "homework never finish". Students express a very stressful experience. Not being able to manage the heavy study load. This finding echoes a previous study on students' life balance by which indicates student's desire a more balanced life than their academic environment allows.

##### ii. Lack of Social Engagement

The analyses show that students have to give up the time for social engagement in order to do homework, and to prepare for classes and study for exams. For example, "I feel like I'm hidden from the world—life of an student". Lack of social engagement is also tangled with the conventional nerdy and anti-social image of students. Some students embrace the anti-social image, while most others desire more social life as the examples above show.

##### iii. Negative Emotion

There are a bunch of negative emotions flowing in the tweets. This category specifically contains negative emotions such as hatred, anger, stress, sickness, depression, disappointment, and hopelessness. Students are mostly stressed with schoolwork. For example, "looking at my grades online makes me sick", "40 hours in the library in the past 3 days. I hate finals", and "I feel myself dying, #nervous". It is necessary for students to get help with how to manage stress and get emotional support.

##### iv. Sleep Issues

Analyses find that sleep issues are widely common among students. Students frequently suffer from lack of sleep and nightmares due to heavy study load and stress. For example, "Napping in the common room because I know I won't sleep for the next three days", "If I don't schedule in sleep time, it doesn't happen", and "I wake up from a nightmare where I didn't finish my physics lab on time". Chronic lack of sleep or low-quality sleep can result in many psychological and physical health issues; therefore this issue needs to be addressed.

##### v. Diversity Issues

The Analyses suggest students perceive a significant lack of females in engineering. For example, "eighty five kids leaving the classroom before mine. of those 85, four are girls. Engineers math class #Stereotypical", and "Keeping up with tradition: 2 girls in a class of 40". Male students in engineering are regarded as bad at talking with female students, because they usually do not have many female students around in their class. For example, "I'm sorry. We're not use to having girls around", "I pity the 1 girl in my lab with 25 guys. The issue here is not lack of diversity, but rather that students have difficulties embracing the diversity, because of many culture conflicts.

##### vi. Positive

A large number of tweets fall under this category. As student faces many issues in their educational life to recognizing as well as reducing those issues this system designed. Student also faces many positive things and emotions in their educational life to recognize such emotions are equally important. By recognizing such emotions which can have positive impact on student life is to make students life happy.

This category is consisting of emotions such as happy, Fine, good, easy, got, holiday, over, etc.

**vii. Struggle:** According to studied dataset conclusion occurs that many student tweets on their struggle in education. So new

class introduced in proposed work. Student Tweets on they face struggle in finding job. They are struggling for passing exam.

### 4.3 Text Preprocessing

Twitter users use some special symbols to convey certain meaning. For example, # is used to indicate a hashtag, @ is used to indicate a user account, and RT is used to indicate a re-tweet. Twitter users sometimes repeat letters in words sothat to emphasize the words, for example, “huuungryyy”, “sooo muuchh”, and “Monnndayyy”. Besides, common stopwords such as “a, an, and, of, he, she, it”, non-letter symbols, and punctuation also bring noise to the text. So we pre-processed the texts before training the classifier:

### 4.4 Naïve Bayes Multi-Label Classification Algorithm

One popular way to implement multi-label classifier is to transform the multi-label classification problem into multiple single-label classification problems. One simple transformation method is called one-versus-all or binary relevance. The basic concept is to assume independence among categories, and train a binary classifier for each category. All kinds of binary classifier can be transformed to multi-label classifier using the one-versus- all heuristic. The following are the basic procedures of the multi- label Naive Bayes classifier. In this implementation, if for a certain document, there is no category with a positive probability larger than T, then assign the one category with the largest probability to this document. In addition, “others” is an exclusive category. A tweet is only assigned to “others” when “others” is the only category with probability larger than T.

### 4.5 Support Vector Machine Algorithm

Transforming a multi-label classification problem into a set of independent binary classification problems via the “one-vs-all” scheme is a conceptually simple and computationally efficient solution for multi-label classification [18]. In this work, multi- label learning conducted under such a

mechanism by using standard support vector machines (SVM) for the binary classification problems associated with each class.

## EXPERIMENTAL RESULT

This system uses approximately 1700 tweet of student downloaded from twitter which contain student emotions. Firstly For a document  $d_i$  in the testing set, there are  $K$  words  $W_{d_i} = \{\omega_1, \omega_2, \dots, \omega_N\}$ , and  $W_{d_i}$  is a subset of  $W$ . The purpose is to classify this document into category  $c$  or  $c'$ . Independence among each word assumed in this document, and any word  $w_{ik}$  conditioned on  $c$  or follows multinomial distribution.

Therefore, according to Bayes' cross validation result of naïve bayes.

| FOLD | TEST<br>SUBJECT | TRAIN<br>SUBJECT | ACCURACY |
|------|-----------------|------------------|----------|
| 0    | 270             | 1373             | 83.18    |
| 1    | 273             | 1370             | 83.34    |
| 2    | 273             | 1370             | 83.18    |
| 3    | 273             | 1370             | 83.26    |
| 4    | 273             | 1370             | 83.12    |

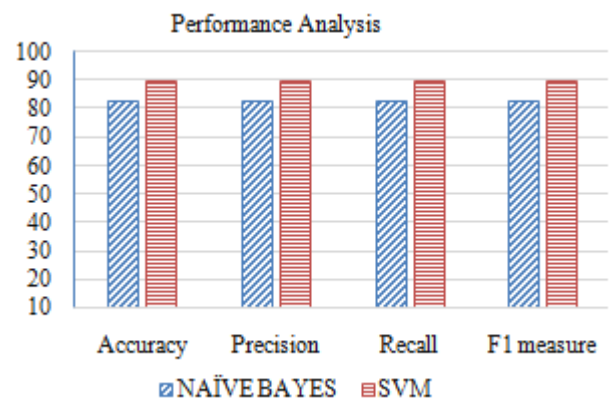
Table I shows the cross validation result of Naïve Bayes Algorithm. Here 1700 dataset is divided in to 5 fold. In first iteration 270 tweets are considered as test set and remaining 1373 are training set accuracy of this iteration is 83.18. In second iteration 273 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 83.34. In third iteration 273 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 83.18. In forth iteration 273 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 83.26. In fifth iteration 272 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 83.12.

In third iteration 273 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 90.48. In forth iteration 273 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 90.16. In fifth iteration 273 tweets are considered as test set and remaining 1370 are training set accuracy of this iteration is 90.35.

**Table.2. Cross validation result of svm.**

| FOLD | TEST<br>SUBJECT | TRAIN<br>SUBJECT | ACCURAC<br>Y |
|------|-----------------|------------------|--------------|
| 0    | 273             | 1370             | 90.35        |
| 1    | 273             | 1370             | 90.14        |
| 2    | 273             | 1370             | 90.48        |
| 3    | 273             | 1370             | 90.16        |
| 4    | 273             | 1370             | 90.35        |

**Table.3. Comparison of parameters using naïve bayes and svm**



**Figure.2.Performance Evaluation of Nb and Svm**

Table III shows the final result by taking average of above cross validation result. This table shows the accuracy, precision, recall and F1-measure of Naïve Bayes is 83.12, 83.35, 83.04 and 83.21 respectively and the accuracy, precision, recall and F1-measure of SVM is 90.35 89.94, 90.00 and 90.30 respectively Figure II shows that Support Vector Machine has better performance than Naïve Bayes algorithm. As accuracy, precision, recall and F1- Measure values of SVM are better as compared to Naïve Bayes.

#### IV. CONCLUSION

This project is classifies the student tweets to understand their problems as well as positive things in their educational life. Database is collection of tweets these tweets which depicts the information of student experiences in their educational life. The Tweets are stored in database. Preprocessing is done on tweets. Then classification techniques (Naïve Bayes and support vector Machine) are applied on data. In classification student tweets classified in to seven categories. That is negative, sleeping problems, positive, diversity issues, Lack of social awareness, heavy study load, Struggle. Cross validation technique is used to evaluate performance of a system. Accuracy of naïve bayes algorithm in five iteration is 83.18, 83.34, 83.18, 83.26, and 83.12 respectively. Accuracy of SVM in five iteration is 90.35, 90.14, 90.48, 90.16, and 90.35 respectively. Finally average of 5 iteration result is calculated to evaluate performance of a system. Accuracy, Precision, Recall And F1-Measure of Naïve Bayes is 83.12, 83.25, 83.04 and 83.21 respectively and the Accuracy, Precision, Recall And F1-Measure of SVM is 90.35, 89.94, 90.00 and 90.30 respectively Comparison of algorithm is done on the values of Accuracy, Precision, Recall, and F1 Measure and came to conclusion that Support Machine Algorithm gives more accurate prediction than Naïve Bayes Algorithm. Accuracy of Naïve Bayes is 83.12 % while Support Vector Machine is 90.35 %. Precision of Naïve Bayes is 83.25% while Support Vector Machine is 89.94 %. Recall of Naïve Bayes is 83.04 % while Support Vector Machine is 90.00 %. F1- Measure of Naïve Bayes is 83.21 % while Support Vector Machine is 90.30 %.

#### V. REFERENCES

- [1]. Xin Chen, Mihaela Vorvoreanu, and Krishna Madhavan, "Mining social media data for understanding students' learning experiences", IEEE Transaction, 2014.
- [2]. Mariam Adedoyin-Olowe,, Mohamed Medhat Gaber and Frederic Stahl," A Survey of Data Mining Techniques
- [3]. E. Goffman, The Presentation of Self in Everyday Life. Lightning Source Inc., 1959.
- [4]. J.M. DiMicco and D.R. Millen, "Identity Management: Multiple Presentations of Self in Facebook," Proc. the Int'l ACM Conf. Supporting Group Work, pp. 383-386, 2007.
- [5]. M. Vorvoreanu and Q. Clark, "Managing Identity Across Social Networks," Proc. Poster Session at the ACM Conf. Computer Supported Cooperative Work, 2010.
- [6]. B. Sriram, D. Fuhry, E. Demir, H. Ferhatosmanoglu, and M. Demirbas, "Short Text Classification in Twitter to Improve Information Filtering," Proc. 33rd Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 841- 842, 2010.
- [7]. Bodong Chen, Xin Chen, Wanli Xing, "Twitter Archeology of Learning Analytics and Knowledge Conferences" Proceedings of fifth international conference on lerning analytic and knowledge, pp. 340-349,2015.
- [8]. A. Go, R. Bhayani, and L. Huang, "Twitter Sentiment Classification Using Distant Supervision," CS224N Project Report, Stanford pp. 1-12, 2009.
- [9]. Hsin-Ying Wu, Kuan-Liang Liu and Charles Trappey "The Theory On User Feedback Analysis."
- [10]. Suleyman Cetintas, Luo Si, Hans Peter Aagard, Kyle Bowen, and Mariheida Cordova-Sanchez," Microblogging in a Classroom: Classifying Students' Relevant and Irrelevant Questions in a Microblogging-Supported Classroom," IEEE Transactions on Learning Technologies, Vol. 4, No. 4, October- December 2011
- [11]. D. Gaffney, "#iranElection: Quantifying Online Activism," Proc. Extending the Frontier of Society On-Line (WebSci10), 2010.
- [12].W. Zhao, J. Jiang, J. Weng, J. He, E.P. Lim, H. Yan, and X. Li,"Comparing Twitter and Traditional Media Using Topic Models," Proc. 33rd European Conf. Advances in Information Retrieval, pp. 338- 349, 2011.
- [15]. Andreas Hotho," A Brief Survey of Text Mining" KDE Group University of Kassel May 13, 2005.

## IMPROVING THE SECURITY ISSUES IN WIRELESS SENSOR USING HETEROGENEOUS ALGORITHM

\*Dr. G. Silambarasan, \*\* Dr. V. Chandrasekar,

\*Assistant Professor, Dept. of Computer Science and Engineering,  
The Kavery College of Engineering, Salem, Tamilnadu, India,

\*\*Associate Professor, Dept. of Computer Science and Engineering,  
Malla Reddy College of Engineering and Technology, Secunderabad, Telangana State, India,

\*\*[drchandru86@gmail.com](mailto:drchandru86@gmail.com), \* [gssilambarasan@gmail.com](mailto:gssilambarasan@gmail.com)

### ABSTRACT

Heterogeneous sensor is a network it will be brought and process for individual applications. So many applications are mixed in these nodes for the operation of distributed application.

Heterogeneous sensor networks are two or more different types of node and varying levels of battery energy. Those nodes are communication capabilities and organized in to cluster to allow the scalability of media access control and routing. After cluster formation any one node act as a cluster head, multihop communication backbone to carry aggregated traffic, and remaining node cluster member which transmit sensing data to head directly. There are many networks layer security protocols are proposed for homogeneous sensor networks and mobile networks are our proposed homogeneous sensor networks. Compared to other networks, they coding are unique properties.

Heterogeneous mobile environment main challenges is to allow the user to access that scheme is service any time, any place and anywhere in any domain.

**Keywords:**—WSN, heterogeneous network, security

### INTRODUCTION

The framework of the wireless networks is to provide the services to the user irrespective of the location or the need of a physical medium. The urge of this flexibility is opted by the users on the run, to access the resources of the concern with provided authenticated identities [3]. The need of fixation of a node in a stated place, communication links via physical cables and limitations to the capacity of the medium used, promoted the concept of deploying wireless networks. The wireless technology eliminates the difficulties of a wired network by allowing the user to access the resources with no limitations. Wireless sensor

networks have a number of thousand nodes or more distributed in remote locations and all nodes are capable of requesting service simultaneously. The important factor is that not all the users are legitimate requestor of a service. There are outsiders who perform activities which disturb the security and integrity of a network. Their goal is to bother the functionality of an intended user and the services provided by the network either by blocking [5], distracting or by flooding the medium. The attacker acts in between two users to prevent them from communicating [7]. The attackers learnt the ways to hide from the detection algorithms by acting as a legitimate user (spoofing)[8] or making the network administrator to believe that there is no attack in the network. There may be attackers internal to a network, that is, a legitimate user could also attack the network functions for his particular reasons. The wireless networks are prone to a much higher rate of attacks than the wired networks. Any attacker who gains access to a network for altering the default activities is a serious threat to the data of high confidentiality. The detection algorithms have not matched to the speed of detecting the attack in earlier stages. This paper studies the jamming attacks of the attacker by various means and analyses the

effects. The motivated study is to make sure that the jammers could be used for constructive mechanisms of conserving the security of a highly important network which cannot be compromised at any cost [9].

## **MAC ROUTING**

MAC Media Access Control Address it is a hardware address that is uniquely identifies .In IEEE 802 networks. In OSI reference layer MAC is the Data link layer .Data link layer consist of two layers one is the Logical Link Control (LLC) layer. On networks that do not conform to the IEEE 802 standards but do conform to the OSI Reference Model, the node address is called the Data Link Control (DLC) address. Routing is the selection of process for traffic in network or across multiple networks. Routing is performed for many types of networks, including circuit – switched networks, such as the PSTN-Public switched telephone network .computer networks, such as the internet as well as in networks used in public and private transportation, such as the system of streets, roads, and highway in national infrastructure.

```

Microsoft Windows [Version 6.3.9600]
(c) 2013 Microsoft Corporation. All rights reserved.

C:\Users\vhudson>route print

Interface List
7...08 11 96 79 59 85 .....Microsoft Wi-Fi Direct Virtual Adapter
4...08 11 96 79 59 84 .....Intel(R) Centrino(R) Advanced-N 6205
3...f0 de f1 95 02 a0 .....Intel(R) 82579LM Gigabit Network Connection
19...08 00 27 00 74 6e .....VirtualBox Host-Only Ethernet Adapter
1.....00 00 00 00 00 00 .....Software Loopback Interface 1
8...00 00 00 00 00 00 e0 .....Microsoft ISAAP Adapter
16...00 00 00 00 00 00 e0 .....Microsoft ISAAP Adapter #2

IPv4 Route Table

Active Routes:
Network Destination        Netmask          Gateway           Interface        Metric
0.0.0.0                    0.0.0.0          10.255.72.1       10.255.77.167    10
10.255.72.0                255.255.248.0    On-link          10.255.77.167    266
10.255.77.167             255.255.255.255  On-link          10.255.77.167    266
10.255.79.255             255.255.255.255  On-link          10.255.77.167    266
127.0.0.0                 255.0.0.0        On-link          127.0.0.1        306
127.0.0.1                 255.255.255.255  On-link          127.0.0.1        306
127.255.255.255          255.255.255.255  On-link          127.0.0.1        306
172.18.72.0              255.255.248.0    10.255.72.10     10.255.77.167    11
192.168.56.0             255.255.255.0    On-link          192.168.56.1     276
192.168.56.1             255.255.255.255  On-link          192.168.56.1     276
192.168.56.255           255.255.255.255  On-link          192.168.56.1     276
  
```

**Fig: MAC Routing Table**

## HETEROGENEOUS ALGORITHM

The Research problems addressed here is about the data transmission from different clusters, each cluster is a system with unique features connected in heterogeneous network and how to configure these cluster with different configuration like hardware and operating system for sharing of information (data). Major scope is to provide the security (authentication).

Key research questions addressed are:

1. What are the various problems that can be possible in heterogeneous network?
2. What will be advantages with clusters in heterogeneous environment for data sharing?
3. How clusters are framed and how these clusters will be connected or mapped.
4. How to provide security Issues to protect the data and network.

As the network is Heterogeneous network is a network connecting computers and other devices with different operating systems and protocols. For Example Local Area Networks (LANs) that connect Microsoft windows and Linux based personal computers with Apple Macintosh Computers are Heterogeneous. The word Heterogeneous network is also used in wireless networks using different access technologies [9]. For Example a wireless network which provides a service through a wireless LAN and is able to maintain the service when switching to a cellular network is called a wireless heterogeneous network. A Wide Area Network (WAN) can use macro cells[11], pico cells and femto cells in order to offer wireless coverage in an environment with variety of wireless coverage zones, It environment range from office building, homes and undergrounds areas. Mobile experts define a Het Net as a network with complex interoperation between macro cell, small cell and in some cases Wi-Fi network elements used together to provide a mosaic of coverage, handoff capability between network elements.

Small Cell Forum defines the Het Net as multi-x environment-multi-technology, multi- domain, and multi-spectrum, multi-operator and multi-vendor.

It must be able to automate the reconfiguration of its operation to deliver assured service quality across the entire network and flexible enough to accommodate changing user needs, business goals and subscriber behaviors.

From that design, that can be encompassing conventional macro radio access networks (RAN) functions, RAN transport capability, small cells and Wi-Fi functionality, that are increasingly being virtualized and delivered in an operational environment where span of control includes data center resources associated with compute, networking and storage. In this framework [12], self –optimizing network (SON) functionality is essential to enable order-of- magnitude network densification with small cells. Self – configuration or ‘plug and play’ reduces time and cost of deployment, while self – optimization ensures the network auto tunes itself for maximum efficiency as condition change. Traffic demand, user movements and service mix will all evolve over time, and the network needs to adapt to keep pace [10]. These enhanced SON capabilities will therefore need to take into account the evolving user needs, business goals and subscriber behaviours. Importantly, functions associated with Het Net operations and

management take earlier SON capability that may have only been targeted at a single domain and technology, and expand it to deliver automated service quality management across the entire Het Net. A Heterogeneous wireless network (HWN) is a special case of a Het Net. Where as a Het Net may consist of a network of computers or devices with different capabilities in term of operating systems, hardware protocols, etc. A HWN is a wireless network which Consist of devices using different underlying radio access technology (RAT).

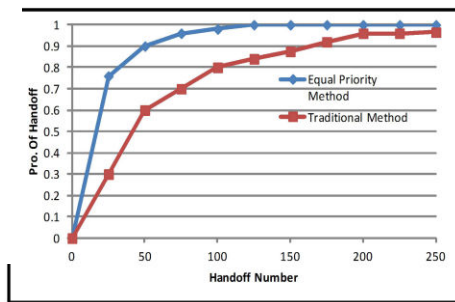
There are several problems still to be solved in heterogeneous wireless networks such as:

They have much solution in this wireless networks for example given below

- Determining the theoretical capacity of HWNs
- Interoperability of technology
- Handover, Mobility
- Quality of service / Quality of Experience
- Interference between RATs

There are several benefits to a HWN as opposed to a traditional homogeneous wireless network including increased reliability, improved spectrum efficiency, and increased coverage. Reliability is improved because when

one particular RAT [13] within the HWN fails, it may still be possible to maintain a connection by failing back to another RAT. Spectrum efficiency is improved by making use of RATs which



**Fig: Handoff Methodology**

may have few users through the use of load balancing across RATs and coverage may be improved because different RATs may fill holes in coverage that any one of the single network alone would not be able to fill. From a semantically point of view, it is very important to note that the Heterogeneous network terminology may have different connotations in wireless telecommunications [11]. For instance, it may refer to the paradigm of seamless and ubiquitous interoperability between various multi-coverage protocols (Het Net). Otherwise, it may refer to the non-uniform spatial distribution of users or wireless nodes (Spatial In homogeneity). Therefore, using the term “heterogeneous network”

without putting it into perspective may result in a source of confusion in scientific literature and during the peer-review cycle. In fact, the confusion may further be aggravated, especially in light of the fact that [12]. “Het Net” paradigm may also be researched from a “geometrical” angle.

## CONCLUSION

Sender sends the data to receiver side with proper protected data or encrypted data. Our proposed work objectives is the receiver side how to encrypted data and protected data improving its security

Using the heterogeneous algorithm. To improving the both sender and receiver side security.

## REFERENCES

- [1] A. Cuevas, P. Serrano, J. I. Moreno, C. J. Bernardos, J. Jähnert, R. L. Aguiar, V. Marques, Usability and Evaluation of a Deployed 4G Network Prototype, Journal of Communications and Networks, Vol. 7 (2), 2008.
- [2] Teo, Joseph Chee Ming; Tan, Chik How; Ng, Jim Mee, Denial-of-service attack resilience dynamic group key agreement for heterogeneous networks, Telecommun. Syst. 35, No. 3-4, 141-160 (2007).

- [3] L. J. LaPadula. State of the Art in Anomaly Detection and Reaction Technical Report MP 99B0000020, Mitre, July 1999.
- [4] G.L.F. Santos, Z. Abdelouahab, R.A. Dias, C.F.L. Lima, E. Nascimento , E.M. Cochra. An Automated Response Approach for Intrusion Detection Security Enhancement, Software Engineering and Applications, 2003.
- [5] M. Petkac and L. Badger, Security agility in response to intrusion detection in 16th Annual Conference on Computer Security Applications (ACSAC '00), 2000.
- [6] C. Feltus, D. Khadraoui, B. de Rémont and A.Rifaut, Business Gouvernance based Policy regulation for Security Incident Response. IEEE Global Infrastructure Symposium, 6 July 2007.
- [7] Gateau, D. Khadraoui, C. Feltus, Multi-Agents System Service based Platform in Telecommunication Security Incident Reaction, IEEE Global Information Infrastructure Symposium, 2009.
- [8] N. Damianou, N. Dulay, E. Lupu, M. Sloman , The Ponder Policy Specification Language, Workshop on Policies for Distributed Systems and Networks (Policy2001), HP Labs Bristol, 29 31. Springer-Verlag.
- [9] Bertino, E., Mileo, A., and Proveti, A. 2005. PDL with Preferences. IEEE international Workshop on Policies For Distributed Systems and Networks, Policy 2005 – Vol. 00, IEEE Computer Society, Washington, DC, 213-222.
- [10] Basile, C.; Lioy, A.; Perez, G. Martinez; C., F. J. Garcia; Skarmeta, A. F. Gomez, POSITIF: A Policy-Based Security Management System, Policies for Distributed Systems and Networks, 2007. POLICY'07, pp. 280 – 280.
- [11] Kim, C.S., Kim, J.I., Han, W.Y., Kwon, O.C., “Development of Open Telematics Service Based on Gateway and Framework”, Proc. of the ICACT, 2006, pp.1349-1352.
- [12] Han, W.Y., Kwon, O.C., Park, J.H., Kang, J.H., “A Gateway and Framework for Interoperable Telematics Systems Independent on Mobile Networks”, ETRI Journal, Vol.27, No.1, 2005, pp.106-109.
- [13] D.W. Lee, H.K. Kang, D.O. Kim, K.J. Han, “Development of a Telematics Service Framework for open Services in the Heterogeneous Network Environment”, Proc. of the International Congress on Anti Cancer Treatment ICACT 2009.
- [14] Shin-Hun Kang, Jae-Hyun Kim, “QoS-Aware Path Selection for Multi- Homed Mobile Terminals in Heterogeneous Wireless Networks”, Proc. Of the IEEE CCNC 2010, Jan, 2010.

[15] Min Liu, Zhongcheng Li, Xiaobing Guo, Eryk Dutkiewicz, “Performance Analysis and Optimization of Handoff Algorithms in Heterogeneous Wireless Networks”, IEEE Transactions on Mobile Computing, 2007.

# Improving Energy Efficient in Wireless Sensor Networks Using Path Algorithm

**Dr.A.Mummoorthy<sup>1</sup>**

Associate Professor

Department of Information Technology,

Malla Reddy College of Engineering & Technology.

e-mail : amummoorthy@gmail.com

**Sudha Pavani. K<sup>2</sup>**

Assistant Professor,

Department of Computer Science of Engineering,

Malla Reddy College of Engineering.

e-mail: sudhapavanil@gmail.com

**Abstract** - Node are using some powered by using of batteries in wireless sensor networks, with limited amount of energy, there are two severe problems in WSN increasing the life span of the network and reducing the usage of energy. We introduce the minimum spanning tree reduce the total energy consumption of WSN. Heavy load of sending the data packets to the node easily reduce the energy. Our proposal work aimed on presenting an Energy Conserved Fast and Secure Data Aggregation Scheme for WSN in time and security logic occurrence data collection application. Invention is finished on Energy Efficient Utilization Path Algorithm (EEUPA), to extend the lifespan by processing the collecting series with path mediators depending on gene characteristics sequencing of node energy drain rate, energy consumption rate, and message overhead together with extended network life span. A mathematical programming technique is designed to improve the lifespan of the network. Simulation experiments carried out among different relating conditions of wireless sensor network by different path algorithms to analyze the efficiency and effectiveness of planned Efficient Energy Utilization Path Algorithm in wireless sensor network (EEUPA).

**Keywords** - WSN, Energy, EEUPA, Data Collection

## 1. INTRODUCTION

A wireless sensor network (WSN) physically consists of large amount of miniature, multifunctional and resource controlled sensors which are self-organized as an informal network to examine the physical world. Sensor networks are often used in applications where it is complicated or impossible to gather wired networks. WSNs in various areas such as examination, disaster liberation, intellectual carrying, surveillance, environmental managing, healthcare, goal tracking, and more. To collect the information and data in unkind or defensive atmosphere, Wireless Sensor Networks are more useful. In a Wireless Sensor Networks, data collected by sensor nodes are preferred to be circulated to destinations (base stations).

## 2. LITERATURE REVIEW

The field of wireless sensor networks (WSN) is now developed in the research community area because of its applications in different fields for instance defense security, civilian applications and medical research. These limitations eliminate the utilization of traditional path protocols planned for other ad hoc wireless networks. Secondly, the base station cannot verify data reliability and accuracy via fixing message digests or signatures to every sensing model. To come across the above two disadvantages, the base station [1] can regain all sensing data yet these data has been combined. The multi level data aggregation method are presented [6] to develop the data collecting path for a mobile destination as Infrastructure based Data Gathering Protocol (IDGP) and a Distributed Data Gathering Protocol (DDGP). A k-hop relay mechanism is established to restrict the number of hops for path data to a mobile sink.

The key point of EEGA scheme [2] is that accurate data aggregation is attained without discharging secret sensor readings and without initiating important overhead on the battery-limited sensors. To come across the delay planned an Efficient Data Collection Aware of Spatio-Temporal Correlation (EAST) that utilizes shortest routes used for forwarding the collected data toward the sink node [3] and fully expand both spatial and temporal correlations to execute near real-time data collection in WSN. Developed [4] easy, least-time, energy-efficient path protocol by means of one-level data aggregation which guarantees improved life span for the network. The energy-efficient data development issue [5] through person packet delay limitations to an energy-efficient service curve construction issue is developed and solves the difficulty by developing a local optimality theorem.

To design a method to mitigate the uneven energy dissipation problem by controlling the mobility of agents, which is achieved by an energy prediction strategy to find their positions using [7] energy balancing cluster routing based on a mobile agent (EBMA) for WSNs. To obtain this while maintaining a good trade-off between the communication overhead of the scheme, the storage space requirements on the nodes, and the ratio between the number of visited nodes  $x$  and the representativeness of the gathered data using [8]. density-based proactive data dissemination Protocol (DEEP), which combines a probabilistic flooding with a probabilistic storing scheme. The QoS of an energy-efficient cluster-based routing protocol called Energy-Aware routing Protocol (EAP) [9] in terms of lifetime, delay, loss percentage, and throughput, and proposes some modifications on it to enhance its performance.

QoS of an energy-efficient cluster-based path protocol called Energy-Aware routing Protocol (EAP) in terms of lifespan, delay, loss percentage, and throughput, and presents some alterations on it to improve its results [10]. DHAC [11] can include both quantitative and qualitative information types in clustering. It prevents the destination to collect a representative investigation of the network's sensed data. Presented efficient and balanced cluster-based data aggregation algorithm (EEBCDA) splits the network into rectangular grids [12] with uneven size and creates cluster heads rotate between the nodes in every grid correspondingly, the grid whose cluster head utilizes more energy by means of offering unbalanced energy dissipation.

To produce system-level behaviors that show life-long adaptively to changes and perturbations in an external environment using [13]. bee-inspired BeeSensor protocol that is energy-aware, scalable and efficient. En Energy Efficient and QoS aware multipath routing protocol (EQSR) [14] that maximizes the network lifetime through balancing energy consumption across multiple nodes, uses the concept of service differentiation to allow delay sensitive traffic to reach the sink node within an acceptable delay, reduces the end to end delay through spreading out the traffic across multiple paths, and increases the throughput through introducing data redundancy. In this approach [15-19] shares intermediate results among queries to reduce the number of messages. When the sink receives multiple queries, it should be propagated these

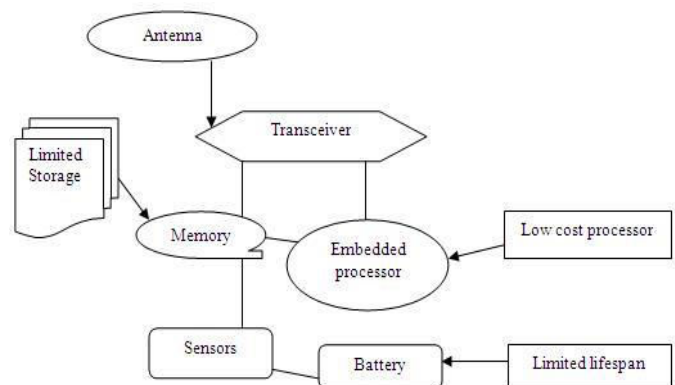
queries to a wireless sensor network via existing routing protocols. The sink could obtain the corresponding topology of queries and views each query as a query tree. With a set of query trees collected at the sink. The above planned research gap, encouraged us to plan a Fast and Secure Energy Efficient Data Aggregation Scheme for Wireless Sensor Network matched to applications.

### 3. METHODOLOGY

WSN is a wireless sensor network relating to the distributed device by means of sensors to observe the environmental conditions at different positions. A WSN is developed by planned sensor nodes in an application area. The Sensor Node is a very important factor of WSN, is prepared with Computation, sensing and wireless Communication unit. Wireless Sensor Networks are created properly for security, environmental monitoring, computerization, habitat groping and creative industries etc. A block diagram of wireless sensor network system architecture is shown in Fig.1.

A characteristic node in a WSN is generally consists of four important units:

- A sensor,
- A processing unit,
- Transceiver, and
- A power supply unit



**Figure 1. Block Diagram of Sensor Node**

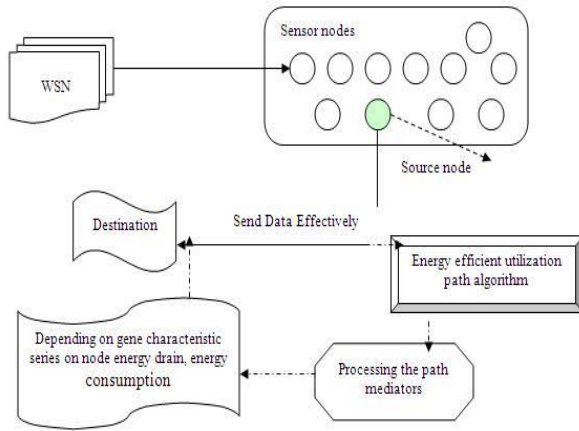
Energy utilization of a node is mainly developed as radio communications. The energy costs of broadcasting and getting a k-bit data packet between two nodes being  $d$  meters separately can be, equally, expressed as

$$E_{Tx}(k, d) \leq k(E_{elec} + \epsilon_{amp} d) \quad (1)$$

and

$$E_{Rx}(k) \leq kE_{elec} \quad (2)$$

$E_{elec}$  denotes the energy utilization due to digital coding, inflection, filtering, and diffusion of the signal, etc and  $\epsilon_{amp}$  is the energy utilized by the source power amplifier. The energy efficiency of WSN is shown in Fig.2



**Figure 2. Architecture diagram of the proposed EEUPA**

The above figure (figure 2) explains the complete process of the EEUPA system. The EEUPA is employed here for sending the sensed set of data from event identified sensor nodes to destination. For sending out the data, efficient energy utilization path algorithm is designed. The efficient energy utilization path algorithm is included by gathering the sequences of data depending on gene characteristics sequencing to improve the life span of the network. The gene characteristic series are estimated depending on the node energy drain rate and the utilization of energy required for sending out the data.

### 3.1. Efficient Energy Utilization Path Algorithm (EEUPA)

In the designed EEUPA, energy utilization of sensed data collection in the network environment is carried out by means of minimum energy utilization path algorithm. Through minimum node energy drain rate, and minimized message overhead, gene contains two properties, Hidden property and Exposed property

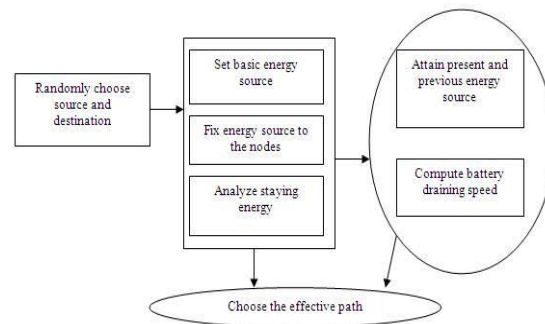
On processor node, drain rate series is sustained in gene property called allele. Path series at

particular occurrence is taken with exposed allele of gene traits. Secondly, hidden allele of gene traits are protected for future sequencing occurrence on exact threshold. Among the frequency sets, we recognize the active path depending on the arrangement of allele in each occurrence of the nodes in the network. Consider a set of sensors planned in a field. The EEUPA system offers the following properties of the Wireless Sensor Networks:

- (1) There are only one destination and sensor nodes in the WSN.
- (2) The destination is inactive and the topology of the WSN collects data generated by the nodes.
- (3) The destination has sufficient power supply while nodes are powered by batteries with restricted energy.
- (4) Nodes do not move after they are planned, to the destination.
- (5) All nodes have similar traits

In the WSNs, recognize the path node mediators from source to sink. After that for every occurrence of node, the utilization of energy and the node energy drain rate is calculated. To offer the effective path, choose the series of gene characteristics which has less energy drain rate and energy utilization. By tracking the system of choosing the path node mediators, the data packet is effectively transferred from source to sink in less amount of time. The algorithm below explains the process of choosing the effective path node mediators.

From Fig.3 primarily packets are sent through the path to reach the destination from source after selecting the S and D from the sensor networks. The staying energy of all the path series of node mediators in the network is subsequently calculated. The staying energy of all nodes for



**Figure 3. Processing the Path Node Mediators**

Particular period of time is accumulated in the file. This is used to calculate the energy draining speed.

The staying energy of nodes and node energy drain rate is calculated by means of the threshold values. Compare the energy utilization of node and the energy drain with their equivalent threshold values to find out the best path. Depending on the calculated value, choose the node which has minimum value of energy utilization and drain rate of energy.

The key task of the efficient energy utilization path algorithm is to choose the maximum network lifespan for a given wireless sensor networks. Consequently, by investigating the properties of gene, the frequency set is recognized. By means of this set, the path for the particular occurrence is recognized and leads it. In the proposed EEUPA work, the path series characteristics are accumulated at every exposed allele and the hidden alleles are preserved for future set of path occurrences.

// Algorithm for Choosing the Effective Path Node Mediators

Input: Source Node S, Destination Node D, routing node Mediator's  $m$   
Start

Choose the S and D in the sensor network to transfer a data packet

Recognize the set of path node mediator's  $m$

Fix threshold value of EU and DR

For each path node  $m$

Calculate the Energy Utilization (EU)

Calculate the Battery Status (BS) in the network

Calculate Current and Prior Energy Status ( $E_{CP}$  and  $E_{PE}$ )

Estimate DR (genes) using Eqn 3

End for each

If EC & DR > threshold value

Transfer the packet via path

Else

Stop the transferring of packets through nodes

End If

End

**Figure 4. Process of Choosing Effective Path Node Mediators**

#### 4. EXPERIMENTAL DISCUSSIONS

We execute experiments to verify the proposed EEUPA system and evaluate its results. In the simulations, we used 1/5 of every data as preparation period, and permit the nodes to create their primary contact history. After the preparation time, we created 1000 messages, all starting from an arbitrary source node to an arbitrary destination node for every t

seconds. The period of the experiment is set as  $t = 300s$ . All messages are consigned a Time-To-Live (TTL) value representing the majority delay restriction.

For the discussions, we have used three set of parameters to calculate the performance of the proposed EEUPA system and compared with the existing works like LEACH (Low Energy Adaptive Cluster Hierarchy), PEGASIS (Power-Efficient Gathering in Sensor Information Systems). The Performance of EEUPA is measured in terms are listed as,

- Node Energy Drain Rate,
- Message Overhead and
- Network Lifespan

#### 5. RESULTS AND DISCUSSION

In this section, we evaluate the attained results with the existing algorithm to calculate the result of the proposed EEUPA system. The below table and graph explains the performance estimation of both the existing and proposed methods.

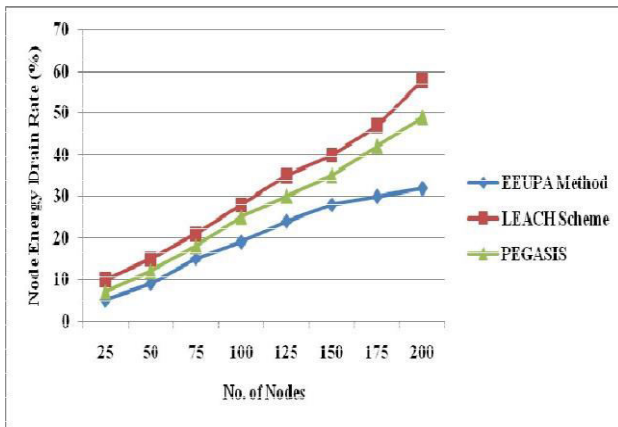
##### 5.1. Node Energy Drain Rate

Table 1. No. of Nodes vs. Node Energy Drain Rate

| No. of Nodes | Node Energy Drain Rate (%) |              |         |
|--------------|----------------------------|--------------|---------|
|              | EEUPA Method               | LEACH Scheme | PEGASIS |
| 25           | 5                          | 10           | 7       |
| 50           | 9                          | 15           | 12      |
| 75           | 15                         | 21           | 18      |
| 100          | 19                         | 28           | 25      |
| 125          | 24                         | 35           | 30      |
| 150          | 28                         | 40           | 35      |
| 175          | 30                         | 47           | 42      |
| 200          | 32                         | 58           | 49      |

The node energy drain rate is calculated based on the number of active nodes in the wireless sensor networks and the values of the proposed Efficient Energy Utilization Path Algorithm is compared with the existing LEACH and PEGASIS schemes and illustrated in Table 1.

Fig.5 explains the drain rate as a method to report for the rate at which energy gets degenerated at a specified node. Every node observes its energy and preserves its battery power drain rate value by taking mean of the amount of energy utilization and calculating the energy dissipation per second. Compared to existing works like LEACH and heuristic model, the proposed EEUPA has small energy drain rate since it used large amount of energy for sending the data between the set of nodes in WSN. However in the proposed work, decrease the surplus dissipation of particular nodes by considering the present traffic situation and by using the drain rate of the staying battery capability.



**Figure 5. No. of Nodes vs. Node Energy Drain Rate**

## 5.2. Message Overhead

The occurrence of message overhead is calculated depending on the number of messages to be sent into the wireless sensor networks and the values of the proposed EEUPA is compared with the existing LEACH and PEGASIS schemes is illustrated in Table 2.

Fig.6 illustrates the existence of message overhead depending on the number of messages to be sent into the network. The message overhead is termed as the number of unsuccessful messages acquired while sending the sensed data from source to destination. Because the path verifying is completed in the proposed EEUPA, results in minimizing the message overhead on the path series for the data collection of sensed events by the destination. By comparing to the other existing works like LEACH and PEGASIS, the proposed EEUPA has fewer messages overhead.

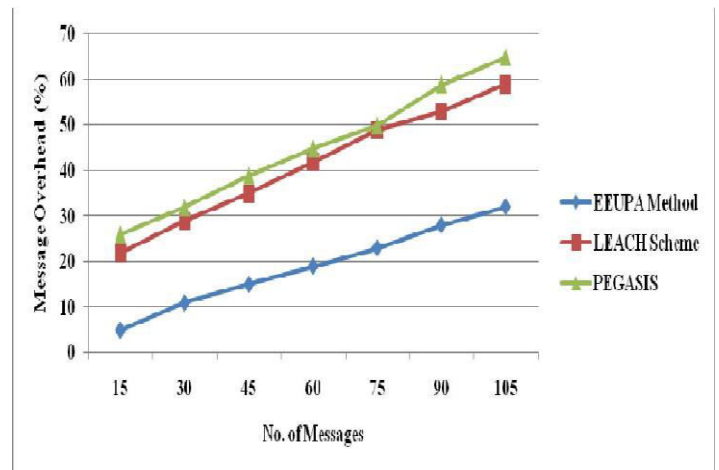
**Table 2. No. Of Messages Vs. Message Overhead**

| No. of Messages | Message Overhead (%) |              |         |
|-----------------|----------------------|--------------|---------|
|                 | EEUPA Method         | LEACH Scheme | PEGASIS |
| 15              | 5                    | 22           | 26      |
| 30              | 11                   | 29           | 32      |
| 45              | 15                   | 35           | 39      |
| 60              | 19                   | 42           | 45      |
| 75              | 23                   | 49           | 50      |
| 90              | 28                   | 53           | 59      |
| 105             | 32                   | 59           | 65      |

## 5.3. Network Lifespan

The lifespan of the network is calculated depending on the number of nodes in the wireless sensor networks and the values of the proposed EEUPA is compared with the existing LEACH and PEGASIS schemes is described in Table 3.

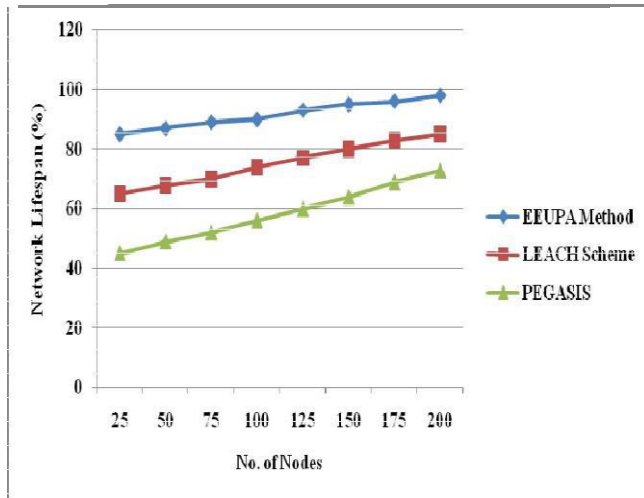
Fig.7 explains the lifespan of the network depending on the number of nodes in the network. Because the energy drain rate of the nodes in the network environment is less, the life span of the network increases correspondingly.



**Figure 6. No. of Messages vs. Message Overhead**

**Table 3: No. Of Nodes vs. Network Lifespan**

| No. of Nodes | Network Lifespan (%) |              |         |
|--------------|----------------------|--------------|---------|
|              | EEUPA Method         | LEACH Scheme | PEGASIS |
| 25           | 85                   | 65           | 45      |
| 50           | 87                   | 68           | 49      |
| 75           | 89                   | 70           | 52      |
| 100          | 90                   | 74           | 56      |
| 125          | 93                   | 77           | 60      |
| 150          | 95                   | 80           | 64      |
| 175          | 96                   | 83           | 69      |
| 200          | 98                   | 85           | 73      |



**Figure 7. No. of Nodes vs. Network Lifespan**

Compared to the existing works like LEACH and PEGASIS, the proposed EEUPA work attains high network lifetime as the utilization of minimal energy drain rate nodes on data collection gene series increases the lifespan of the node and the network.

## 6. CONCLUSION

In this paper we proposed a new method which is employed to calculate the lifespan of nodes constant with present traffic conditions. We described a mechanism, called the EEUPA which cannot be utilized in any of the present wireless sensor path protocol as a way organization principle. This metric is good at reflecting the current dissipation of energy without reflecting on other traffic measurements, like queue length and the number of links passing via the nodes. The main objective of EEUPA is not only to increase the lifespan of every node, but also to extend the lifespan of every link by means of selecting the consistent path from source to destination. With NS-2 simulator, the EEUPA method is compared with the existing works like LEACH and PEGASIS, and the results explained that the EEUPA evades over dissipation of energy, since the sequence path for the particular node is processed with allele of gene traits.

## REFERENCES

- [1].Chien-Ming Chen., Yue-Hsun Lin., Ya-Ching Lin., and Hung-Min Sun., "RCDA: Recoverable Concealed Data Aggregation for Data Integrity in Wireless Sensor Networks," IEEE Transactions on parallel and distributed systems Vol. 23, no. 4, 2012.
- [2].Hongjuan Li., Kai Lin., Keqiu Li., "Energy-efficient and high-accuracy secure data aggregation in wireless sensor networks," Journal on Computer Communications 2011.
- [3].Leandro A. Villas., Azzedine Boukerche, Daniel L. Guidoni., Horacio A.B.F. de Oliveira., Regina Borges de Araujo., Antonio A.F. Loureiro., "An energy-aware spatio-temporal correlation mechanism to perform efficient data collection in wireless sensor networks," Journal on Computer Communications 2012.
- [4].Sudip Misra., P. Dias Thomasinos., "A simple, least-time, and energy-efficient routing protocol with one-level data aggregation for wireless sensor networks," The Journal of Systems and Software, 2010.
- [5].Haitao Zhang., Huadong Maa., Xiang-Yang Li., Shaojie Tang., Xiaohua Xu., "Energy-efficient scheduling with delay constraints for wireless sensor networks: A calculus-based perspective," Journal on Computer Communications 2012.
- [6].Jang-Ping Sheu., Prasan Kumar Sahoo., Chang-Hsin Su., Wei-Kai Huc., "Efficient path planning and data gathering protocols for the wireless sensor network," Journal on Computer Communications 2010.
- [7].Kai Lin a., Min Chenb., Sherali Zeadally., Joel J.P.C. Rodrigues., "Balancing energy consumption with mobile agents in wireless sensor networks," Future Generation Computer Systems 2012.

- [8].Massimo Vecchio., Aline Carneiro Viana., Artur Ziviani., Roy Friedman., "DEEP: Density-based proactive data dissemination protocol for wireless sensor networks with uncontrolled sink mobility," Journal on Computer Communications 2010.
- [9].Basma M. Mohammad El-Basioni., Sherine M. Abd El-kader., Hussein S. Eissa.,Mohammed M. Zahra., "An Optimized Energy-aware Routing Protocol for Wireless Sensor Network," Egyptian Informatics Journal, 2011.
- [10].Chung-Horng Lung., Chenjuan Zhou., "Using hierarchical agglomerative clustering in wireless sensor networks: An energy-efficient and flexible approach," Ad Hoc Networks 2010.
- [11].Manish Kumar Jhaa., Atul Kumar Pandeyb., Dipankar Pala., Anand Mohanc., " An energy-efficient multi-layer MAC (ML-MAC) protocol for wireless sensor networks," International Journal of Electronics and Communications (AEÜ), 2011.
- [12].Jun Yuea., Weiming Zhang., Weidong Xiao., Daquan Tang., Jiuyang Tang., "Energy Efficient and Balanced Cluster-Based Data Aggregation Algorithm for Wireless Sensor Networks," International Workshop on Information and Electronics Engineering (IWIEE), 2012.
- [13].Muhammad Saleem., Israr Ullah., Muddassar Farooq., "BeeSensor: An energy-efficient and scalable routing protocol for wireless sensor networks," Journal on Information Sciences 2012.
- [14].Jalel Ben-Othman., Bashir Yahya.," Energy efficient and QoS based routing protocol for wireless sensor networks," J. Parallel Distrib. Computation, 2010.
- [15].Chih-Chieh Hung., Wen-Chih Peng., "Optimizing in-network aggregate queries in wireless sensor networks for energy saving," Data & Knowledge Engineering 2011.
- [16].Bhoopathy, V. and Parvathi, R.M.S. "Energy Efficient Secure Data Aggregation Protocol for Wireless Sensor Networks", European Journal of Scientific Research, Vol. 50, Issue 1, pp.48-58, 2011.
- [17]. Bhoopathy, V. and Parvathi, R.M.S. "Secure Authentication Technique for Data Aggregation in Wireless Sensor Networks" Journal of Computer Science, Vol. 8, Issue 2, pp 232-238, 2012.
- [18]. Bhoopathy, V. and Parvathi, R.M.S. "Energy Constrained Secure Hierarchical Data Aggregation in Wireless Sensor Networks" American Journal of Applied Sciences, Vol. 9, Issue 6, pp. 858-864, 2012.
- [19]. Bhoopathy, V. and Parvathi, R.M.S. "Securing Node Capture Attacks for Hierarchical Data Aggregation in Wireless Sensor Networks" International Journal of Engineering Research and Applications, Vol. 2, Issue 2, pp. 458-466, 2012.

# Authorized Auditing of Dynamic Big-Data on Cloud

A.S. Gousia Banu, Research Scholar, Department of CSE

Pramod Kumar Singh, Global Program Manager, MBA (Systems), IBM India Pvt. Ltd.

**Abstract**— Cloud computing is widely spreading era. It includes it companies, business line , all online shopping sites including cell phone service providers etc... but in other hand storage capacity and security are increasing issues. Cloud user has no longer direct control over their data, which makes data security, one of the major concerns of using cloud. Previous research work already allows data integrity to be verified without possession of the actual data file. The trusted third party known as auditor. And verification done by this auditor is known as authorized auditing. The Previous system has many drawbacks regarding third party like any one can challenge to the cloud service provider for proof of data integrity. Also in it includes research in Best Least Squares Solution (BLSS) signature algorithm to supporting fully dynamic data updates. This algorithm is used to update an only fixed-sized block known as coarse-grained updates. Though this system takes more time for updating data. In our paper, we are providing a system which support authorized auditing and fine-grained update request. Thus, our system dose not only increases security and flexibility but also providing a new big data application to all cloud service providers for large data frequent small updates.

**Keywords**— Cloud computing, big data, data security, authorized auditing, fine-grained dynamic data update.

## 1. Introduction

Although previous data auditing schemes already have different properties potential risks and inefficiency such as security risks in unauthorized auditing requests and inefficiency in processing small updates still exist. We will focus on better support for small dynamic updates, which benefits the scalability and efficiency of a cloud storage server. To achieve this, our strategy utilizes a flexible data segmentation strategy. Meanwhile, we will address a potential security problem in supporting public verifiability to make the strategy more secure and robust, which is achieved by adding an additional authorization process among the three participating parties of client, Client self-services (CSS) and a third-party auditor (TPA). For providing more security we are using TPA(third party authenticator). This is able to verify our data from cloud and check our data's integrity.

We are providing authenticity to the TPA using md5 hashing algorithm which is going to perform main function in our system .it will allow achieving us the security of our data from TPA also. MD5 hashing algorithm gives 128 bit hash key which is allocate to every TPA which should be given at the time of verifying data at cloud.

## 2. Related Work

[9] "Addressing cloud computing security issues"

The recent emergence of cloud computing has drastically altered everyone's perception of infrastructure architectures, software delivery and development models. Projecting as an evolutionary step, following the transition from mainframe computers to client/server deployment models, cloud computing encompasses elements from grid computing, utility computing and autonomic computing, into an innovative deployment architecture. This rapid transition towards the clouds, has fuelled concerns on a critical issue for the success of information systems, communication and information security. From a security perspective, a number of uncharted risks and challenges have been introduced from this relocation to the clouds, deteriorating much of the effectiveness of traditional protection mechanisms.

As a result the aim of this paper is twofold; firstly to evaluate cloud security by identifying unique security requirements and secondly to attempt to present a viable solution that eliminates these potential threats. This paper proposes introducing a

Trusted Third Party, tasked with assuring specific security characteristics within a cloud environment.

The proposed solution calls upon cryptography, specifically Public Key Infrastructure operating in concert with SSO and LDAP, to ensure the authentication, integrity and confidentiality of involved data and communications. The solution, presents a horizontal level of service, available to all implicated entities, that realizes a security mesh, within which essential trust is maintained.

[3]"a digital signature based on a conventional encryption function"

A new digital signature based only on a conventional encryption function (such as DES) is described which is as secure as the underlying encryption function -- the security does not depend on the difficulty of factoring and the high computational costs of modular arithmetic are avoided.

The signature system can sign an unlimited number of messages, and the signature size increases logarithmically as a function of the number of messages signed. Signature size in a 'typical' system might range from a few hundred bytes to a few kilobytes, and generation of a signature might require a few hundred to a few thousand computations of the underlying conventional encryption function.

[1] "PORs: Proofs of retrievability for Large Files"

In this paper, we define and explore proofs of irretrievability (PORs). A POR scheme enables an archive or back-up service (Prover) to produce a concise proof that a user (verifier) can retrieve a target file  $F$ , that is, that the archive retains and reliably transmits file data sufficient for the user to recover  $F$  in its entirety. A POR may be viewed as a kind of cryptographic proof of knowledge (POK), but one specially designed to handle a large file (or bit string)  $F$ .

We explore POR protocols here in which the communication costs, number of memory accesses for the Prover, and storage requirements of the user (verifier) are small parameters essentially independent of the length of  $F$ . In addition to

proposing new, practical POR constructions, we explore implementation considerations and optimizations that bear on previously explored, related schemes.

In a POR, unlike a POK, neither the Prover nor the verifier need actually have knowledge of  $F$ . PORs give rise to a new and unusual security definition whose formulation is another contribution of our work. We view PORs as an important tool for semi-trusted online archives. Existing cryptographic techniques help users ensure the privacy and integrity of files they retrieve. It is also natural, however, for users to want to verify that archives do not delete or modify files prior to retrieval. The goal of a POR is to accomplish these checks without users having to download the files themselves. A POR can also provide quality-of-service guarantees, i.e., show that a file is retrievable within a certain time bound.

[2] Compact Proofs of Retrievability

In a proof-of-retrievability system, a data storage center must prove to a verifier that he is actually storing all of a client's data. The central challenge is to build systems that are both efficient and provably secure — that is, it should be possible to extract the client's data from any Prover that passes a verification check. In this paper, we give the first proof-of-retrievability schemes with full proofs of security against arbitrary adversaries in the strongest model.

Our first scheme, built from BLS signatures and secure in the random oracle model, features a proof-of-retrievability protocol in which the client's query and server's response are both extremely short. This scheme allows public verifiability: anyone can act as a verifier, not just the file owner. Our second scheme, which builds on pseudorandom functions (PRFs) and is secure in the standard model, allows only private verification. It features a proof-of-retrievability protocol with an even shorter server's response than our first scheme, but the client's query is long. Both schemes rely on Homomorphic properties to aggregate a proof into one small authenticator value.

### 3. Motivation

1. Cost-efficiency brought by elasticity is one of the most important reasons why cloud is being widely adopted. For example, Vodafone Australia is currently using Amazon cloud to provide their users with mobile online-video watching services. Without cloud computing, Vodafone cannot avoid purchasing computing facilities that can process 700 rps, but it will be a total waste for most of the time.

2. Other two large companies who own news.com.au and realestate.com.au, respectively, are using Amazon cloud for the same reason. We can see through these cases that scalability and elasticity, thereby the capability and efficiency in supporting data dynamics, are of extreme importance in cloud computing.

### 4. Purpose and Scope

For providing more security we are using TPA (third party authenticator). This is able to verify our data from cloud and check our data's integrity. We are providing authenticity to the TPA using md5 hashing algorithm which is going to perform main function in our system .it will allow achieving us the security of our data from TPA also. Md5 hashing algorithm gives 128 bit hash key which is allocate to every TPA which should be given at the time of verifying data at cloud.

#### ALGORITHM USED:

##### 1. Message Digestion (MD5):

- i. It Is Designed To Run Effectively On 32-Bit Processor.
- ii. Generate Unique Hash Value For Each Input.
- iii. It Produce Fixed Length 128-Bit Hash Value With No Limit Of Input Message.
- iv. Advantage Is Fast Computing And Uniqueness.
- v. Also Known As Hashing Function.

##### 2. Advanced Encryption Standards (AES)

- I. Secrete Key Generation Algo.
- II. AES Work By Repeating The Same Defined Steps Multiple Times For Encryption & Decryption.
- III. It Operates On Fixed Number Of Bytes.
- IV. Block Size: 128-Bit

V. Key Length: 128,192,256-Bits

VI. Encryption Primitives: Substitution, Shift, Bit Mixing.

### 5. Problem Statement

The challenge/verification process of our strategy, we try to secure the strategy against a malicious CSS who tries to cheat the verifier TPA about the integrity status of the client's data, which is the same as previous work on both PDP and por. In this step, aside from the new authorization process (which will be discussed in detail later in this section), the only difference compared to is the and variable-sectored blocks. Therefore, the security of this phase can be proven through a process highly similar with using the same framework, adversarial model and interactive games defined in. A detailed security proof for this phase is therefore omitted here.

### 6. Proposed System

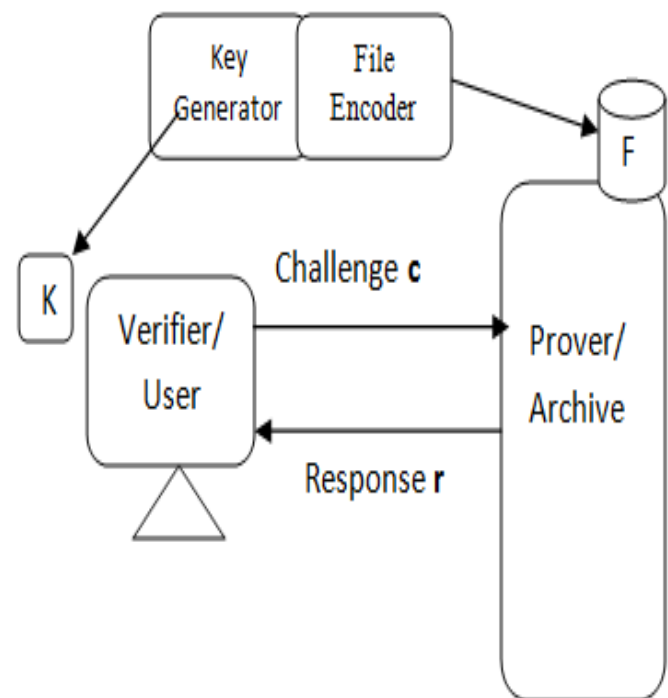


Fig:- Architecture of propose system

#### Authorities of Components:

1. Client will create account
  - select a file

- upload a file to CSS
- updates in file

## 2. Cloud Service Provider (CSP)

- get file
- store file
- convert it in blocks

## 3. Third Party Authenticator (TPA)

- get a file request
- verify file integrity
- challenge to C

As a result, every small update will cause re-computation and updating of the authenticator for an entire file block, which in turn causes higher storage and communication overheads. In this paper, we provide a formal analysis for possible types of fine-grained data updates and propose a strategy that can fully support authorized auditing and fine-grained update requests. Based on our strategy, we also propose an enhancement that can dramatically reduce communication overheads for verifying small updates. Theoretical analysis and experimental results demonstrate that our strategy can offer not only enhanced security and flexibility, but also significantly lower overhead for big data applications with a large number of frequent small updates.

## 7. Result analysis

The below tables and graph shows the comparison of AES with key addition and Key multiplication the number of CPU cycles taken by encryption functions take:

| AES Algorithm           | Encrypt 128 (Cycles) | Encrypt 192 (Cycles) | Encrypt 256 (Cycles) |
|-------------------------|----------------------|----------------------|----------------------|
| With Key addition       | 88.8                 | 86.2                 | 86.06                |
| With Key multiplication | 87.26                | 86.88                | 86.65                |

Fig:- Encryption Table

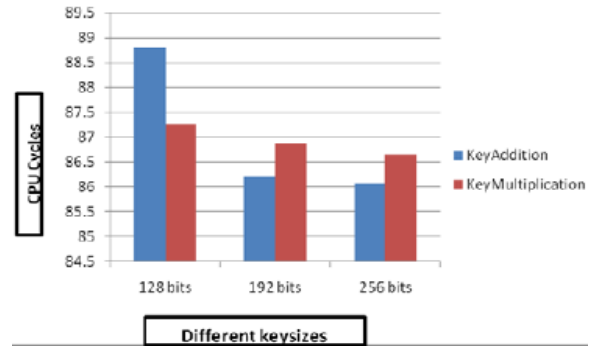


Fig:- Graph of encryption function with different key sizes.

## 8. Conclusion

Thus, in our paper we are providing a formal analysis and fine-grained data updating. Purpose of our strategy is that fully support authorized auditing & fine-grained data updating as per request. Based on our strategy we have also proposed modification that is dramatically reduce communication overheads for verification of small updates. We also plan that for further investigate on the next step how to improve server side protection methods for data security. Hence, in our paper data security, storage and computation, efficient security plays important role under cloud computing context.

## 9. References

1. Juels And B.S. Kaliski Jr., "Pors: Proofs Of Retrieability For Large Files," In Pro. fourteenth Acm Conf. On Comput. What's more, Commun.Security (Ccs), 2007, Pp. 584-597.
2. H. Shacham And B. Waters, "Compact Proofs Of Retrieability,"In Proc. fourteenth Int'l Conf. On Theory And Appl. Of Cryptol. What's more, Inf.Security (Asiacrypt), 2008, Pp. 90-107.
3. R.C. Merkle, "A Digital Signature Based On A Conventional Encryption Function," In Proc. Int'l Cryptol. Conf. On Adv. In Cryptol. (Crypto), 1987, Pp. 369-378.
4. Hadoop Mapreduce. [Online]. Accessible: [Http://Hadoop.Apache.Org](http://Hadoop.Apache.Org)
5. Openstack Open Source Cloud Software, Accessed On: March 25,2013. [Online]. Accessible: [Http://Openstack.Org/](http://Openstack.Org/)
6. Armbrust, A.Fox, R.Griffith, A.D.Joseph, R.Katz, A.Konwinski, G.Lee,D.Patterson, A.Rabkin,I.Stocia,

- And M Zaharia "A View Of Cloud Computing  
."Commum,Acm, Vol.53,No.4,Pp.50-58,Apr.2010
7. Client Presentation Of Amazom Summit Australia,  
Sydney,2012, Accessed On:March  
25,2013.[Online].Available  
:Http://Aws.Amazon.Com/Apac/Awssummit-Au/  
8. D.Boneh, H. Shachhan, And B. Lynn, "Short  
Signatures From The Weil Pairing," J. Cryptoll., Vol.  
17, No. 4, Pp. 297-319, Sept. 2004.
9. D. Zissis And D. Lekkass, "Addressing Coud  
Computing Issues," Future Gen. Comuting Syst., Vol.  
28, No. 3, Pp. 583-592, Mar. 2011.
10. R. Lu et al., —EPPA: An Efficient and Privacy-  
Preserving Aggregation Sheme for Secure Smart  
Grid Communicationsl, IEEE Trans. Parallel  
Distributed System, vol. 23, no. 9, 2012.
11. Certicom, Standards for Efficient Cryptography,  
SEC 1: Elliptic Curve Cryptography, Version 1.0,  
September 2009. pp. 64–76,Apr. 2011
- 12 R. Cramer and V. Shoup. Signature schemes  
based on the strong RSA assumption. ACM Trans.  
Info. & System Security, 3(3):161–85, 2000.
13. R. Cramer and V. Shoup. Design and analysis of  
practical public-key encryption schemes secure  
against adaptive chosen ciphertext attack. SIAM J.  
Computing, 33(1):167–226, 2003.
- 14 Y. Deswarte, J.-J. Quisquater, and A. Saïdane.  
Remote integrity checking. In S. Jajodia and L.  
Strous, editors, Proceedings of IICIS 2003, volume  
140 of IFIP, pages 1–11. Kluwer Academic, Jan.  
2004.
15. Y. Dodis, S. Vadhan, and D. Wichs. Proofs of  
retrievability via hardness amplification. In O.  
Reingold, editor, Proceedings of TCC 2009, volume  
5444 of LNCS, pages 109–27. SpringerVerlag, Mar.  
2009.
- 16 D. Freeman, M. Scott, and E. Teske. A taxonomy  
of pairing-friendly elliptic curves. J. Cryptology,  
23(2):224–80, Apr. 2010.
17. D. Gazzoni Filho and P. Barreto. Demonstrating  
data possession and uncheatable data transfer.  
Cryptology ePrint Archive, Report 2006/150, 2006.  
<http://eprint.iacr.org/>.

#### First A. Author



A. S. Gousia Banu  
Research Scholar  
Department of CSE  
[gbanuzia@gmail.com](mailto:gbanuzia@gmail.com)

#### Second B. Author



Pramod Kumar Singh  
Global Program Manager,  
MBA(Systems),  
IBM India Pvt. Ltd.  
[ziauddinb17@gmail.com](mailto:ziauddinb17@gmail.com)

# Smart Agriculture through IOT

***Hari Krishna,***  
CSE Department,  
mnpt.harikrishna@gmail.com  
Malla Reddy College of Engineering

***Dr.T.Sunil,***  
Prof. & Dean  
Malla Reddy College Engineering,  
sunil.tekale2010@gmail.com

**Abstract:** Agriculture plays vital role in the development of agricultural country. In India about 70% of population depends upon farming and one third of the nation's capital comes from farming. Issues concerning agriculture have been always hindering the development of the country. The only solution to this problem is smart agriculture by modernizing the current traditional methods of agriculture. Hence the project aims at making agriculture smart using automation and IoT technologies. The highlighting features of this project includes smart GPS based remote controlled robot to perform tasks like weeding, spraying, moisture sensing, bird and animal scaring, keeping vigilance, etc. Secondly it includes smart irrigation with smart control and intelligent decision making based on accurate real time field data. Thirdly, smart warehouse management which includes temperature maintenance, humidity maintenance and theft detection in the warehouse. Controlling of all these operations will be through any remote smart device or computer connected to Internet and the operations will be performed by interfacing sensors, Wi-Fi or ZigBee modules, camera and actuators with micro-controller and raspberry pi.

**Keywords:** IoT, automation, Wi-Fi

## I. INTRODUCTION

Agriculture is considered as the basis of life for the human species as it is the main source of food grains and other raw materials. It plays vital role in the growth of country's economy. It also provides large ample employment opportunities to the people. Growth in agricultural sector is necessary for the development of economic condition of the country. Unfortunately, many farmers still use the traditional methods of farming which results in low yielding of crops and fruits. But wherever automation had been implemented and human beings had been replaced by automatic machineries, the yield has been improved. Hence there is need to implement modern science and technology in the agriculture sector for increasing the yield. Most of the papers signifies the use of wireless sensor network which collects the data from different types of sensors and then send it to main server using wireless protocol. The collected data provides the information about different environmental factors which in turns helps to monitor the system. Monitoring environmental factors is not enough and complete solution to improve the yield of the crops. There are number of other factors that affect the productivity to great extent. These factors include attack of insects and pests which can be controlled by spraying the crop with proper insecticide and pesticides. Secondly, attack of wild animals and birds when the crop grows up. There is also possibility of thefts when crop is at the stage of harvesting. Even after harvesting, farmers also face problems in storage of harvested crop. So, in order to provide solutions to all such problems, it is necessary to develop integrated system which will take care of all factors affecting the productivity in every stages like; cultivation, harvesting and post harvesting storage. This paper therefore proposes a system which is useful in monitoring the field data as well as controlling the field operations which provides the

flexibility. The paper aims at making agriculture smart using automation and IoT technologies. The highlighting features of this paper includes smart GPS based remote controlled robot to perform tasks like; weeding, spraying, moisture sensing, bird and animal scaring, keeping vigilance, etc. Secondly, it includes smart irrigation with smart control based on real time field data. Thirdly, smart warehouse management which includes; temperature maintenance, humidity maintenance and theft detection in the warehouse. Controlling of all these operations will be through any remote smart device or computer connected to Internet and the operations will be performed by interfacing sensors, Wi-Fi or ZigBee modules, camera and actuators with micro-controller and raspberry pi.

## II. LITERATURE REVIEW

The newer scenario of decreasing water tables, drying up of rivers and tanks, unpredictable environment present an urgent need of proper utilization of water. To cope up with this use of temperature and moisture sensor at suitable locations for monitoring of crops is implemented in. [1] An algorithm developed with threshold values of temperature and soil moisture can be programmed into a microcontroller-based gateway to control water quantity. The system can be powered by photovoltaic panels and can have a duplex communication link based on a cellular-Internet interface that allows data inspection and irrigation scheduling to be programmed through a web page. [2] The technological development in Wireless Sensor Networks made it possible to use in monitoring and control of greenhouse parameter in precision agriculture. [3] After the research in the agricultural field, researchers found that the yield of agriculture is decreasing day by day. However, use of technology in the field of agriculture

plays important role in increasing the production as well as in reducing the extra man power efforts. Some of the research attempts are done for betterment of farmers which provides the systems that use technologies helpful for increasing the agricultural yield.

A remote sensing and control irrigation system using distributed wireless sensor network aiming for variable rate irrigation, real time in field sensing, controlling of a site specific precision linear move irrigation system to maximize the productivity with minimal use of water was developed by Y. Kim . The system described details about the design and instrumentation of variable rate irrigation, wireless sensor network and real time in field sensing and control by using appropriate software. The whole system was developed using five in field sensor stations which collects the data and send it to the base station using global positioning system (GPS) where necessary action was taken for controlling irrigation according to the database available with the system. The system provides a promising low cost wireless solution as well as remote controlling for precision irrigation. [4]

In the studies related to wireless sensor network, researchers measured soil related parameters such as temperature and humidity. Sensors were placed below the soil which communicates with relay nodes by the use of effective communication protocol providing very low duty cycle and hence increasing the life time of soil monitoring system. The system was developed using microcontroller, universal asynchronous receiver transmitter (UART) interface and sensors while the transmission was done by hourly sampling and buffering the data, transmit it and then checking the status messages. The drawbacks of the system were its cost and deployment of sensor under the soil which causes attenuation of radio frequency (RF) signals. [5]

### III. SYSTEM OVERVIEW

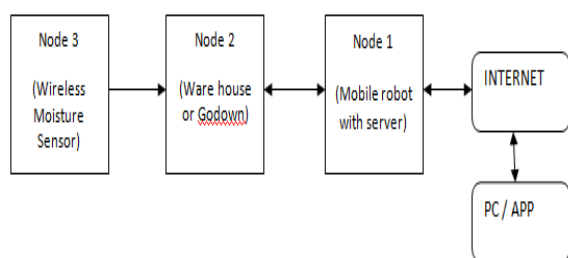


Figure 1: System overview

The paper consist of four sections; node1, node2, node3 and PC or mobile app to control system. In the present system, every node is integration with different sensors and devices and they are interconnected to one central server via wireless communication modules. The server sends and receives information from user end using internet connectivity. There are two modes of operation of the system; auto mode and manual mode. In auto mode system takes its own decisions and controls the installed devices whereas in manual mode user can control the operations of system using android app or PC commands.

### IV. ARCHITECTURE OF THE SYSTEM

Node 1:

Node1 is GPS based mobile robot which can be controlled remotely using computer as well as it can be programmed so as to navigate autonomously within the boundary of field using the co-ordinates given by GPS module.

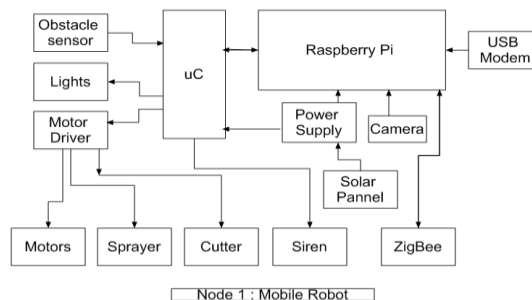


Figure 2: Node 1

The Remote controlled robot have various sensors and devices like camera, obstacle sensor, siren, cutter, sprayer and using them it will perform tasks like; Keeping vigilance, Bird and animal scaring, Weeding, and Spraying

Node 2:

Node2 will be the warehouse. It consists of motion detector, light sensor, humidity sensor, temperature sensor, room heater, cooling fan altogether interfaced with AVR microcontroller. Motion detector will detect the motion in the room when security mode will be ON and on detection of motion, it will send the alert signal to user via Raspberry pi and thus providing theft detection.

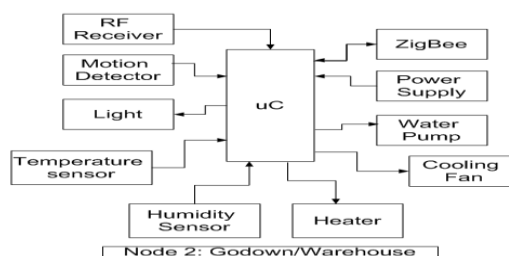


Figure 3: Node 2

Temperature sensor and Humidity sensor senses the temperature and humidity respectively and if the value crosses the threshold then room heater or cooling fan will be switched ON/OFF automatically providing temperature and humidity maintenance. Node2 will also controls water pump depending upon the soil moisture data sent by node3.

Node 3:

Node3 is a smart irrigation node with features like ; Smart control of water pump based on real time field data i.e. automatically turning on/off the pump after attaining the required soil moisture level in auto mode, Switching water pump on/off remotely via mobile or computer in manual mode, and continuous monitoring of soil moisture.

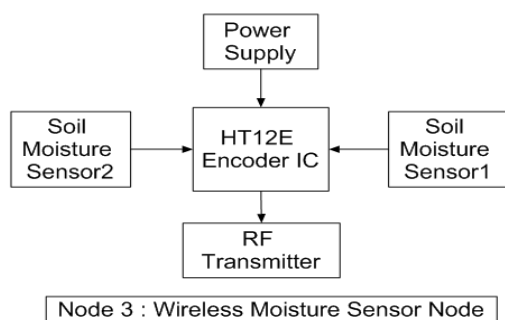


Figure 4: Node 3

In node3, moisture sensor transmits the data using HT12E Encoder IC and a RF transmitter. The transmitted data is received by node2 and there it is processed by microcontroller in order to control the operation of water pump.

Hardware used:

a) AVR Microcontroller Atmega 16/32:

The microcontroller used is, Low-power AVR® 8-bit Microcontroller, having 8K Bytes of In-System Self-programmable Flash program memory, Programmable Serial USART, 8-channel, 10-bit ADC, 23 Programmable I/O Lines.

b) ZigBee Module:

ZigBee is used for achieving wireless communication between Node1 and Node2. The range for Zigbee is roughly 50 meters and it can be increased using high power modules or by using network of modules. It operates on 2.4 GHz frequency. Its power consumption is very low and it is less expensive as compared to other wireless modules like Wi-Fi or Bluetooth. It is usually used to establish wireless local area networks.

c) Temperature Sensor LM35:

The LM35 is precision IC temperature sensor. Output voltage of LM35 is directly proportional to the Centigrade/Celsius of temperature. The LM35 does not need external calibration or trimming to provide accurate temperature range. It is very low cost sensor. It has low output impedance and linear output. The operating temperature range for LM35 is  $-55^{\circ}$  to  $+150^{\circ}\text{C}$ . With rise in temperature, the output voltage of the sensor increases linearly and the value of voltage is given to the microcontroller which is multiplied by the conversion factor in order to give the value of actual temperature.

d) Moisture sensor:

Soil moisture sensor measures the water content in soil. It uses the property of the electrical resistance of the soil. The relationship among the measured property and soil moisture is calibrated and it may vary depending on environmental factors such as temperature, soil type, or electric conductivity. Here, It is used to sense the moisture in field and transfer it to microcontroller in order to take controlling action of switching water pump ON/OFF.

Humidity sensor:

The DHT11 is a basic, low-cost digital temperature and humidity sensor. It gives out digital value and hence there is no need to use conversion algorithm at ADC of the microcontroller and hence we can give its output directly to data pin instead of ADC. It has a capacitive sensor for measuring humidity. The only real shortcoming of this sensor is that one can only get new data from it only after every 2 seconds.

e) Obstacle sensor (Ultra-Sonic):

The ultra-sonic sensor operates on the principle of sound waves and their reflection property. It has two parts; ultra-sonic transmitter and ultra-sonic receiver. Transmitter transmits the 40 KHz sound wave and receiver receives the reflected 40 KHz wave and on its reception, it sends the electrical signal to the microcontroller. The speed of sound in air is already known.

Hence from time required to receive back the transmitted sound wave, the distance of obstacle is calculated. Here, it is used for obstacle detection in case of mobile robot and as a motion detector in ware house for preventing thefts. The ultra-sonic sensor enables the robot to detect and avoid obstacles and also to measure the distance from the obstacle. The range of operation of ultra-sonic sensor is 10 cm to 30 cm.

f) Raspberry Pi :

The Raspberry Pi is small pocket size computer used to do small computing and networking operations. It is the main element in the field of internet of things. It provides access to the internet and hence the connection of automation system with remote location controlling device becomes possible. Raspberry Pi is available in various versions. Here, model Pi 2 model B is used and it has quad-core ARM Cortex-A53 CPU of 900 MHz, and RAM of 1GB. It also has: 40 GPIO pins, Full HDMI port, 4 USB ports, Ethernet port, 3.5mm audio jack, video Camera interface (CSI), the Display interface (DSI), and Micro SD card slot.

Softwares used:

a) AVR Studio Version 4:

It is used to write, build, compile and debug the embedded c program codes which are needed to be burned in the microcontroller in order to perform desired operations. This software directly provides .hex file which can be easily burned into the microcontroller.

b) Proteus 8 Simulator:

Proteus 8 is one of the best simulation software for various circuit designs of microcontroller. It has almost all microcontrollers and electronic components readily available in it and hence it is widely used simulator.

It can be used to test programs and embedded designs for electronics before actual hardware testing. The simulation of programming of microcontroller can also be done in Proteus. Simulation avoids the risk of damaging hardware due to wrong design.

c) Dip Trace:

Dip Trace is EDA/CAD software for creating schematic diagrams and printed circuit boards. The developers provide multi-lingual interface and tutorials (currently available in English and 21 other languages). DipTrace has 4 modules: Schematic Capture Editor, PCB Layout Editor with built-in shape-based auto router and 3D Preview & Export, Component Editor, and Pattern Editor.

d) SinaProg:

SinaProg is a Hex downloader application with AVR Dude and Fuse Bit Calculator. This is used to download code/program and to set fuse bits of all AVR based microcontrollers.

e) Raspbian Operating System:

Raspbian operating system is the free and open source operating system which Debian based and optimized for Raspberry Pi. It provides the basic set of programs and utilities for operating Raspberry Pi. It comes with around 35,000 packages which are pre-compiled softwares that are bundled in a nice format for hassle free installation on Raspberry Pi. It has good community of developers which runs the discussion forms and provides solutions to many relevant problems. However, Raspbian OS is still under consistent development with a main focus on improving the performance and the stability of as many Debian packages as possible.

## V. EXPERIMENTATION AND RESULTS



Figure 5: experimental setup for Node1

As shown in figure 5, experimental setup for node1 consists of mobile robot with central server, GPS module, camera and other sensors. All sensors are successfully interfaced with microcontroller and the microcontroller is interfaced with the raspberry pi. GPS and camera are also connected to raspberry pi. Test results shows that the robot can be controlled remotely using wireless transmission of PC commands to R-Pi. R-Pi forwards the commands to microcontroller and microcontroller gives signals to motor driver in order to drive the Robot. GPS module provides the co-ordinates for the location of the robot

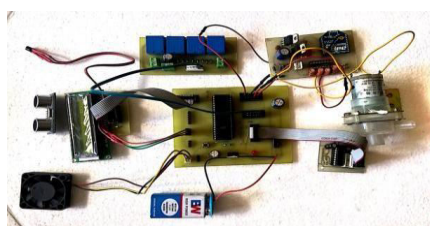


Figure 6: experimental setup for Node2

As shown in above figure, node2 consists of motion detector, temperature sensor, humidity sensor, cooling fan, water pump, etc. connected to the microcontroller board.

The sensors give input to the controller and according to that microcontroller controls the devices in auto mode and also sends the value of sensors to R-Pi and R-Pi forwards it to user's smart device using internet. Test results shows that when temperature level increases above preset threshold level then cooling fan is started automatically in auto mode.

The water pump also gets turned ON if moisture level goes below fixed threshold value. In manual mode, microcontroller receives the controlling signals from R-Pi through ZigBee and accordingly takes the control action.

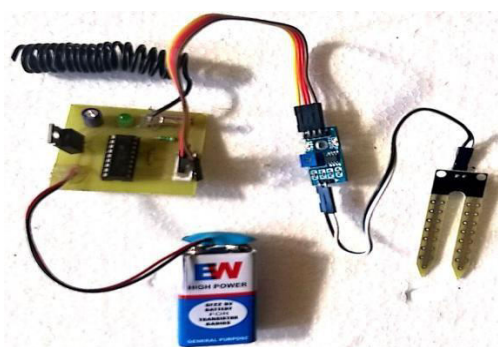


Figure 7: experimental setup for Node3

As shown in above figure, node3 consists of a moisture sensor connected to HT12E. Moisture sensor transmits the data using HT12E Encoder IC and a RF transmitter to the Node2 where it is processed by microcontroller and accordingly water pump is switched ON/OFF.

## VI. CONCLUSION

The sensors and microcontrollers of all three Nodes are successfully interfaced with raspberry pi and wireless communication is achieved between various Nodes.

All observations and experimental tests prove that project is a complete solution to field activities, irrigation problems, and storage problems using remote controlled robot, smart irrigation system and a smart warehouse management system respectively. Implementation of such a system in the field can definitely help to improve the yield of the crops and overall production.

## ACKNOWLEDGMENT

I am sincerely thankful to all my teachers for their guidance for my seminar. Without their help it was tough job for me to accomplish this task. I am especially very thankful to my guide **Dr. R.S. Kawitkar** for his consistent guidance, encouragement and motivation throughout the period of this work. I also want to thank our Head of the Department (E&TC) **Dr. M. B. Mali** for providing me all necessary facilities.

## REFERENCES

- [1] S. R. Nandurkar, V. R. Thool, R. C. Thool, "Design and Development of Precision Agriculture System Using Wireless Sensor Network", IEEE International Conference on Automation, Control, Energy and Systems (ACES), 2014
- [2] Joaquín Gutiérrez, Juan Francisco Villa-Medina, Alejandra Nieto-Garibay, and Miguel Ángel Porta-Gándara, "Automated Irrigation System Using a Wireless Sensor Network and GPRS Module", IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT, 0018-9456, 2013
- [3] Dr. V. Vidya Devi, G. Meena Kumari, "Real-Time Automation and Monitoring System for Modernized Agriculture", International Journal of Review and Research in Applied Sciences and Engineering (IJRRASE) Vol3 No.1. PP 7-12, 2013
- [4] Y. Kim, R. Evans and W. Iversen, "Remote Sensing and Control of an Irrigation System Using a Distributed Wireless Sensor Network", IEEE Transactions on Instrumentation and Measurement, pp. 1379–1387, 2008.
- [5] Q. Wang, A. Terzis and A. Szalay, "A Novel Soil Measuring Wireless Sensor Network", IEEE Transactions on Instrumentation and Measurement, pp. 412–415, 2010
- [6] Yoo, S.; Kim, J.; Kim, T.; Ahn, S.; Sung, J.; Kim, D. A2S: Automated agriculture system based on WSN. In ISCE 2007. IEEE International Symposium on Consumer Electronics, 2007, Irving, TX, USA, 2007
- [7] Arampatzis, T.; Lygeros, J.; Manesis, S. A survey of applications of wireless sensors and Wireless Sensor Networks. In 2005 IEEE International Symposium on Intelligent Control & 13<sup>th</sup> Mediterranean Conference on Control and Automation. Limassol, Cyprus, 2005, 1-2, 719-724
- [8] Orazio Mirabella and Michele Brischetto, 2011. "A Hybrid Wired/Wireless Networking Infrastructure for Greenhouse Management", IEEE transactions on instrumentation and measurement, vol. 60, no. 2, pp 398-407.
- [9] N. Kotamaki and S. Thessler and J. Koskiahio and A. O. Hannukkala and H. Huitu and T. Huttula and J. Havento and M. Jarvenpaa (2009). "Wireless in-situ sensor network for agriculture and water monitoring on a river basin scale in Southern Finland: evaluation from a data users perspective". Sensors 4, 9: 2862-2883. doi:10.3390/s90402862 2009.
- [10] Liu, H.; Meng, Z.; Cui, S. A wireless sensor network prototype for environmental monitoring in greenhouses. International Conference on Wireless Communications, Networking and Mobile Computing (WiCom 2007), Shanghai, China; 21-25 September 2007.
- [11] Baker, N. ZigBee and bluetooth - Strengths and weaknesses for industrial applications. Comput. Control. Eng. 2005, 16, 20-25.
- [12] IEEE, Wireless medium access control (MAC) and physical layer (PHY) specifications for low rate wireless personal area networks (LR-WPANs). In The Institute of Electrical and Electronics Engineers Inc.: New York, NY, USA, 2003.

# Spectrum Sensing For Performance Improvement

**V.Uday Kiran**

CSE Department

udaykiran.velpula@gmail.com

Malla Reddy College of Engineering

**Dr.V.Bhoopathy**

Professor, CSE Department

v.bhoopathy@gmail.com

Malla Reddy College of Engineering

**Abstract**— Cognitive radio has emerged as one of the most promising candidate solutions to improve spectrum utilization in next generation cellular networks. A crucial requirement for future cognitive radio networks is wideband spectrum sensing: secondary users reliably detect spectral opportunities across a wide frequency range. In this article, various wideband spectrum sensing algorithms are presented, together with a discussion of the pros and cons of each algorithm and the challenging issues. Special attention is paid to the use of sub-Nyquist techniques, including compressive sensing and multi-channel sub-Nyquist sampling techniques.

**Index Terms**— Cellular network, cognitive radio, compressive sensing, spectrum sensing, sub-Nyquist sampling, wideband spectrum sensing.

## I. INTRODUCTION

Radio frequency (RF) spectrum is a valuable but tightly regulated resource due to its unique and important role in wireless communications. With the proliferation of wireless services, the demands for the RF spectrum are constantly increasing, leading to scarce spectrum resources. On the other hand, it has been reported that localized temporal and geographic spectrum utilization is extremely low [1]. Currently, new spectrum policies are being developed by the Federal Communications Commission (FCC) that will allow secondary users to opportunistically access a licensed band, when the primary user (PU) is absent. Cognitive radio [2], [3] has become a promising solution to solve the spectrum scarcity problem in the next generation cellular networks by exploiting opportunities in time, frequency, and space domains.

Cognitive radio is an advanced software-defined radio that automatically detects its surrounding RF stimuli and intelligently adapts its operating parameters to network infrastructure while meeting user demands. Since cognitive radios are considered as secondary users for using the licensed spectrum, a crucial requirement of cognitive radio networks is that they must efficiently exploit under-utilized spectrum (denoted as spectral opportunities) without causing harmful interference to the PUs. Furthermore, PUs have no obligation to share and change their operating parameters for sharing spectrum with cognitive radio networks. Hence, cognitive radios should be able to independently detect spectral

opportunities without any assistance from PUs; this ability is called spectrum sensing, which is considered as one of the most critical components in cognitive radio networks.

Many narrowband spectrum sensing algorithms have been studied in the literature [4] and references therein, including matched-filtering, energy detection [5], and cyclostationary feature detection. While present narrowband spectrum sensing algorithms have focused on exploiting spectral opportunities over narrow frequency range, cognitive radio networks will eventually be required to exploit spectral opportunities over wide frequency range from hundreds of megahertz (MHz) to several gigahertz (GHz) for achieving higher opportunistic throughput. This is driven by the famous Shannon's formula that, under certain conditions, the maximum theoretically achievable bit rate is directly proportional to the spectral bandwidth. Hence, different from narrowband spectrum sensing, wideband spectrum sensing aims to find more spectral opportunities over wide frequency range and achieve higher opportunistic aggregate throughput in cognitive radio networks. However, conventional wideband spectrum sensing techniques based on standard analog-to-digital converter (ADC) could lead to unaffordably high sampling rate or implementation complexity; thus, revolutionary wideband spectrum sensing techniques become increasingly important.

In the remainder of this article, we first briefly introduce the traditional spectrum sensing algorithms for narrowband sensing in Section II. Some challenges for realizing wideband spectrum sensing are then discussed in Section III. In addition, we categorize the existing wideband spectrum sensing algorithms based on their implementation types, and review the state-of-the-art techniques for each category. Future research challenges for implementing wideband spectrum sensing are subsequently identified in Section IV, after which concluding remarks are given in Section V.

## II. NARROWBAND SPECTRUM SENSING

The most efficient way to sense spectral opportunities is to detect active primary transceivers in the vicinity of cognitive radios. However, as primary receivers may be passive, such as TVs, some receivers are difficult to detect in practice. An alternative is to detect the primary transmitters by using

traditional narrowband sensing algorithms, including matched-filtering, energy detection, and cyclostationary feature detection as shown in Fig. 1. Here, the term “narrowband” implies that the frequency range is sufficiently narrow such that the channel frequency response can be considered flat. In other words, the bandwidth of our interest is less than the coherence bandwidth of the channel. The implementation of these narrowband algorithms requires different conditions, and their detection performance are correspondingly distinguished. The advantages and disadvantages of these algorithms are summarized in Table I.

The matched-filtering method is an optimal approach for spectrum sensing since it maximizes the signal-to-noise ratio (SNR) in the presence of additive noise. This advantage is achieved by correlating the received signal with a template for detecting the presence of a known signal in the received signal. However, it relies on prior knowledge of the PUs and requires cognitive radios to be equipped with carrier synchronization and timing devices, leading to increased implementation complexity. Energy detection [5] is a non-coherent detection method that avoids the need for prior knowledge of the PUs and the complicated receivers required by a matched filter. Both the implementation and the computational complexity are relatively low. A major drawback is that it has poor detection performance under low SNR scenarios and cannot differentiate between the signals from PUs and the interference from other cognitive radios. Cyclostationary feature detection method detects and distinguishes between different types of primary signals by exploiting their cyclostationary features. However, the computational cost of such an approach is relatively high, because it requires to calculate a two-dimensional function dependent on both frequency and cyclic frequency.

### III. WIDEBAND SPECTRUM SENSING

Against narrowband techniques as mentioned above, wideband spectrum sensing techniques aim to sense a frequency bandwidth that exceeds the coherence bandwidth of the channel. For example, for exploiting spectral opportunities in the whole ultra-high frequency (UHF) TV band (between 300 MHz and 3 GHz), wideband spectrum sensing techniques should be employed. We note that narrowband sensing techniques cannot be directly used for performing wideband spectrum sensing, because they make a single binary decision for the whole spectrum and thus cannot identify individual spectral opportunities that lie within the wideband spectrum. As shown in Table II, wideband spectrum sensing can be broadly categorized into two types: Nyquist wideband sensing and sub-Nyquist wideband sensing. The former type processes digital signals taken at or above the Nyquist rate, whereas the latter type acquires signals using sampling rate lower than the Nyquist rate. In the rest of this article, we will provide an overview of the state-of-the-art wideband spectrum sensing algorithms and discuss the pros and cons of each algorithm.

#### A. Nyquist Wideband Sensing

A simple approach of wideband spectrum sensing is to directly acquire the wideband signal using a standard ADC and then use digital signal processing techniques to detect spectral opportunities. For example, Quan et al. [6] proposed a multi-band joint detection algorithm that can sense the primary signal over multiple frequency bands. As shown in Fig. 2(a), the wideband signal  $x(t)$  was firstly sampled by a high sampling rate ADC, after which a serial to parallel conversion circuit (S/P) was used to divide sampled data into parallel data streams. Fast Fourier transform (FFT) was used to convert the wideband signals to the frequency domain. The wideband spectrum  $X(f)$  was then divided into a series of narrowband spectra  $X_1(f), \dots, X_v(f)$ . Finally, spectral opportunities were detected using binary hypotheses tests, where  $H_0$  denotes the absence of PUs and  $H_1$  denotes the presence of PUs. The optimal detection threshold was jointly chosen by using optimization techniques. Such an algorithm can achieve better performance than the single band sensing case.

Furthermore, by also using a standard ADC, Tian and Giannakis proposed a wavelet-based spectrum sensing algorithm in [7]. In this algorithm, the power spectral density (PSD) of the wideband spectrum (denoted as  $S(f)$ ) was modeled as a train of consecutive frequency subbands, where the PSD is smooth within each subband but exhibits discontinuities and irregularities on the border of two neighboring subbands. The wavelet transform was then used to locate the singularities of the wideband PSD, and the wideband spectrum sensing was formulated as a spectral edge detection problem as shown in Fig. 2(b).

However, special attention should be paid to the signal sampling procedure. In these algorithms, sampling signals should follow Shannon’s celebrated theorem: the sampling rate must be at least twice the maximum frequency present in the signal (known as Nyquist rate) in order to avoid spectral aliasing. Suppose that the wideband signal has frequency range  $0 \sim 10$  GHz, it should be uniformly sampled by a standard ADC at or above the Nyquist rate 20 GHz which will be unaffordable for next generation cellular networks. Therefore, sensing wideband spectrum presents significant challenges on building sampling hardware that operates at a sufficiently high rate, and designing high-speed signal processing algorithms. With current hardware technologies, high-rate ADCs with high resolution and reasonable power consumption (e.g., 20 GHz sampling rate with 16 bits resolution) are difficult to implement. Even if it comes true, the real-time digital signal processing of sampled data could be very expensive.

One naive approach that could relax the high sampling rate requirement is to use superheterodyne (frequency mixing) techniques that “sweep” across the frequency range of interest as shown in Fig. 2(c). A local oscillator (LO) produces a sine wave that mixes with the wideband signal and down-converts it to a lower frequency. The down-converted signal is then filtered by a bandpass filter (BPF), after which existing narrowband spectrum sensing techniques in Section II can be

applied. This sweep-tune approach can be realized by using either a tunable BPF or a tunable LO. However, this approach is often slow and inflexible due to the sweep-tune operation. Another solution would be the filter bank algorithm presented by Farhang-Boroujeny [8] as shown in Fig. 2(d). A bank of prototype filters (with different shifted central frequencies) was used to process the wideband signal. The base-band can be directly estimated by using a prototype filter, and other bands can be obtained through modulating the prototype filter. In each band, the corresponding portion of the spectrum for the wideband signal was down-converted to base-band and then low-pass filtered. This algorithm can therefore capture the dynamic nature of wideband spectrum by using low sampling rates. Unfortunately, due to the parallel structure of the filter bank, the implementation of this algorithm requires a large number of RF components.

### B. Sub-Nyquist Wideband Sensing

Due to the drawbacks of high sampling rate or high implementation complexity in Nyquist systems, sub-Nyquist approaches are drawing more and more attention in both academia and industry. Sub-Nyquist wideband sensing refers to the procedure of acquiring wideband signals using sampling rates lower than the Nyquist rate and detecting spectral opportunities using these partial measurements. Two important types of sub-Nyquist wideband sensing are compressive sensing-based wideband sensing and multi-channel sub-Nyquist wideband sensing. In the subsequent paragraphs, we give some discussions and comparisons regarding these sub-Nyquist wideband sensing algorithms.

1) *Compressive Sensing-based Wideband Sensing:* Compressive sensing is a technique that can efficiently acquire a signal using relatively few measurements, by which unique representation of the signal can be found based on the signal's sparseness or compressibility in some domain. As the wideband spectrum is inherently sparse due to its low spectrum utilization, compressive sensing becomes a promising candidate to realize wideband spectrum sensing by using subNyquist sampling rates. Tian and Giannakis firstly introduced compressive sensing theory to sense wideband spectrum in [9]. This technique used fewer samples closer to the information rate, rather than the inverse of the bandwidth, to perform wideband spectrum sensing. After reconstruction of the wideband spectrum, wavelet-based edge detection was used to detect spectral opportunities across wideband spectrum.

Furthermore, to improve the robustness against noise uncertainty, Tian et al. [10] studied a cyclic feature detection-based compressive sensing algorithm for wideband spectrum sensing. It can successfully extract second-order statistics of wideband signals from digital samples taken at sub-Nyquist rates. The 2-D cyclic spectrum (spectral correlation function) of a wideband signal can be directly reconstructed from the compressive measurements. In addition, such an algorithm is

also valid for reconstructing the power spectrum of wideband signal, which is useful if the energy detection algorithm is used for detecting spectral opportunities.

For further reducing the data acquisition cost, Zeng et al. [11] proposed a distributed compressive sensing-based wideband sensing algorithm for cooperative multi-hop cognitive radio networks. By enforcing consensus among local spectral estimates, such a collaborative approach can benefit from spatial diversity to mitigate the effects of wireless fading. In addition, decentralized consensus optimization algorithm was proposed that aims to achieve high sensing performance at a reasonable computational cost.

However, compressive sensing has concentrated on finite-length and discrete-time signals. Thus, innovative technologies are required to extend the compressive sensing to continuoustime signal acquisition, i.e., implementing compressive sensing in analog domain. To realize the analog compressive sensing, Tropp et al. [12] proposed an analog-to-information converter (AIC), which could be a good basis for the above-mentioned algorithms. As shown in Fig. 3(a), the AIC-based model consists of a pseudo-random number generator, a mixer, an accumulator, and a low-rate sampler. The pseudo-random number generator produces a discrete-time sequence that demodulates the signal  $x(t)$  by a mixer. The accumulator is used to sum the demodulated signal for  $1/w$  seconds, while its output signal is sampled using a low sampling rate. After that, the sparse signal can be directly reconstructed from partial measurements using compressive sensing algorithms. Unfortunately, it has been identified that the performance of AIC model can be easily affected by design imperfections or model mismatches.

2) *Multi-channel Sub-Nyquist Wideband Sensing:* To circumvent model mismatches, Mishali and Eldar proposed a modulated wideband converter (MWC) model in [13] by modifying the AIC model. The main difference between MWC and AIC is that MWC has multiple sampling channels, with the accumulator in each channel replaced by a general low-pass filter. One significant benefit of introducing parallel channel structure in Fig. 3(b) is that it provides robustness against the noise and model mismatches. In addition, the dimension of the measurement matrix is reduced, making the spectral reconstruction more computationally efficient. An alternative multi-channel sub-Nyquist sampling approach is the multi-coset sampling as shown in Fig. 3(c). The multi-coset sampling is equivalent to choosing some samples from a uniform grid, which can be obtained using a sampling rate  $f_s$  higher than the Nyquist rate. The uniform grid is then divided into blocks of  $m$  consecutive samples, and in each block  $v$  ( $v < m$ ) samples are retained while the rest of samples are skipped. Thus, the multi-coset sampling is often implemented by using  $v$  sampling channels with sampling rate of  $f_s/m$ , with different sampling channels having different time offsets. To obtain a unique solution for the wideband spectrum from these partial measurements, the sampling pattern should be carefully designed. In [14], some sampling patterns were proved to be

valid for unique signal reconstruction. The advantage of multi-coset approach is that the sampling rate in each channel is  $m$  times lower than the Nyquist rate. Moreover, the number of measurements is only  $v$ -mth of that in the Nyquist sampling case. One drawback of the multi-coset approach is that the channel synchronization should be met such that accurate time offsets between sampling channels are required to satisfy a specific sampling pattern for a robust spectral reconstruction.

To relax the multi-channel synchronization requirement, asynchronous multi-rate wideband sensing approach was studied in [15]. In this approach, sub-Nyquist sampling was induced in each sampling channel to wrap the sparse spectrum occupancy map onto itself; the sampling rate can therefore be significantly reduced. By using different sampling rates in different sampling channels as shown in Fig. 3(d), the performance of wideband spectrum sensing can be improved. Specifically, in the same observation time, the numbers of samples in multiple sampling channels are chosen as different consecutive prime numbers. Furthermore, as only the magnitudes of subNyquist spectra are of interest, such a multi-rate wideband sensing approach does not require perfect synchronization between multiple sampling channels, leading to easier implementation.

#### IV. RESEARCH CHALLENGES

In this section, we identify the following research challenges that need to be addressed for implementing a feasible wideband spectrum sensing device for future cognitive radio networks.

##### A. Sparse Basis Selection

Nearly all sub-Nyquist wideband sensing techniques require that the wideband signal should be sparse in a suitable basis. Given the low spectrum utilization, most of existing wideband sensing techniques assumed that the wideband signal is sparse in the frequency domain, i.e., the sparsity basis is a Fourier matrix. However, as the spectrum utilization improves, e.g., due to the use of cognitive radio techniques in future cellular networks, the wideband signal may not be sparse in the frequency domain any more. Thus, a significant challenge in future cognitive radio networks is how to perform wideband sensing using partial measurements, if the wideband signal is not sparse in the frequency domain. It will be essential to study appropriate wideband sensing techniques that are capable of exploiting sparsity in any known sparsity basis. Furthermore, in practice, it may be difficult to acquire sufficient knowledge about the sparsity basis in cognitive radio networks, e.g., when we cannot obtain enough prior knowledge about the primary signals. Hence, future cognitive radio networks will be required to perform wideband sensing when the sparsity basis is not known. In this context, a challenging issue is to study “blind” sub-Nyquist wideband sensing algorithms, where we do not require prior knowledge about the sparsity basis for the sub-Nyquist sampling or the

spectral reconstruction.

##### B. Adaptive Wideband Sensing

In most of sub-Nyquist wideband sensing systems, the required number of measurements will proportionally change when the sparsity level of wideband signal varies. Therefore, sparsity level estimation is often required for choosing an appropriate number of measurements in cognitive radio networks. However, in practice, the sparsity level of wideband signal is often time-varying and difficult to estimate, because of either the dynamic activities of PUs or the timevarying fading channels between PUs and cognitive radios. Due to this sparsity level uncertainty, most of sub-Nyquist wideband sensing systems should pessimistically choose the number of measurements, leading to more energy consumption in cellular networks. As shown in Fig. 4, more measurements (i.e.,  $0.38N$  rather than  $0.25N$  measurements for achieving the success recovery rate 0.9) are required for the sparsity uncertainty between 10 and 20, which does not fully exploit the advantages of using sub-Nyquist sampling technologies. Hence, future cognitive radio networks should be capable of performing wideband sensing, given the unknown or timevarying sparsity level. In such a scenario, it is very challenging to develop adaptive wideband sensing techniques that can intelligently/quickly choose an appropriate number of compressive measurements without the prior knowledge of the sparsity level.

##### C. Cooperative Wideband Sensing

In a multipath or shadow fading environment, the primary signal as received at cognitive radios may be severely degraded, leading to unreliable wideband sensing results in each cognitive radio. In this situation, future cognitive radio networks should employ cooperative strategies for improving the reliability of wideband sensing by exploiting spatial diversity. Actually, in a cluster-based cognitive radio network, the wideband spectra as observed by different cognitive radios could share some common spectral components, while each cognitive radio may observe some innovative spectral components. Thus, it is possible to fuse compressive measurements from different nodes and exploit the spectral correlations among cognitive radios in order to save the total number of measurements and thus the energy consumption in cellular networks. Such a data fusion-based cooperative technique, however, will lead to heavy data transmission burden in the common control channels. It is therefore challenging to develop data fusion-based cooperative wideband sensing techniques subject to relaxed data transmission burden. An alternative is to develop decision fusion-based wideband sensing techniques, if each cognitive radio is able to detect wideband spectrum independently. Due to the limited computational resource in cellular networks, the challenge that remains in the decision fusion-based cooperative approach is how to appropriately combine information in real time.

## V. CONCLUSION

In this article, we first addressed the challenges in the design and implementation of wideband spectrum sensing algorithms for the cognitive radio-based next generation cellular networks. Then, we categorized the existing wideband spectrum sensing algorithms based on their sampling types and discussed the pros and cons of each category. Moreover, motivated by the fact that wideband spectrum sensing is critical for reliably finding spectral opportunities and achieving opportunistic spectrum access for next generation cellular networks, we made a brief survey of the state-of-the-art wideband spectrum sensing algorithms. Finally, we presented several open research issues for implementing wideband spectrum sensing.

## REFERENCES

- [1] M. McHenry, "NSF spectrum occupancy measurements project summary," Shared Spectrum Company, Tech. Rep., Aug. 2005.
- [2] C.-X. Wang, X. Hong, H.-H. Chen, and J. S. Thompson, "On capacity of cognitive radio networks with average interference power constraints," *IEEE Trans. Wireless Commun.*, vol. 8, no. 4, Apr. 2009, pp. 1620–1625.
- [3] X. Hong, C.-X. Wang, H.-H. Chen, and Y. Zhang, "Secondary spectrum access networks: recent developments on the spatial models," *IEEE Vehi. Technol. Mag.*, vol. 4, no. 2, June 2009, pp. 36–43.
- [4] T. Yucek and H. Arslan, "A survey of spectrum sensing algorithms for cognitive radio applications," *IEEE Commun. Surveys and Tutorials*, vol. 11, no. 1, Jan. 2009, pp. 116–130.
- [5] H. Sun, D. Laurenson, and C.-X. Wang, "Computationally tractable model of energy detection performance over slow fading channels," *IEEE Commun. Letters*, vol. 14, no. 10, Oct. 2010, pp. 924–926.
- [6] Z. Quan, S. Cui, A. H. Sayed, and H. V. Poor, "Optimal multiband joint detection for spectrum sensing in cognitive radio networks," *IEEE Trans. Signal Processing*, vol. 57, no. 3, Mar. 2009, pp. 1128–1140.
- [7] Z. Tian and G. Giannakis, "A wavelet approach to wideband spectrum sensing for cognitive radios," in *Proc. IEEE Cognitive Radio Oriented Wireless Networks and Commun.*, Mykonos Island, Greece, June 2006, pp. 1–5.
- [8] B. Farhang-Boroujeny, "Filter bank spectrum sensing for cognitive radios," *IEEE Trans. Signal Processing*, vol. 56, no. 5, May 2008, pp. 1801–1811.
- [9] Z. Tian and G. Giannakis, "Compressive sensing for wideband cognitive radios," in *Proc. IEEE International Conference on Acoustics, Speech, and Signal Processing*, Honolulu, HI, USA, April 2007, pp. 1357–1360.
- [10] Z. Tian, Y. Tefesse, and B. M. Sadler, "Cyclic feature detection with sub-Nyquist sampling for wideband spectrum sensing," *IEEE J. Sel. Topics in Signal Processing*, vol. 6, no. 1, Feb. 2012, pp. 58–69.
- [11] F. Zeng, C. Li, and Z. Tian, "Distributed compressive spectrum sensing in cooperative multihop cognitive networks," *IEEE J. Sel. Topics in Signal Processing*, vol. 5, no. 1, Feb. 2011, pp. 37–48.
- [12] J. A. Tropp, J. N. Laska, M. F. Duarte, J. K. Romberg, and R. G. Baraniuk, "Beyond Nyquist: Efficient sampling of sparse bandlimited signals," *IEEE Trans. Information Theory*, vol. 56, no. 1, Jan. 2010, pp. 520–544.
- [13] M. Mishali and Y. C. Eldar, "Blind multiband signal reconstruction: Compressive sensing for analog signals," *IEEE Trans. Signal Processing*, vol. 57, no. 3, March 2009, pp. 993–1009.
- [14] R. Venkataramani and Y. Bresler, "Perfect reconstruction formulas and bounds on aliasing error in sub-Nyquist nonuniform sampling of multiband signals," *IEEE Trans. Information Theory*, vol. 46, no. 6, Sep. 2000, pp. 2173–2183.
- [15] H. Sun, W.-Y. Chiu, J. Jiang, A. Nallanatah, and H. V. Poor, "Wideband spectrum sensing with sub-Nyquist sampling in cognitive radios," *IEEE Trans. Signal Processing*, vol. 60, no. 11, Nov. 2012, pp. 6068–6073.

# **SURVEY ON CRIME ANALYSIS AND PREDICTION USING DATA MINING TECHNIQUES**

**G.Soundarya,**  
CSE Department,  
nandugudipudi@gmail.com  
Malla Reddy College of Engineering

**Dr.P.Mani Kandan**  
Proffessor, CSE Department,  
mani.p.mk@gmail.com  
Malla Reddy Engineering College for Women,

## **Abstract**

*Data Mining is the procedure which includes evaluating and examining large pre-existing databases in order to generate new information which may be essential to the organization. The extraction of new information is predicted using the existing datasets. Many approaches for analysis and prediction in data mining had been performed. But, many few efforts has made in the criminology field. Many few have taken efforts for comparing the information all these approaches produce. The police stations and other similar criminal justice agencies hold many large databases of information which can be used to predict or analyze the criminal movements and criminal activity involvement in the society. The criminals can also be predicted based on the crime data. The main aim of this work is to perform a survey on the supervised learning and unsupervised learning techniques that has been applied towards criminal identification. This paper presents the survey on the Crime analysis and crime prediction using several Data Mining techniques.*

## **Keywords:**

*Criminology, Crime Analysis, Crime Prediction, Data Mining*

## **1. INTRODUCTION**

Historically solving crimes has been the right of the criminal justice and law enforcement specialists. With the increase in the use of the computerized systems to track crimes and trace criminals, computer data analysts have started lending their hands in helping the law enforcement officers and detectives to speed up the process of solving crimes. Criminology is process that is used to identify crime and criminal characteristics. The criminals and the crime occurrence possibility can be assessed with the help of criminology techniques. The criminology aids the police department, the detective agencies and crime branches in identifying the true characteristics of a criminal. The criminology department has been used in the proceedings of crime tracking ever since 1800. Crimes are a social nuisance and cost our society dearly in several ways. Even, the Indian Government has taken steps to develop applications and software for the use of State and Central Police in relation with the National Crime Records Bureau (NCRB) [27]. Any research that can help in solving crimes faster will pay for itself. About 10% of the criminals commit about 50% of the crimes [15]. People who study criminology will be able to identify the criminals based on the traces, characteristics and methods of crime which can be collected from the crime scene. In the middle of 1990s, data mining came into existence as a strong tool to extract useful information from large datasets and find the relationship between the attributes of the data [11]. Data mining originally came from statistics and machine learning as an interdisciplinary field, but then it was grown a lot that in 2001 it was considered as one of the top 10 leading technologies which will change the world [12]. According to many researchers such

as Nath [23], solving crimes is a difficult and time consuming task that requires human intelligence and experience and data mining is one technique that can help us with crime detection problems. For solving crimes faster we have to develop a data mining paradigm that performs an interdisciplinary approach between computer science and criminal justice. As said earlier, the Criminology is a process that aims to identify crime characteristics and it is one of the most important fields for applying data mining. By using this, data mining algorithms will be able to produce crime reports and help in the identification of criminals much faster than any human could. Because of this remarkable feature, there is a growing demand for data mining in criminology. Actually, Crime analysis is a process which includes exploring the behavior of the crimes, detecting crimes and their relationships with criminals. The huge volume of crime and criminal datasets and the complexity of relationships between these kinds of information have made criminology an appropriate field for applying data mining techniques. Identifying crime characteristics is the first step for proceeding with any further analysis. The quality of data analysis depends greatly on background knowledge of analyst. A criminal can range from civil infractions such as illegal driving to terrorism mass murder such as the 9/11 attacks, therefore it is difficult to model the perfect algorithm to cover all of them [21]. The knowledge that is gained from Data Mining approaches is a very useful and this can help and support, the police. More specifically, we can use classification and clustering based models to help in identification of crime patterns and criminals. The wide range of data mining applications in the criminology has made it an important field of research. Data mining systems have played as a key role in assisting humans in this forensic domain and criminology domain. This makes it one of the most challenging decision-making environments for research.

The motivation for proceeding with this survey work is to aid a helping hand to the young researchers who are performing their research in criminal analysis and crime prediction areas. The paper is organized in such a manner to provide insights about the crime analysis procedure and then produce different types of crime analysis operations and those which can be applied together for producing an end user product which can be applied to the crime analysis in any police stations and detective agencies. This work will be a valuable reference to those who precede their research work in the crime analysis and Crime prediction using data mining techniques.

This survey paper is organized in such a manner for easy understanding of the concepts. The general crime analysis procedure is discussed in section 2. The Criminal analysis methods are discussed in the section 3 which will include all the different types of methods grouped under their own categories. Finally, section 4 gives the Qualitative analysis of the Crime Analysis and Prediction techniques and section 5 gives the

Quantitative analysis of the Crime Analysis and Prediction techniques.

## 2. CRIME ANALYSIS PROCEDURE

Usually, the crime analysis tasks can be a tedious process for the police or the investigation team to work with. The criminals when leaving the crime scene does leave some traces which can be used as a clue to identify the criminals. The crime sequence and the patterns which several criminals follow when committing a crime make it easy for analyzing the crime. This process includes several procedures to be followed in order to identify the criminals and getting more information based only on the clues or information given by the local people. The criminal can be analyzed based on the information from the crime scene which is tested against the previous crime patterns and judging by the method which is implied to test and proceed with the information that can affect the prediction results. The prediction can be further made useful for detecting the crimes in advance or by adding more cops to the sensitive areas which are identified by the system. The police stations can put up special force when there are chances for crime ahead of time. This type of the system will ensure there are peace and prosperity among the citizens.

The crime analysis can be performed procedure which is similar to figure Fig.1 which specifies each module which is used for machine learning to predict the crime or form group of clusters of criminals according to crime records. The criminals can hold certain properties and their crime characteristics and crime careers may vary from one criminal to another. Such a type of information can be taken as the input dataset. The input dataset is given to a pre-processor which performs the preprocessing based on the requirements. Once the pre processing is completed the features or attributes from those information are extracted which may be in the form of text content from emails, the crime factors for a day, criminal characteristics, geo-location of the criminal, etc., The pre processed result is further given to the classification algorithm or the clustering algorithm based on the requirements. The requirements may be anything from selecting the crime prone areas to predicting the criminal based on the previous crime records.

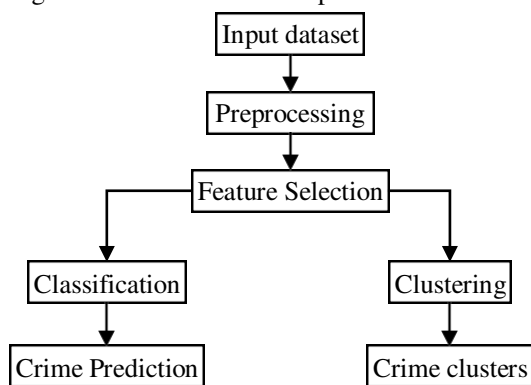


Fig.1. Crime Prediction and Crime clustering based on the input dataset

The classification algorithm works in a supervised learning manner in which the training and testing phase is required in order to train the classifier to identify the new unknown crime record. This is known as prediction. Whereas the clustering algorithm works in an un-supervised learning manner which automatically

separates the crime records based on the number of groups to be created. The groups created in such a manner are known as clusters. Such a type of design can be a general template for applying crime prediction and crime analysis based on data mining algorithms.

## 3. CRIMINAL ANALYSIS METHODS

### 3.1 TEXT, CONTENT AND NLP-BASED METHODS

Sharma [1] proposed a concept which depicts zero crime in the society. For detecting the suspicious criminal activities, he has concentrated on the importance of data mining technology and designed a proactive application for that purpose. In his paper, he proposed a tool which applies an enhanced Decision Tree Algorithm to detect the suspicious e-mails about the criminal activities. An improved ID3 Algorithm with an enhanced feature selection method and attribute-importance factor is applied to produce a better and faster Decision Tree based on the information entropy which is explicitly derived from a series of training data sets from several classes. He proposed a new algorithm which is a combination of Advanced ID3 classification algorithm and enhanced feature selection method for the better efficiency of the algorithm.

Hamdy et al. [8] described an approach based on the people's interaction with social networks and mobile usage such as location markers and call logs. Their work also introduced a model for detecting suspicious behavior based on social network feeds and it not only describes a new method using the social interaction of people but, their work proposes a new system to help crime analysis create faster and precise decisions. The suspicious movement of the entity can be determined using the sequence of inference rules. Their constructed model is able to predict and characterize human behavior from reality data sources

### 3.2 CRIME PATTERNS AND EVIDENCE-BASED METHODS

Bogahawatte and Adikari [2] proposed an approach in which they highlighted the usage of data mining techniques, clustering and classification for effective investigation of crimes and criminal identification by developing a system named Intelligent Crime Investigation System (ICSIS) that could identify a criminal based up on the evidence collected from the crime location. They used clustering to identify the crime patterns which are used to commit crimes knowing the fact that each crime has certain patterns. The database is trained with a supervised learning algorithm, Naïve Bayes to predict possible suspects from the criminal records. His approach includes developing a multi-agent for crime pattern identification. There are agents for the place, time, role trademark and substance of criminals which separates the role of the criminals in components. The system is a multi-agent system and made with managed Java Beans. It makes it easy to encapsulate the requested entities in the work into objects and returns it to the bean for exposing properties. Classifying the criminals/ suspects is based on the Naïve Bayes classifier for identifying most possible suspects from crime data. Clustering the criminals is based on the model to help to identify patterns of committing crimes.

Agarwal et al. [3] used the rapid miner tool for analyzing the crime rates and anticipation of crime rate using different data

mining techniques. Their work done is for crime analysis using the K-Means Clustering algorithm. The main objective of their crime analysis work is to extract the crime patterns, predict the crime based on the spatial distribution of existing data and detection of crime. Their analysis includes the tracking homicide crime rates from one year to the next

Kiani et al. [4] performed a crime analysis work based on the clustering and classification techniques. Their work includes the extraction of crime patterns by crime analysis based on available criminal information, prediction of crimes based on the spatial distribution of existing data and crime recognition. They proposed a model in which the analysis and prediction of crimes are done through the optimization of outlier detection operator parameters which is performed through the Genetic Algorithm. The features are weighted in this model and the low-value features were deleted through selecting a suitable threshold. After which the clusters are clustered by the k-means clustering algorithm for classification of crime dataset.

Satyadevan et al. [5] has done a work which will display high probability for crime occurrence and can visualize crime prone areas. Instead of just focusing on the crime occurrences, they are focusing mainly on the crime factors of each day. They used the Naïve Bayes, Logistic Regression and SVM classifiers for classification of crime patterns and crime factors of each day. Their method consists of a pattern identification phase which can identify the trends and patterns in crime using the Apriori Algorithm. The prediction of crime spots is done with the help of Decision Tree algorithm which will detect the crime possible areas and their patterns.

Bruin et al. [7] proposed a technique which is used to determine the clustering of criminals based on the criminal careers. The criminal profile per offense per year is extracted from the database and a profile distance is calculated. After that, the distance matrix in profile per year is created. The distance matrix including the frequency value is made to form clusters by using naïve clustering algorithm. They made a criminal profile which is established in a way of representing the crime profile of an offender for a single year. With this information, the large group of criminals is easily analyzed and they predicted the future behavior of individual suspects. It will be useful for establishing the clear picture on different existing types of criminal careers. They tested the tool on actual Dutch National Criminal Record Database for extracting the factors for identifying the criminal careers of a person.

### 3.3 SPATIAL AND GEO-LOCATION BASED METHODS

Huang et al. [6] focused on a different approach for criminal activity prediction based on mining location based Social Network interactions. By using these interactions, they can collect information using the geographical interactions and data collections from the people. They devised a working procedure in which a series of features are categorized from the Foursquare and Gowalla used in the San Francisco Bay area. The crime patterns and the crime occurrences are tracked with the geographical features which are extracted from the map and they are analyzed to detect the urban areas with high crime activities. Their work aims at exploiting the location-based social network data to investigate the criminal activities in urban areas. By using the

Haversine formula the distance between the two points i.e. the crime location and venue location is calculated and shown in the Google Maps API and OpenStreetMap.

Chen [19] have presented a general framework for crime data mining that draws on experience gained with the Coplink project with the researchers at Arizona and their work mainly focuses on showing the relationships between crime types and the link between the criminal organizations. They used a concept space approach which will extract criminal from the incident summaries.

Yu [20] have discussed the preliminary results of a crime forecasting model developed in collaboration with the police department of a United States city in the Northeast. Their approach is to architect datasets from original crime records. The datasets contain aggregated counts of crime and crime-related events categorized by the police department. The location and time of these events is embedded in the data. Additional spatial and temporal features are harvested from the raw data set. Second, an ensemble of data mining classification techniques is employed to perform the crime forecasting. Then they analyzed a variety of classification methods to determine which is best for predicting crime “hotspots”. They even investigated classification on increase or emergence. Last, they have proposed the best forecasting approach which is aimed at achieving the most stable outcomes.

Rizwan et al. [22] have performed classification of crime dataset to predict Crime Category for different states of the United States of America. The crime dataset that they used in this research is real in nature. That is, it was collected from socio-economic data from 1990 US Census, law enforcement data from the 1990 US LEMAS survey, and crime data from the 1995 FBI UCR. Their work compared the two different classification algorithms namely, Naïve Bayesian and Decision Tree for predicting Crime Category for different states in USA. The results from their experiment showed that, Decision Tree algorithm out performed Naïve Bayesian algorithm and achieved 83.9519% Accuracy in predicting Crime Category for different states of USA.

Donald [24] have proposed a system for Crime Analysis which was named by them as The Regional Crime Analysis Program (ReCAP) system. It was designed by them as a computer application designed to aid local police forces (e.g. University of Virginia (UVA), City of Charlottesville, and Albemarle County) in the analysis and prevention of crime. ReCAP works in cooperation with the Pistol 2000 records management system, which aggregated and housed all of the crime information from a region. Their research and development was primarily focused on the individual components of the system which includes a database, geographic information system (GIS), and data mining tools which consisted of data mining algorithms which produced spatial mining results over the crime hotspots. Their system consists of the seamless integration of all the components in the system.

### 3.4 PRISONER BASED METHODS

Sheehy et al. [10] came up with a research idea which was geared towards the treatment of the mentally ill people inside the prison. According to their work, the mentally ill criminals are identified using their Social Security Number (SSN) with all the criminal personal records and their crime career records attached. As the outcome, the Criminals are classified into “high”, “medium” and “low” levels of recidivism risk potential according

to their mental health. Their objective was to describe and classify the criminals into a misdemeanor and a felony which can be referred and not referred based on the mental health of the criminals. Their ill activities are monitored and data collection is continuous. By these, the criminals can be separated from other criminals who are hazardous and those who can cause damage to other inmates along with them. Further, their study also involves the classification of the mental health of the criminals into two categories i.e. “referred” and “not-referred”. This helps the guards to identify the prisoners who are referred for the mental health check-up. The research work they had undergone will provide a summary of the inmates who are seriously mentally ill and those who are to be separated from the other inmates.

### 3.5 COMMUNICATION BASED METHODS

Taha et al. [9] has developed a forensic investigation tool for identifying the influential members who create an impact in a criminal organization. The immediate leaders can also be identified in a criminal organization. Removing these influential members can weaken the strength of the criminal organization. Their work is based on this methodology. They proposed a new work which is known as SIIMCO which first constructs the graph representing the criminal group or organization as a network from either mobile communication data of the criminal organization or based on the crime records. The system works on the basis of the created networks. These networks represent the criminal organization or crime incident reports. The vertex represents the individual criminals and the link represents the relationships or communication link between those two criminals. They employed certain formulas that quantify the degree of influence/ importance of each vertex in the network relative to all other vertices i.e. criminals in the graph. Based on this their system identifies the immediate leaders with the weighted graph which connects the criminals and identify them for further processing.

## 4. QUALITATIVE ANALYSIS OF CRIME ANALYSIS AND PREDICTION APPROACHES

The prediction can be made based on the Textual information or the Geospatial information or even the prisoner records which were manually recorded. By using the real open data such as internet, social feeds and messages the researcher can use the text processing or NLP techniques to mine information from the data

and categorize the e-mails, messages or posts into a suspicious or a non-suspicious record [1]. Whereas in the Spatial mining area, the extraction of features from SNAP Gowalla dataset, DataSF criminal dataset up to February 2015 provides the way to plot the crime occurrences on the Google Map which is interpreted easily. The communication based methods describe the identification of the leaders in a criminal organization may be a tedious process. Kamal Taha et al. [9] produced an approach through the phone calls and other communication data such as call logs and records, the influential members on a crime organization can be tracked. Kevin Sheehy, Thomas Rehbreger, Andrew O’Shea, William Hammond, Charlotte Blais, Michael Smith K., Preston White, Jr., Neal Goodloe [10] introduced an approach to categorize and identify the mentally ill prisoners among the prisoners and keep them separate from other prisoners to avoid conflict and injuries between them. Even though there are many methods for analyzing the crimes, this paper concludes many results based on the qualitative analysis. When considering the Text/NLP based methods, Hamdy et al. [8] overcame the defects from the work of Sharma [1] based on many factors such as implementation of preprocessing for the data and extraction of relevant features. Both the paper labels the outcome based on suspicious activity. Mugdha Sharma [1] used an enhanced ID3 algorithm whereas the work produced by Ehab Hamdy, Ammar adl, Aboul Ella Hassanien, Osman Hegazy and Tai-Hoon Kim. [8] does not specify the classification algorithm. The weakness of this paper is mostly about not giving the clear view of the pre processing and classification algorithm. When considering crime patterns and evidence based methods, there are clustering and classification based papers. Bogawatte and Adikari [2] concentrated on using the Naïve Bayes for finding out most possible suspect. Jyoti Agarwal et al. [3] on the other hand focused on crime analysis by implementing the K-Means clustering algorithm on crime dataset using rapid miner tool and the author had performed the crime analysis by considering the homicide crimes and plotting it with respect to year. Kiani et al. [4] concentrated on using the Genetic Algorithm to optimize the distance operator parameter of the decision tree using GINI index. The clustering of the criminal careers has been effectively done in the work [7]. Whereas, Shiju Satyadevan, et al. [5] have performed a comparison of the Naïve Bayes, SVM, Logistic Regression and Decision Tree. This paper presents the crime prone regions and represented as heatmaps which indicate the level of heat. When considering the Spatial and Geolocation based methods, all these methods are analysed based on qualitative manner and the analysis information is described in the below mentioned table Table.1.

Table.1. Qualitative Analysis of Crime Analysis and Prediction

| METHOD                  | INPUT                            | DATASET USED                                                      | PRE PROCESSING                          | FEATURE EXTRACTION                                                                              | CLASSIFICATION/ CLUSTERING                            | STRENGTH                                                                                  | WEAKNESS                                  | OUTCOME                                                             |
|-------------------------|----------------------------------|-------------------------------------------------------------------|-----------------------------------------|-------------------------------------------------------------------------------------------------|-------------------------------------------------------|-------------------------------------------------------------------------------------------|-------------------------------------------|---------------------------------------------------------------------|
| Text/ NLP-based methods | [1] E-mail messages              | Real and open emails sent by terrorists and some are dummy emails | Nil                                     | Selection of a subset of the original text containing “kill”, “death”, “bomb”, “guns”, “blasts” | Enhanced ID3 Decision Tree algorithm                  | Introducing attribute importance as a factor before information gain in the decision tree | Nil                                       | Labeling email as Suspicious, Non-suspicious, and May be suspicious |
|                         | [8] Crime history, age, previous | Device sensors, Security                                          | Structuring collective data into {Time, | Similarity matching for sensory images                                                          | A trained classification model is used to predict the | Consideration of location feeds and                                                       | Not giving a clear view of the processing | Suspicious behavior to three levels                                 |

|                                           |     |                                                                                                                                                                       |                                                                                      |                                                        |                                                                                                                                                        |                                                                                                                                         |                                                                                                                                                    |                                                                                                                                   |                                                                                                                              |
|-------------------------------------------|-----|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------|--------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------|
|                                           |     | arrests, Modus Operandi, countries visited, place of birth, Average use of ATM, Types of crimes, Entrance with respect to Time of Day, Crime areas, Victims' mistakes | camera information, Messages, Audio feeds, Social network posts and messages         | Final Movement, Frequency rate, Video, Images, Audio } | using sliding window. Text semantic Analysis of the text information performed using Lexical processing, Natural Language Processing (NLP).            | similarity of a given input to the suspicious item or location.                                                                         | mobile usage information                                                                                                                           | and comparison of criminal behavior.                                                                                              | such as "High", "Medium" and "Low"                                                                                           |
| Crime patterns and Evidence-based methods | [2] | Crime evidences including many attributes like crime scene, day, month, offense, resources used, time, role in crime, transportation etc.,                            | Colombo crime and criminal records                                                   | Nil                                                    | Extraction of evidence                                                                                                                                 | Clustering based model to identify patterns of committing crimes.<br><br>Naïve Bayes classifier applied to find most possible suspect   | Uses Naïve Bayes so this can be even suitable for small datasets.                                                                                  | No clear view of clustering method and Prisoner verification                                                                      | Finding Categories as robbery, burglary, and theft<br>Classifying person as "suspect" and after judgment "criminal"          |
|                                           | [3] | Homicide crimes and their occurrences                                                                                                                                 | Crime dataset for crime analysis by polices in England and Wales from 1990 – 2011-12 | Nil                                                    | Extraction of crime patterns based on the available crime and criminal data                                                                            | K-means clustering algorithm                                                                                                            | Produces year wise clusters of homicide crimes committed                                                                                           | Concentration is only on clustering of homicide crimes                                                                            | Year and analysis of variation in clusters formed                                                                            |
|                                           | [4] | Burglary, Robbery, and Homicide                                                                                                                                       | Crime dataset for crime analysis by polices in England and Wales from 1990 – 2011    | Nil                                                    | Filtering of dataset, Outlier detection using distance operator (k-NN), Genetic Algorithm used for optimizing of outlier detection operator parameters | Classification was done using Decision Tree using GINI index and the testing and training done using Sample Stratified                  | Use of GA to optimize the distance operator parameters in Clustering and Predict the cluster's members based on classification using Decision Tree | The number of clusters in the clustering process needs to be optimized and further optimization of the technique needs to be done | The results for the optimized and non-optimized parameters were compared to show the difference in quality and effectiveness |
|                                           | [5] | location, date, type of crime data extracted from Websites, Blogs, Social Media, RSS Feeds                                                                            | Websites, Spatial Information, and date about crimes                                 | Nil                                                    | Extraction of the following crime data related to "vandalism", "murder", "robbery", "burglary", "sex abuse", "gang rape", "arson", "armed robbery",    | Naïve Bayes, SVM, Logistic regression<br><br>Crime prediction was done using decision tree which is done using sample police complaints | Comparison of Naïve Bayes with SVM. Decision Tree is easy to interpret and understand for crime spot identification.                               | Not predicting the time in which the crime is happening.                                                                          | The crime-prone areas (regions) are graphically represented using a heat map which indicates the level of crimes             |

|                                        |      |                                                                                                                                                                                            |                                                                                                   |                                                                                                                                                |                                                                                                                                         |                                                                                                                                                                            |                                                                                                             |                                                                                 |                                                                                                                                                                                                                                                                                                  |
|----------------------------------------|------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|                                        |      |                                                                                                                                                                                            |                                                                                                   |                                                                                                                                                | “highway robbery”, “snatching”                                                                                                          |                                                                                                                                                                            |                                                                                                             |                                                                                 |                                                                                                                                                                                                                                                                                                  |
|                                        | [7]  | Crime database and criminal information                                                                                                                                                    | National Crime Record Database                                                                    | Nil                                                                                                                                            | Crime nature, frequency, duration, severity                                                                                             | Crime profile of offender for single year is determined for comparison and he                                                                                              | Development of new distance measures with combination of profile distance with crime frequency of criminals | The runtime of the chosen approach is not optimal                               | Clustering of criminal careers based on the nature. One time criminal, severe criminals and minor career criminals                                                                                                                                                                               |
| Spatial and Geo-location based methods | [6]  | Geo-location and Crime Type                                                                                                                                                                | SNAP Gowalla dataset, DataSF criminal dataset up to February 2015                                 | Extraction of crime type like Assault, Robbery, Theft, Vandalism, Drug                                                                         | Geographical features, Popularity, Location category, Neighbor entropy, Social Tightness density, crime location, venue from Foursquare | Random Forest(RF), Linear Regression (LR) and Support Vector Machine (SVM)                                                                                                 | Random Split method utilized with 80% for training and 20% for testing in classification                    | Nil                                                                             | Crime Areas plotted using Google Map API and OpenStreetMap in San Francisco Bay area and Criminal pattern discovery according to the context of user activity and location-based social networks. Predict crime frequency and find which crime is to be more difficult or easier to be predicted |
| Communication based methods            | [9]  | Flow of communications/information links between two criminals (e.g., phone call records, messages, etc.), names of criminals/suspects, the type of crime, location and date of the crime. | Real-world communication records (DBLP, Enron email dataset, Nodobo mobile phone records dataset) | Creating the graph based on the data and then assigning weight to a vertex based on its number of communication attempts in the criminal graph | The immediate leaders of lower-level criminals and the lower-level criminals themselves are extracted.                                  | Evaluation of the accuracy of the three systems by measuring their Recall, Precision, and Euclidean Distance.                                                              | Evaluated SIIMCO by comparing it experimentally with CrimeNet Explorer and LogAnalysis                      | Nil                                                                             | System can identify the influential members of a criminal organization and the immediate leaders of a given list of lower-level criminals                                                                                                                                                        |
| Prisoner based methods                 | [10] | The Social Security Number (SSN) with all the criminal personal and crime career records.                                                                                                  | Albemarle-Charlottesville Regional Jail (ACRJ), Jefferson Area Community Corrections (JACC) and   | A combination which includes the Social Security Number (SSN) and date was used to link the databases together.                                | age, criminal history, employment history, crime type := “assault”, “larceny”, “supervision violations”, “narcotics                     | Offenders are classified into three classes namely “high”, “medium”, and “low” as levels of recidivism risk potential. Further, the mental health status of the inmates is | Analysis for the identification of the mentally ill felony.                                                 | Statistical classification of criminals missing. Could have taken more features | “Referred” individuals can be made to have a longer stay in jail longer than “not-referred” individuals.                                                                                                                                                                                         |

|  |  |  |                                               |  |                                                                       |                                                                       |  |  |  |
|--|--|--|-----------------------------------------------|--|-----------------------------------------------------------------------|-----------------------------------------------------------------------|--|--|--|
|  |  |  | Region Ten<br>Community<br>Services<br>Board. |  | charges”, “traffic<br>violations”,<br>“driving while<br>intoxicated”, | categorized into two<br>categories “referred.”<br>and “not-referred.” |  |  |  |
|--|--|--|-----------------------------------------------|--|-----------------------------------------------------------------------|-----------------------------------------------------------------------|--|--|--|

## 5. QUANTITATIVE ANALYSIS OF CRIME ANALYSIS AND PREDICTION APPROACHES

For performing the quantitative analysis of the methods taken, the performance metric value needed to be computed and they are to be compared with the other. Hence, for performing the calculations of the performance metric there are a few formulas which can be utilized for achieving the performance value from the dataset. The formulae for the calculation of the performance metrics are given below in Table.2.

Table.2. Metrics and their formula

Although many papers were studied in the literature review all the papers were irrelevant to the crime prediction and criminal analysis domain. Hence a few papers in the crime analysis and prediction domain has been taken and their results have been reproduced as given originally in the reference papers. The below given table Table 3 provides quantitative analysis of the three tools and the Decision Tree algorithm which is supported with the Genetic Algorithm for optimization of the parameters. When the parameters are optimized, the classification accuracy of the Decision Tree is increased a bit further. This shows that although the Decision tree performs well, when it used with the Genetic Algorithm for optimization of the decision tree parameters, the results shown show significant improvement in the accuracy and further more the tools given below have the metric value, which is purely based on the dataset and the records and the performance values are taken as it is in the reference paper. The quantitative analysis produced results which show the increase in the accuracy level of classification because of using the GA to optimize the parameters. This occurs because of the ability of the GA to learn the optimal values and then it is applied to set the parameter to optimal value when performing calculation. Also, the Precision, Recall and F-value varies from the dataset and the system. This shows the SIIMCO performing well when defined in terms of the metrics.

Table.3. Quantitative Analysis of Crime Analysis & Prediction

|     | NAME OF THE METHOD/ SYSTEM                                             | PERFORMANCE METRIC     | PERFORMANCE VALUE       |        |
|-----|------------------------------------------------------------------------|------------------------|-------------------------|--------|
| [4] | Decision Tree classification with GA for optimizing the the parameters | Accuracy of Prediction | Optimized parameter     | 91.64% |
|     |                                                                        |                        | Non-Optimized parameter | 85.74% |
|     |                                                                        | Classification Error   | Optimized parameter     | 8.36%  |
|     |                                                                        |                        | Non-Optimized parameter | 13.26% |
|     |                                                                        | Fitness Function       | Optimized parameter     | 72.28% |
|     |                                                                        |                        | Non-Optimized parameter | 72.48% |
| [9] | SIIMCO                                                                 | Recall                 | 0.62                    |        |
|     |                                                                        | Precision              | 0.56                    |        |
|     |                                                                        | F-Value                | 0.59                    |        |
|     | CrimeNet Explorer                                                      | Recall                 | 0.36                    |        |
|     |                                                                        | Precision              | 0.41                    |        |
|     |                                                                        | F-value                | 0.38                    |        |
|     | Log Analysis                                                           | Recall                 | 0.53                    |        |
|     |                                                                        | Precision              | 0.51                    |        |
|     |                                                                        | F-value                | 0.52                    |        |

## 6. CONCLUSION

In this paper, we have studied some known approaches for crime analysis and prediction concerned with data mining. Although many papers have been studied, only those papers with background in the crime prediction and criminal identification papers are compared with a theoretical study. Each paper has their own advantages and disadvantages. Each paper has its own individual approach for solving the crimes and criminal prediction. This is a theoretical study for several methods in identification of crime and criminals which includes Text/ NLP based methods, crime patterns and crime evidence based methods, spatial and geo location based methods, communication based methods and finally Prisoner based methods. The data mining techniques studied from this survey can be applied for identifying the criminals in the society and also for providing a better future to live in.

## REFERENCES

- [1] Mugdha Sharma, "Z-Crime: A Data Mining Tool for the Detection of Suspicious Criminal Activities based on the Decision Tree", *International Conference on Data Mining and Intelligent Computing*, pp. 1-6, 2014.
- [2] Kaumalee Bogahawatte and Shalinda Adikari, "Intelligent Criminal Identification System", *Proceedings of 8<sup>th</sup> IEEE International Conference on Computer Science and Education*, pp. 633-638, 2013.
- [3] Jyoti Agarwal, Renuka Nagpal and Rajni Sehgal, "Crime Analysis using K-Means Clustering", *International Journal of Computer Applications*, Vol. 83, No. 4, pp. 1-4, 2013.
- [4] Rasoul Kiani, Siamak Mahdavi and Amin Keshavarzi, "Analysis and Prediction of Crimes by Clustering and Classification", *International Journal of Advanced Research in Artificial Intelligence*, Vol. 4, No. 8, pp. 11-17, 2015.
- [5] Shiju Sathyadevan, M.S. Devan and S. Surya Gangadharan, "Crime Analysis and Prediction using Data Mining", *Proceedings of IEEE 1<sup>st</sup> International Conference on Networks and Soft Computing*, pp. 406-412, 2014.
- [6] Yu-Yueh Huang, Cheng-Te Li and Shyh-Kang Jeng, "Mining Location-based Social Networks for Criminal Activity Prediction", *Proceedings of 24<sup>th</sup> IEEE International Conference on Wireless and Optical Communication*, pp. 185-190, 2015.
- [7] Jeroen S. De Bruin, Tim K. Cocx, Walter A. Kusters, Jeroen F. J. Laros and Joost N. Kok, "Data Mining Approaches to Criminal Career Analysis", *Proceedings of 6<sup>th</sup> IEEE International Conference on Data Mining*, pp. 1-7, 2006.
- [8] Ehab Hamdy, Ammar Adl, Aboul Ella Hassanien, Osman Hegazy and Tai-Hoon Kim, "Criminal Act Detection and Identification Model", *Proceedings of 7<sup>th</sup> International Conference on Advanced Communication and Networking*, pp. 79-83, 2015.
- [9] Kamal Taha and Paul D. Yoo, "SIIMCO: A Forensic Investigation Tool for Identifying the Influential Members of a Criminal Organization", *IEEE Transactions on Information Forensics and Security*, Vol. 11, No. 4, pp. 811-822, 2016.
- [10] Kevin Sheehy et al., "Evidence-based Analysis of Mentally 111 Individuals in the Criminal Justice System", *Proceedings of IEEE Systems and Information Engineering Design Symposium*, pp. 250-254, 2016.
- [11] David J. Hand, Heikki Mannila and Padhraic Smyth, "*Principles of Data Mining*", MIT Press, 2001.
- [12] 10 Emerging Technologies That Will Change Your World, Available at: [http://www.rle.mit.edu/thz/documents/10\\_emerging\\_tech.pdf](http://www.rle.mit.edu/thz/documents/10_emerging_tech.pdf).
- [13] U. Fayyad, G. Piatetsky-Shapiro and P. Smyth, "The KDD Process for Extracting Useful Knowledge from Volumes of Data", *Communications of the ACM*, Vol. 39, No. 11, pp. 27-34, 1996.
- [14] Illhoi Yoo, Patricia Alafaireet, Miroslav Marinov, Keila Pena-Hernandez, Rajitha Gopidi, Jia-Fu Chang and Lei Hua, "Data Mining in Healthcare and Biomedicine: A Survey of the Literature", *Journal of Medical Systems*, Vol. 36, No. 4, pp. 2431-2448, 2011.
- [15] Shyam Varan Nath, "Crime Pattern Detection using Data Mining", *Proceedings of IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology Workshops*, pp. 1-4, 2006.

## REACH CENTROID ALGORITHM IDENTIFY THE WIRELESS LOCALIZATION IN WIRELESS SENSOR NETWORKS

\*Dr. G. Silambarasan, \*\* Dr. V. Chandrasekar,

\*Assistant Professor, Dept. of Computer Science and Engineering,  
The Kavery College of Engineering, Salem, Tamilnadu, India,

\*\*Associate Professor, Dept. of Computer Science and Engineering,

Malla Reddy College of Engineering and Technology, Secunderabad, Telangana State, India,

\*\*[drchandru86@gmail.com](mailto:drchandru86@gmail.com), \* [gssilambarasan@gmail.com](mailto:gssilambarasan@gmail.com)

### ABSTRACT

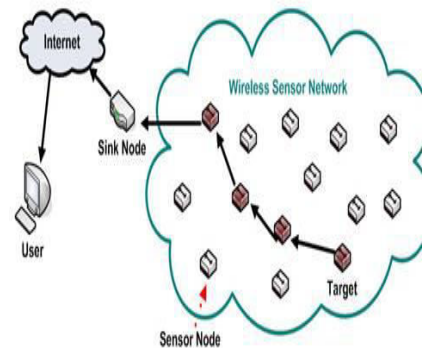
Position and accurate is the important for the localization in wireless sensor networks. This technique it the low cost technique in wireless sensor networks. To manage the cost and, few location aware node, it is called as anchors in wireless sensor networks environment. Other anchors of free node estimate the own positions. If you have change the any position of the node or modify from centroid localization algorithm called as reach centroid localization algorithm. This method mainly focuses on anchor nodes position validation methods. Every anchor node with the near free node search and validate the position of actual or near node will received the signal strength. This work reduce the multipath effects of radio waves, particularly enclosed the environment. Our proposed work is localization become more significant, particularly indoor environments, the result show a significant improvement in localization accuracy when compare with the original centroid localization algorithm.

**Keyword:** WSN, Centroid, Localization

### INTRODUCTION

Wireless sensor networks (WSN), are similar to wireless ad hoc networks in the sense that they rely on wireless connectivity and spontaneous formation of networks so that sensor data can be transported wirelessly. Sometimes they are called dust networks, referring to minute sensors as small as dust. Smart dust [1][2][3] is a U C Berkeley project sponsored by DARPA. Networking. Is one of the early companies that produced wireless sensor network products. WSNs are spatially distributed autonomous sensors to monitor physical or environmental conditions, such as temperature, sound, pressure, etc. and to cooperatively pass their data through the network to other locations.[4] The more modern

networks are bi-directional, also enabling control of sensor activity. The development of wireless sensor networks was motivated by military applications such as battlefield surveillance; today such networks are used in many industrial and consumer applications, such as industrial process monitoring and control, machine health monitoring, and so on.



**Fig1:** Wireless sensor Networks

In this section, we are going to derive a novel node localization method for 3D WSNs.

### ALGORITHM DEVELOPMENT

Before proceeding, we review the centroid algorithm. 2.1 Review of Centroid Algorithm Bulusu and Heidemann [2] have proposed the centroid localization algorithm, which is a range-free, proximity-based, coarse-grained localization algorithm. The algorithm implementation contains three core steps. First, all anchors send their positions to all sensor nodes within their transmission range. Each unknown node listens for a fixed time period  $t$  and collects all the beacon signals it receives from various reference points. Second, all unknown sensor nodes calculate their own positions by a centroid determination from all  $n$  positions of the anchors in range. The centroid localization algorithm, which uses anchor nodes (reference nodes), containing location information  $(x_i, y_i)$ , to estimate node position. After receiving these beacons, a node estimates its location using the following centroid formula:

## RANGE-FREE ALGORITHMS

The nature of radio wave propagation is such that the attenuation of radio signal increases as distance between the transmitter and receiver increases. Radio propagation models [11] in various environments are well documented and have often focused on estimating the average Received Signal Strength (RSS) at a particular distance of the transmitter, as well as the variability of the signal strength in close spatial proximity to the location.

The Departure Test Definition shows that when incrementally increasing the distance between anchors and receiving nodes, the RSS monolithically decreases with distance [8]. However, there are instances where there are burst in signal strength due to disturbance effects, such as, reflection leading to signal amplification or sudden loss of signal, due to absorption as a result of environmental conditions. Nevertheless, the test does not make

any assumption about the correlation between absolute distance and signal strength.

Relevant to our research are some Range-free algorithms such as DV-hop, Amorphous, APIT and CLA algorithms. This is because these algorithms operate on the same fundamental principle; they all attempt to select the anchors with the most significant characteristics for location estimation. These algorithms will be briefly discussed.

## DV-HOP AND AMORPHOUS ALGORITHMS

DV-Hop and Amorphous algorithms both use a form of distance vector exchange so that all the nodes within the network get estimated distances in hops to the anchors [12] as against the linear distance between the free node and the anchor. A node estimates its position by assuming the average distance of the closest anchor to it. It then uses the distance in hop count to estimate its position from at least two other anchors using the same distance average received from the anchor closest to it. After which, triangulation is performed to estimate the position of the free node. This procedure is appropriate for nodes with limited capabilities and lacks the ability to process the image of the entire network.

## REACH CENTROID LOCALIZATION ALGORITHM

The relationship between transmitted and received signal is not symmetric. Two nodes transmitting and receiving signals between each other are likely not to receive each other's signal at the same strength even if their power of transmission is the same. Many factors such as multipath effects account for this discrepancy. This situation or condition often leads to the erroneous localization of objects, particularly indoors. ReachCLA seeks to manage these

conditions towards a more reliable estimation of the position of objects.

The algorithm establishes an authentication or feedback process between the free node and the anchors within its reach. In a global WSN environment, a free node is likely to receive signals from anchors within its immediate locality. These anchors are the ones within its reach and they are the ones that will participate in the localization process. ReachCLA [9] has five main phases, they are summarized as follows:

The free node selects all the anchors that are within its reach.

The second phase is the feedback or handshake phase [6]. This is to establish how the anchors read each other. This process is key as it allows the anchors establish their proximity to one another, and by so doing, they are able to reasonably ascertain the true position of the free node.

Selection of the anchors with the strongest reach. This is done as each selected anchor by the free node compare the anchors within their reach to other anchors selected by the free node.

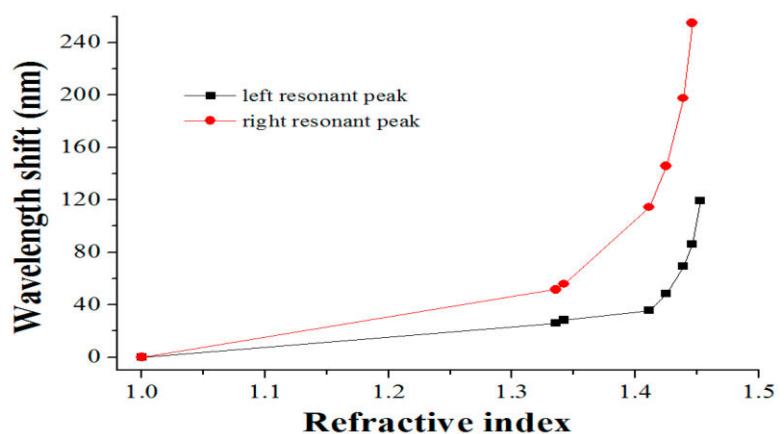
Selection of the anchors with the three highest links or reach. Selection is from the group of anchors within the range of the free nodes, these are the anchors that have more anchors in common when matched with the anchors within the reach of the free node.

There are probable instances where an anchor within the reach of a free node will not be able to read the transmitted signal from the free node in a two way communication. However, because the anchor in case is able to communicate with other anchors that are within the reach of the free node, it therefore assures the likelihood or authenticates that the free node

is within its own proximity. It is this authentication process that makes ReachCLA a unique and reliable algorithm for localization. An experimental tested was set up to analyze the performance of ReachCLA alongside CLA and APIT algorithms

## FINDINGS AND DISCUSSION

The simulations of ReachCLA, CLA and APIT algorithms respectively. The simulations show the different estimated positions of the sensor node using each algorithm. Presents the various estimated locations of the sensor node using ReachCLA, CLA and APIT on a single plot[13]. The various plots indicate that the fundamental principle of ReachCLA is similar to that of CLA and APIT as they all embrace the formation of triangles within areas of discernible radio signal, and use the triangles to estimate the location of the intended object. However, the main distinction lies in the process of selecting the most likely triangle or triangles to be used for localization. This is where ReachCLA, as indicated in Table 3, outperforms both CLA and APIT.



**Fig:** Comparison Graph for Left Peak & Right

Peak

## CONCLUSION

We present a new localization method that improves the basic centroid localization algorithm significantly. It is shown in the simulation results that the proposed algorithm can improve location accuracy than the conventional centroid localization algorithm. The performance of our proposed scheme has identified that it has potential application advantage of centroid algorithm [3].

## REFERENCES

- [1]. Sharma, R. and Malhotra, S. (2015) Approximate Point in Triangulation (Apit) Based Localization Algorithm in Wireless Sensor Network. International Journal for Innovative Research in Science and Technology, 2, 39-42. [Citation Time(s):1]
- [2]. He, T., Huang, C., Blum, B.M., Stankovic, J.A. and Abdelzaher, T. (2003) Range-Free Localization Schemes for Large Scale Sensor Networks. Proceedings of the 9th Annual International Conference on Mobile Computing and Networking, San Diego, 14-19 September 2003, 81-95.
- [3]. Dargie, W. and Poellabauer, C. (2010) Fundamentals of Wireless Sensor Networks: Theory and Practice. John Wiley & Sons, Hoboken.
- [4]. Kumar, A., Kumar, V. and Kapoor, V. (2011) Range Free Localization Schemes for Wireless Sensor Networks. 10th WSEAS International Conference on Software Engineering, Cambridge, 22-24 February 2011, 101-106.
- [5]. Karl, H. and Willig, A. (2007) Protocols and Architectures for Wireless Sensor Networks. John Wiley & Sons, Hoboken.
- [6]. Kim, S.-Y. and Kwon, O.-H. (2005) Location Estimation Based on Edge Weights in Wireless Sensor Networks. The Journal of Korean Institute of Communications and Information Sciences, 30, 938-948
- [7]. Chong, C. and Kumar, S.P. (2003) Sensor Networks: Evolution, Opportunities, and Challenges. Proceedings of the IEEE, 91, 1247-1256.
- [8]. Li, X.-Y., Wan, P.-J. and Frieder, O. (2003) Coverage in Wireless Ad Hoc Sensor Networks. IEEE Transactions on Computers, 52, 753-763.
- [9]. Lazos, L. and Poovendran, R. (2004) Serloc: Secure Range-Independent Localization for Wireless Sensor Networks. Proceedings of the 3rd ACM Workshop on Wireless Security, Philadelphia, 01 October 2004, 21-30.
- [10]. Nagpal, R., Shrobe, H. and Bachrach, J. (2003) Organizing a Global Coordinate System from Local Information on an Ad Hoc Sensor Network. In: Zhao, F. and Guibas, L., Eds., Information Processing in Sensor Networks, Springer, Berlin, 333-348.
- [11]. Priyantha, N.B., Chakraborty, A. and Balakrishnan, H. (2000) The Cricket Location-Support System. Proceedings of the 6th Annual International Conference on Mobile Computing and Networking, Boston, 06-11 August 2000, 32-43.
- [12]. Bulusu, N., Heidemann, J. and Estrin, D. (2000) Gps-Less Low-Cost Outdoor Localization for Very Small Devices. IEEE Personal Communications, 7, 28-34.
- [13]. Long, S., Wang, F., Duan, W. and Ren, F. (2004) Range-Free Self-Localization Mechanism and Algorithm for Wireless Sensor Networks. Computer Engineering and Applications, 23, 39.
- [14]. He, T., Huang, C., Blum, B.M., Stankovic, J.A. and Abdelzaher, T.F. (2005) Range-Free Localization and Its Impact on

Large Scale Sensor Networks. ACM  
Transactions on Embedded Computing Systems,  
4, 877-906

[15]. Rappaport, T.S. (1996) Wireless  
Communications: Principles and Practice. Vol.  
2, Prentice Hall, Upper Saddle River.

# Cybercrime: A threat to Network Security

**Ch.Vijaya Kumari,**

Associate Professor & HOD

Department of CSE,

Malla Reddy College of Engineering,

hodcse@mrce.in

**Ch.Vengaiah**

Assistant Professor

Department of CSE,

Malla Reddy College of Engineering,

vengaiah19@gmail.com

## ABSTRACT

This research paper discusses the issue of cyber crime in detail, including the types, methods and effects of cyber crimes on a network. In addition to this, the study explores network security in a holistic context, critically reviewing the effect and role of network security in reducing attacks in information systems that are connected to the internet. As, all this adversely affects the efficiency of information security of any kind of security that exists and is used in information systems. Since hackers and other offenders in the virtual world are trying to get the most reliable secret information at minimal cost through viruses and other forms of malicious soft-wares, then the problem of information security - the desire to confuse the attacker: Service information security provides him with incorrect information; the protection of computer information is trying to maximally isolate the database from outside tampering. In other words, the Internet is a large computer network, or a chain of computers that are connected together. This connectivity allows individuals to connect to countless other computers to gather and transmit information, messages, and data. Unfortunately, this connectivity also allows criminals to communicate with other criminals and with their victims.

### **Keywords**

*Security, Network Security, Computer, Privacy, Cyber Crimes.*

## 1. INTRODUCTION

The advent of computers and the expansion of the Internet made likely the accomplishment of large improvement in research, surgery, expertise, and communication. Unfortunately, computers and the Internet have

furthermore supplied a new natural environment for crime. As Janet Reno, U.S. advocate general throughout the Clinton management, put it, "While the Internet and other data technologies are conveying tremendous advantages to humanity, they furthermore supply new possibilities for lawless individual behavior" (Dasey, Pp. 5-19).

Cybercrime is roughly characterized as committing a misdeed through the use of a computer or the Internet. The Internet has been characterized as "collectively the myriad of computer and telecommunications amenities, encompassing gear and functioning programs, which comprise the interconnected worldwide mesh of systems that provide work the Transmission Control Protocol/Internet Protocol, or any predecessor or successor protocols to such protocol, to broadcast data of all types by cable or radio"

(Internet Tax Freedom Act of 1998: 112 Stat. 2681-719).

In other phrases, the Internet is a large computer mesh, or a string of connections of computers that are attached together. This connectivity permits persons to attach to countless other computers to accumulate and convey data, notes, and data. Unfortunately, this connectivity furthermore permits lawless individuals to broadcast with other lawless individuals and with their victims. Although no unanimously acknowledged delineation of cybercrime lives, a distinction is often made between a customary misdeed that is perpetrated through the use of a computer or the Internet and a misdeed that engages expressly aiming at computer technology (Richards, Pp. 21-54).

This paper provides an understanding of how network security protection can help a firm to keep its information

safe from potential losses. The research builds upon extensive research and literature related to network security and protection. The paper gives a comprehensive account of some most important security tools (like firewalls) which can help companies to secure their information networks from unauthorized use. A brief account of challenges faced in Network Security Management is also provided to identify the potential areas for research.

## 2. AIMS AND OBJECTIVES

To determine the impact of cybercrime on networks. To determine the advent of cyber-crime.

## 3. RESEARCH DESIGN

The search will base on secondary data accumulation. The data will be pressed out from various journals, articles and books. Secondary research depicts information assembled by literature, broadcast media, publications, and other nonhuman origins. This type of research does not necessitate human fields. The research accession used is qualitative and used the case study methodology. Qualitative research is practically more immanent than quantitative one.

Methods of research in order to achieve intended tasks various theoretical and empirical methods are invoked in the master thesis. Theoretical methods of research: Historical method is applied to provide knowledge about cyber-crimes and network security. Logical methods (generalization, induction, deduction) are invoked to generalize the used literature and to draw inferences.

## 4. MATERIALS AND METHODS

The measurements of choice for literature were relevancy to research topic and year of publishing. Both public and individual libraries, as well as online libraries, were chaffered to approach the data. Some of online databases that were accessed are SAGE, Questia, emerald, proudest and so on. Data collection establishments, for example, Gallup and AC Nielsen carry on researches on a repeated basis homing in on a wide lay of subjects. A library is an assemblage of services, resources and sources.

It is organized for the functioning of and is maintained by a

Thesis Statement:

There are a number of adverse impacts of cyber-crime on networks, and the network security reduces them to a significant extent.

Purpose:

The purpose of the study is to determine the impact of cybercrimes on network security and to determine at what level network security is able to reduce cyber-crimes.

To determine the pros and corn of network security.

To determine how network security reduces the treat of cyber-crimes.

public body, an organization, or even an individual. In the formal sense, a library is a collection of books. The term can mean the aggregation, the construction that homes such an assemblage, or both. Public and committed accumulations and services of process may be designated for use by individuals who prefer not to or cannot spend to purchase a huge collection, who require significant material that no person can fairly be anticipated to bear, or who demand professional help with his/her probe.

## 6. DATA ANALYSIS

Understanding the nature and function of cyber-crimes and network security; the qualitative descriptive mechanism is the most ideal means of collecting and analyzing data due to the flexibility, adaptiveness, and immediacy of the topic. This brings inherent biases, but another characteristic of such research is to identify and monitor these biases, thus including their influence on data collection and analysis rather than trying to eliminate them. Finally, data analysis in an interpretive qualitative research design is an inductive process. Data are richly descriptive and contribute significantly to this research.

## 7. DEFINITION OF QUALITATIVE RESEARCH

Qualitative research is much more subjective as compared to quantitative research and employs very unlike methods for data accumulation. These methods are primarily

interpersonal, in-depth consultations and focus groups. The nature of this type of research is explorative and open-ended. Small counts of people are questioned in-depth, and a comparatively small number of focus groups are conducted. Participants are called for to reply to general questions.

## **8. CONSIDERATION OF QUALITATIVE RESEARCH**

Qualitative research measures consider a universal law in the

## **5. RESEARCH APPROACH**

An effective strategy will be used to collect most of the information and data from various sources. The research will contain two parts, firstly "Secondary Research" in which the researcher will go through various research papers, electronic journals and database. Whereas for "Qualitative Research", different methods that will be used to collect information from companies.

## **9. RESULT**

Computer geniuses, usually in their twenties, are thrown challenges to break one or another security program, capture the passwords to remote computers and use their accounts to travel the cyberspace, enter data networks, airline reservation systems, banking, or any other "cave" more or less dangerous. Managers of all systems have tools to control that "all is well", if the processes are normal or if there is suspicious activity, a user is using to access roads which is not authorized. All movements are recorded in system files, high operators review daily (Farmer & Charles, Pp. 46). Furthermore, the network is becoming the ideal place for criminals and terrorists to carry out their actions and activities. Hence, cybercrime and cyber terrorism have become two of the most serious threats seem to haunt Western societies. Moreover, the impact of the crimes on the victims and their measures to cope up with such crimes in the future will also be a part of the paper. This paper will also discuss the how network security is critically important in preventing the recurrence of these types of cyber-attacks in the future.

## **10. DISCUSSION AND ANALYSIS**

After analyzing the results of the study through qualitative

hope of developing a static reality; on the other hand, qualitative research is an assumption of what the reality is a dynamic exploration. To solve the current problems of the market, it is important for retailers to adopt a more strategic approach to decision making, taking into account the great importance of intellectual property management, it does not say what is found in the process of universality, therefore, promotion and tolerance is often cited as a study of these types.

analysis it can be said that computers and the Internet are now a familiar part of our lives. You may not see them often, but they are involved in some way in most of our daily activities in the business, educational institutions, and government. Without the support of any of these tools, we would be able to handle the overwhelming amount of information that seems to characterize our society. But the problem of security limits the integrity of information and computer systems. More people need to know the use of computers and the protections that are daily offered for the safe handling of information (Roland, Pp. 638-645).

Cyber warfare has been defined as the process of nation-state to introduce computers of other countries or networks to cause damage or destruction. Cyber warfare is a form of warfare that occurs on computers and the Internet, by electronic means rather than physical. Moreover, the Internet is a means of easy access, where any person, remaining anonymous, can proceed with an attack that is difficult to associate, virtually undetectable and difficult to smuggle, let alone reaching a high impact such action directly hitting the opponent (network) by surprise. The term network security refers to protection against attacks and intrusions on corporate resources by intruders who are not allowed access to these resources.

During 1997, 54 percent of American companies were attacked by hackers in their systems. The incursions of hackers caused total losses of \$ 137 million that year. The Pentagon, CIA, UNICEF, the UN and other world bodies have been subjected to interference by these people who have much knowledge on the subject and a great ability to solve the obstacles they face. A hacker can take months to violate a system and increasingly sophisticated methodologies are used by the present day hackers (Ogut, Menon & Ragunathan, Pp. 14-28). A hacker enters a prohibited area to

gain access to confidential or unauthorized information.

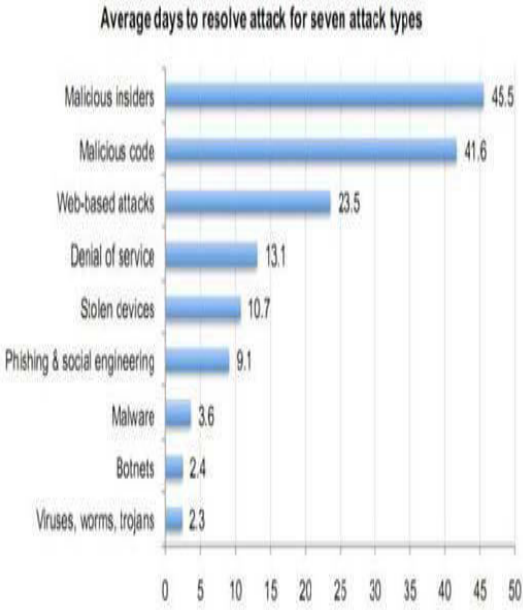
The mass media prefer to characterize them as criminals to intercept credit card codes and use them for personal gain. There are also those who intrude into airport systems producing chaos in flight schedules and aircraft. Today hacking is a prime concern of businessmen, legislators and officials.

A similar but different form of information intrusion is cracking. Crackers are people who disturb others; pirate software protected by law, destroying complex systems by transmitting powerful virus, and so on. Restless teenagers quickly learn this complex craft. They differ with hackers because they lack any kind of ideology when they do their "jobs". Instead, the main objective of hackers is not to become criminals, but "fight against an unjust system" used as a weapon system itself (Katz & Shapiro, Pp. 822-841).

Every organization should be at the forefront of change processes. Where continuous information is available, reliability and time is a key advantage. In fact, hacking is so easy that if you have an on-line and know how to send and read e-mail, you can start hacking immediately. Here you can find a guide where you can download programs especially appropriate for the hacker on Windows. Usually, these programs are free. They try to explain some easy hacker tricks that can be used without causing intentional damage (Tipton & Krause, Pp. 320-386).

The threats to the network security are not just for organizations, but can be observed all over the world in different countries and the degree to which each of them are exposed to threats. The statistics to which can be seen in the table illustrated below:

Table 1:



The Pirates of the cyber age considered as a sort of modern Robin Hood and demand a free and unrestricted access to electronic media (Whitman & Mattord, Pp. 205-249).

Another common attack on a computer system is the creation and distribution of malicious computer code, called "viruses". Computer viruses are computer programs written specifically to damage other computer systems. Sometimes these malicious programs are contained within another program, known as a "Trojan horse," and are copied by a user without his or her knowledge (Richards, Pp. 21-54).

The approximate time to resolve some categories of attacks on networks can see in the following table:

Table 2:

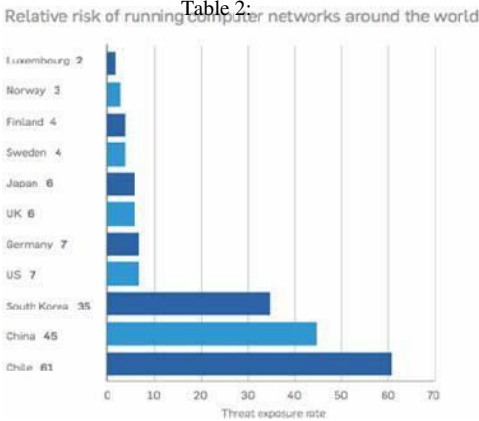


Table 3:

| <i>Name</i>            | <i>Description</i>                                                                                                                                                                                                                      |
|------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| BlackIce               | There are several levels of protection and warns us when we scanned through a sound and a flashing icon. It offers a wealth of information about the attackers and attacks following statistics, broken down by hours, days and months. |
| Conseal PC             | It is a bit outdated and tends to disappear when there is an installation of several old files.                                                                                                                                         |
| Tiny                   | It comes configured with a medium security level, suitable for normal Internet browsing. The network works fine without having to set special rules.                                                                                    |
| Protect X              | Leave some open ports and others as closed. It notes from where you've connected the port. Facilitates IP registration information.                                                                                                     |
| Sygate Firewall        | Interface comes configured with a high level of security.                                                                                                                                                                               |
| Win Route Pro          | It is a proxy server. Its source addresses filtering and destination of both incoming and outgoing. It does not consume system resources and is not necessary to install an sw.                                                         |
| Zone Alarm             | It not only detects all access from the Internet unless it gives control of the programs try to access the Internet. You can select levels of protection and block access to Internet after a certain time.                             |
| At Guard               | Allows you to define rules for everything, is fast and gives you a control of what is happening. Blocks unwanted advertising has a log of date, time, URL, IP, bytes sent and received and time of all Web connections.                 |
| Esafe Desktop          | Consumes many system resources, you can prohibit total or partial access to your computer.                                                                                                                                              |
| Freedom                | It is very easy to install, pass ports invisible so you can surf the Internet anonymously.                                                                                                                                              |
| Hack tracer            | You can identify the attacker, since it has a program that includes a world map with the route of the attacker's computer. Easy to install and default pass all ports invisible.                                                        |
| Internet Firewall 2000 | Does not work on local network, you can see active connections                                                                                                                                                                          |

## 12. BENEFITS OF NETWORK SECURITY

1. problems.
2. Provides a comprehensive system of warning alarms attempt to access your network.

The following table lists the most common types of Information Security Network Protection ISNP packages (Network Security Packages) in use by large organizations:

Today, for the implementation of effective measures for protecting information requires not only protection of information networks and mechanisms for a model of with an

1. Prevents unauthorized users from accessing your network.
2. Provides transparent access to Internet-enabled users.

3. Ensures that sensitive data is transferred safely by the public network.

Help your managers to find and fix security

## 13. CONCLUSION

In conclusion, it can be said that attacks on machines connected to the Internet have increased by 260% since 1994,

advances in software technology allowed the birth of a virtual world whose ultimate expression is the Internet.

estimated loss of 1,290 million dollars annually in the U.S.

In the era of information, ideas, data and files on your network are probably more valuable than your entire company. Think about your customer lists and records of

shareholders, trading and marketing materials, marketing strategies and product design, the loss of which could mean

the significant loss for your firm. With advances in technology, no one is safe from an attack by "hackers.

Currently it is relatively easy to gain control of a machine on the Internet that has not been adequately protected.

Companies invest a significant portion of their money in protecting their information, since the loss of irreplaceable

data is a real threat to their business. The technology boom in the development of networks, digital communications and

network security and implementation of a systematic

## REFERENCES

- [1] Casey, E. Digital Evidence and Computer Crime: Forensic Science, Computers and the Internet. London: Academic Press, 2011: Pp. 5-19.
- [2] Farmer, Dan. & Charles, Mann C. Surveillance nation. Technology Review; Vol. 106, No. 4, 2003: Pp. 46.
- [3] Harrison, A. Privacy group critical of release of carnivore data. Computerworld; Vol. 34, No. 41, 2006: Pp. 24
- [4] Internet Tax Freedom Act of 1998: 112 Stat. 2681-2719. Retrieved from: (<http://www.cbo.gov/doc.cfm?index=608&type=0>). Accessed on : 29th January, 2012.
- [5] Katz, Mira L. & Shapiro, Carl. Technology Adoption in the Presence of Network Externalities. Journal of Political Economy; Vol. 94, No.4, 1986: Pp. 822-841.
- [6] Ogut, Hulusi. Menon, Nirup. & Ragnathan, Srinivasia. Cyber Insurance and IT Security Investment: Impact of Independent Risk. Proceedings of the Workshop on the Economics of Information Security (WEIS), Cambridge, MA: Harvard University, 2005: Pp. 14-28.
- [7] Richards, James. Transnational Criminal Organizations, Cybercrime, and Money Laundering: A Handbook for Law Enforcement Officers, Auditors, and Financial Investigators. Boca Raton, FL: CRC Press, 1999: Pp. 21-54.
- [8] Roland, Sarah E. The Uniform Electronic Signatures in Global and National Commerce Act: Removing Barriers to E-Commerce or Just Replacing Them with Privacy and Security Issues?. Suffolk University Law Review; Vol. 35, 2001: Pp. 638-45.

# ROUTING INFORMATION PROTOCOL FOR WIRELESS SENSOR NETWORKS

Dr. V. Chandrasekar,  
Associate Professor,  
Dept. of Computer Science and Engineering,  
Malla Reddy College of Engineering & Tech.,  
drchandru86@gmail.com

P.Pavani  
Assistant Professor,  
Dept. of Computer Science and Engineering,  
Malla Reddy College of Engineering,  
Pavani20891@gmail.com

**Abstract** - A routing algorithm is a method for determining the routing of packets in a node. For each node of a network, the algorithm determines a routing table, which in each destination, matches an output line. The algorithm should lead to a consistent routing, that is to say without loop. This means that you should not route a packet a node to another node that could send back the package. Our contribution in this paper is e3D, diffusion based routing protocol that prolongs the system lifetime, evenly distributes the power dissipation throughout the network, and incurs minimal overhead for synchronizing communication. We compare e3D with other algorithms in terms of system lifetime, power dissipation distribution, cost of synchronization, and simplicity of the algorithm.

**Keyword:** Routing Algorithm, OSPF, RIP, LSP, Authentication

## INTRODUCTION

In this paper, we attempt to overcome limitations of the wireless sensor networks such as: limited energy resources, varying energy consumption based on location, high cost of transmission, and limited processing capabilities. Besides maximizing the lifetime of the sensor nodes, it is preferable to distribute the energy dissipated throughout the wireless sensor network in order to minimize maintenance and maximize overall system performance. For more in depth understanding of the problem statement and proposed algorithm, we refer the reader to the full length version of this paper [1].

Any communication protocol that involves synchronization between peer nodes incurs some overhead of setting up the communication. We attempt to calculate this overhead and to come to a conclusion whether the benefits of more complex routing algorithms overshadow the extra control messages each node needs to communicate. Obviously, each node could make the most informed decision regarding its communication options if they had complete knowledge of the entire network

topology and power levels of all the nodes in the network[2]. This indeed proves to yield the best performance if the synchronization messages are not taken into account. However, since all the nodes would always have to know everything, it should be obvious that there will be many more synchronization messages than data messages, and therefore ideal case algorithms are not feasible in a system where communication is very expensive. For both the diffusion and clustering algorithms, we will analyze both realistic and optimum schemes in order to gain more insight in the properties of both approaches[3]. The benefit of introducing these ideal algorithms is to show the upper bound on performance at the cost of an astronomical prohibitive synchronization costs.

## RIP (ROUTING INFORMATION PROTOCOL)

RIP is the most widely used protocol in the TCP / IP environment to route packets between the gateways of the Internet. It is a protocol IGP [4] (Interior Gateway Protocol),

which uses an algorithm to find the shortest path. y the way, refers to the number of nodes crossed, which must be between 1 and 15. The value 16 indicates impossibility. In other words, if the path to get from one point to another of the Internet is above 15, the connection can not be established. RIP messages to establish the routing tables are sent approximately every 30 seconds. If a RIP message does not reach its neighbor after three minutes, the latter considers that the link is no longer valid; the number of links is greater than 15 [5]. RIP is based on a periodic distribution of states network from a router to its neighbors. The release includes a RIP2 routing subnet, message authentication, multipoint transmission, etc.

## OSPF (Open Shortest Path First)

OSPF is part of the second generation of routing protocols. Much more complex than RIP, but at higher performance rates, it uses a distributed database that keeps track of the link state. This information forms a description of the network topology and the status of nodes, which defines the routing algorithm by calculating the shortest paths. The algorithm allows OSPF, from a node, to calculate the shortest path, with the constraints specified in the content associated with each link. OSPF routers communicate with each other via the OSPF protocol, placed on top of IP[ 6]. Now look at this protocol a bit more detail. The assumption for link state protocols is that each node can detect link status with its neighbors (on or off) and the cost of this link. We must give to each node enough information to enable him to find the cheapest route to any destination. Each node must have knowledge of its neighbors. If each node to the knowledge of other nodes, a complete map of the network can be established. An algorithm based on the state of the neighboring requires two mechanisms: the dissemination of reliable information on the state of the links and the calculation of routes by summing the accumulated knowledge of the link state [7]. One solution is to provide a reliable flood of information, to ensure that each node receives his copy of the information from all other nodes. In fact, each node floods its neighbors, which, in turn, flood their own neighbors. Specifically, each node creates its own update packets, called LSP (Link-State

Packet), containing the following information

Identity of the node that creates the LSP.

- List of neighboring nodes with the cost of the associated link.
- Sequence Number.
- Timer (Time to Live) for this message.

two information is needed to calculate routes. The last two aim to make reliable flooding. The sequence number allows putting in order the information that would have been received out of order. The protocol has error detection and retransmission elements [8]. The route calculation is performed after receiving all the information on the links. From the complete map of the network and costs of links, it is possible to calculate the best route. The calculation is performed using Dijkstra's algorithm on the shortest path. In the acronym OSPF (Open Shortest Path First) Open the word indicates that the algorithm is open and supported by the IETF. Using the mechanisms outlined above, the OSPF protocol adds the following additional properties

## AUTHENTICATION OF ROUTING MESSAGES

Malfunction can lead to disasters. For example, a node that, following the receipt of wrong messages, intentionally or not, or a striker messages modifying its routing table, calculates a routing table in which all nodes can be achieved at a cost zero automatically receives all network packets. These problems can be avoided by authenticating issuer's messages. Early versions had a OSPF authentication password of 8 bytes. The latest versions have much stronger authentication.

## NEW HIERARCHY

This hierarchy allows for better scalability. OSPF introduces another level of hierarchy by partitioning the areas into eras (area). This means that a router within a domain does not need to know how to reach all the networks in the field. Just that he knows how to reach the right age. This results in a reduction of information to be transmitted and stored.

There are several types of OSPF messages, but they all use the same header, which is shown in Figure.

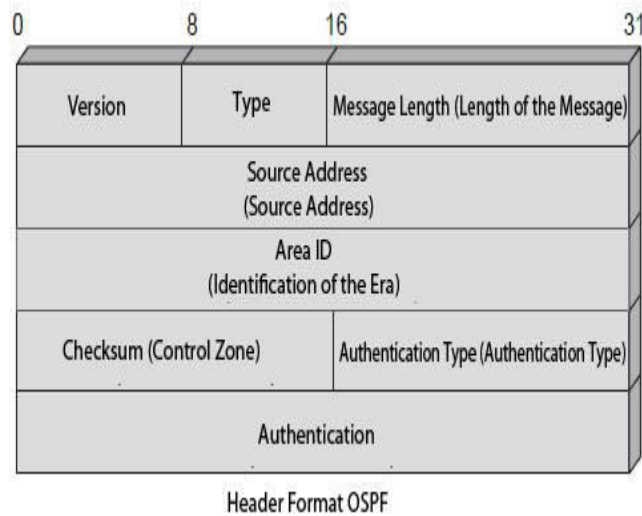
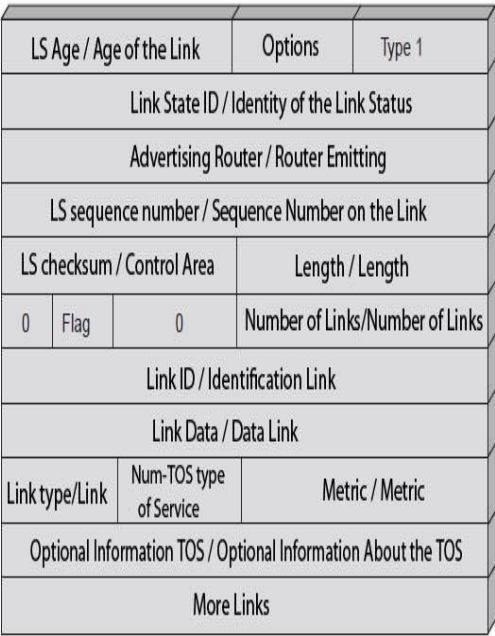


Fig: OSPF messages

The current version is 2. Five types were defined with values from 1 to 5. The source address indicates the sender of the message. The identification of the era indicates the era in which lies the sending node. The authentication type has the value 0 if there is no authentication, 1 if the authentication password and 2 if an authentication technique is implemented and described in the following 4 bytes.

The five types of messages have the Hello message as Type 1. This message is sent by a node to its neighbors to tell them that it is always present and not broken. The four other types are used to send information such as

queries, shipments or acquittals LSP messages. These messages mainly carrying LSA (Link-State Advertisement), that is to say, information about link state. A message can contain several OSPF LSA.



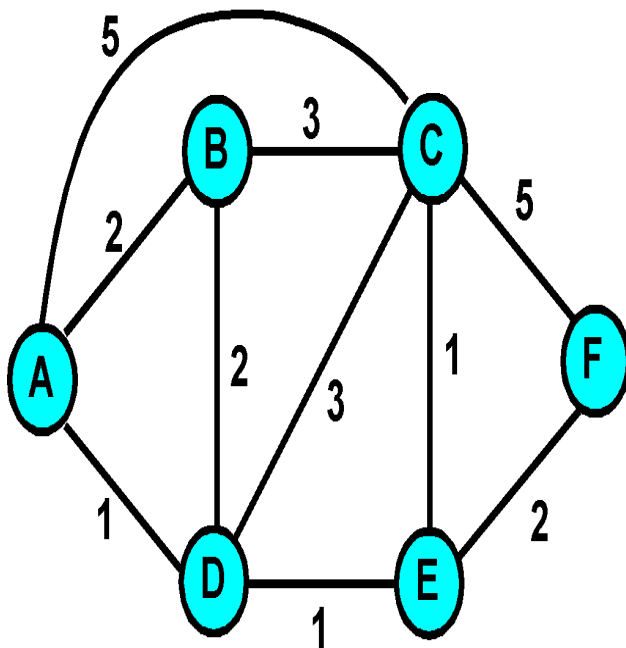
OSPF message with an LSA

Fig: type of OSPF LSA carrying a message

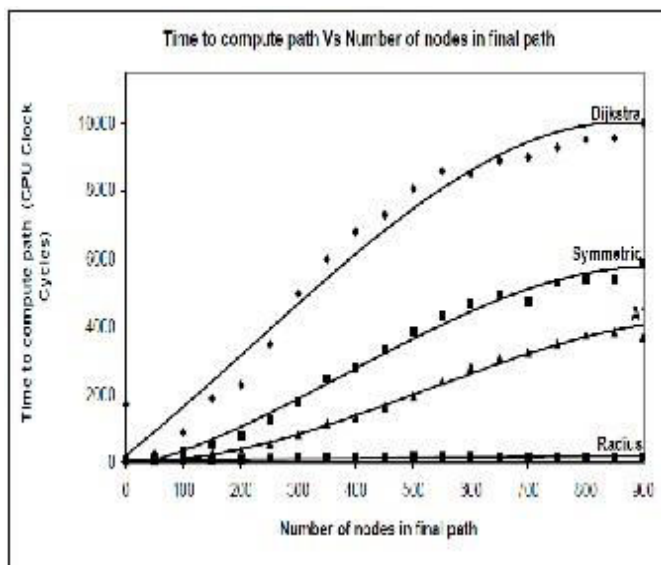
e3D: EXPERIMENTAL RESULTS

We introduce a new algorithm, e3D, and compare it to two other algorithms, namely directed, and random clustering communication. We take into account the setup costs and analyze the energy-efficiency and the useful lifetime of the system [10]. We compare the algorithms in terms of system lifetime, power dissipation distribution, cost of synchronization, and simplicity of the algorithm. Our simulation results show that e3D performs only slightly worse than its optimal counterpart while having much less overhead. Therefore, our contribution is diffusion based routing protocol that prolongs the system lifetime, evenly distributes the power dissipation throughout the network, and incurs minimal overhead for synchronizing communication. In our simulation, we use a data collection problem in which the system is driven by rounds of communication, and each sensor node has a packet to send to the distant base station [11]. The diffusion algorithm is based on location, power levels, and load on the node, and its goal is to distribute the power consumption throughout the network so that the majority of the nodes consume their power supply at relatively the same rate regardless of physical location. This leads to better maintainability of

the system, such as replacing the batteries all at once rather than one by one, and maximizing the overall system performance by allowing the network to function at 100% capacity throughout most of its lifetime instead of having a steadily decreasing node population.



**Fig:** Time to Compute Path vs Number of Node in Final Path



**Fig:** Routing Graph

## CONCLUSION

In summary, we showed that energy-efficient distributed dynamic diffusion routing is possible at very little overhead cost. The most significant outcome is the near optimal performance of e3D when compared to its ideal counterpart in which global knowledge is assumed between the network nodes. We therefore conclude that complex clustering techniques are not necessary in order to achieve good load and power usage balancing. Previous work suggested random clustering as a cheaper alternative to traditional clustering; however, random clustering cannot guarantee good performance according to our simulation results.

## REFERENCES

- [1] Ioan Raicu, Scott Fowler, Loren Schwiebert, Sandeep K.S. Gupta. "Energy-Efficient Distributed Dynamic Diffusion Routing Algorithm in Wireless Sensor Networks: e3D Diffusion vs. Clustering" submitted for review to ACM MOBICOM 2017;
- [2]. Nilesh P. Bobade, Performance Evaluation of Ad Hoc on Demand Distance Vector in MANETs with varying Network Size using NS-2 Simulation, (IJCS) International Journal on Computer Science and Engineering, 02(08), 2731-2735 (2015)
- [3]. Revathi Venkataraman, Pushpalatha .M, and Rama Rao.T, Performance Analysis of Flooding Attack Prevention Algorithm in MANETs, World Academy of Science, Engineering and Technology, 32 (2014)
- [4]. Venkateshwara Rao .K, Dynamic Search Algorithm used in Unstructured Peer- to-Peer Networks, International Journal of Engineering Trends and Technology, 2(3), (2014)
- [5]. Cheolgi Kim, Young-Bae Koy and Nitin H.Vaidya, LinkState Routing Protocol for Multi-Channel Multi-Interface Wireless Networks, IEEE (978-1-4244-2677- 5/08)/\$25.00c (2008)

[6]. Kiruthika R., "An exploration of count-to-infinity problem in networks" *International Journal of Engineering Science and Technology*, 2(12), 7155-7159 (2010)

[7]. Mohammad reza soltan aghaei, "A hybrid algorithm for finding shortest path in network routing," *Journal of Theoretical and Applied Information Technology* © 2005- 2009 (2009)

[8]. Taehwan Cho, "A Multi-path Hybrid Routing Algorithm in Network Routing," *International Journal of Hybrid Information Technology*, 5(3), (2012)

[9]. M. Abolhasan, J. Lipman, and J. Chicharo, "A routing strategy for heterogeneous mobile ad hoc networks," in *Proceedings of the 6th IEEE Circuits and Systems Symposium on Emerging Technologies: Frontiers of Mobile and Wireless Communication*, Vol. 1, 2004, pp. 13-16

[10]. A. D. Amis and R. Prakash, "Load-balancing clusters in wireless ad hoc networks," in *Proceedings of the 3rd IEEE Symposium on Application-Specific Systems and Software Engineering Technology*, 2000, pp. 25-32

[11]. A. D. Amis, R. Prakash, T. H. P. Vuong, and D. T. Huynh, "Max-min d-cluster formation in wireless ad hoc networks," in *Proceedings of the 19th IEEE Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM)*, Vol. 1, 2000, pp. 32-41.

[12]. T. Asano, H. Unoki, and H. Higaki, "LBSR: routing protocol for MANETs with unidirectional links," in *Proceedings of the 18th International Conference on Advanced Information Networking and Applications*, Vol. 1, 2004, pp. 219-224

[13]. M. Chatterjee, S. K. Sas, and D. Turgut, "An on-demand weighted clustering algorithm (WCA) for ad hoc networks," in *Proceedings of the IEEE Global Telecommunications Conference*, Vol. 3, 2000, pp. 1697-1701

[14]. Y. S. Chen, Y. C. Tseng, J. P. Sheu, and P. H. Kuo, "An on-demand, link-state, multi-path QoS routing in a wireless mobile ad hoc network," *Computer Communications*, Vol. 27, 2004, pp. 27-40

[15]. C. C. Chiang, H. K. Wu, W. Liu, and M. Gerla, "Routing in clustered multihop, mobile wireless networks with fading channel," in *Proceedings of the IEEE Singapore International Conference on Networks*, 1997, pp. 197-211.

# Realistic Future Flying Car

**Mr.U.Nagaiah**  
Assistant Prof.

Dept. of CSE  
MRCE-Hyderabad

[nagaiah1212@gmail.com](mailto:nagaiah1212@gmail.com)

**Ms.Razia sultana**  
Assistant Prof.

Dept. of CSE  
MRCE-Hyderabad

[razia.sultana13@gmail.com](mailto:razia.sultana13@gmail.com)

**Mr.M.Amarnath**  
Assistant Prof.

Dept. of CSE  
MRCE-Hyderabad

[amar.sap16@gmail.com](mailto:amar.sap16@gmail.com)

**Dr.Sunil Tekale**  
Professor

Dept of CSE  
MRCE-Hyderabad

[Sunil.tekale2010@gmail.com](mailto:Sunil.tekale2010@gmail.com)

**Abstract:** As the number of vehicles are increasing at a very high rate on the roads and it is almost becoming impossible to travel, there needs to have a solution for the traffic congestion. Many methods were tried and almost every method has some or the other drawback[7]. The only solution for reducing the traffic congestion is to have triple mode car, which should run on road, water and should fly in the sky[3]. This will definitely reduce the traffic congestion and will also provide or perhaps helps to start a thought process in designing the same concept car. These cars will be useful for different section of people in terms of commercial and personnel use. People can travel on their own or can use the same for delivery of goods from one place to another and even for lifting patients from a place to nearby hospital. Which means this car can be a life saver vehicle. Many issues needs to be addressed in order to make it a safe triple mode car. Here we have to design a car with road safety measures ,safety measures necessary to fly and also we need to take care of safety measures to run on water It has become inexcusably obvious that our technology has exceeded our humanity[8].---

Key words: Traffic congestion, Design car, Safety measures, Radar, Flying ,Futuristic.

**Albert Einstein.**

**Introduction:** 'Flying car', 'Street car', 'swimming car' a triple mode car will help to fulfill the long pending dreams of aviation, automobile, and navy enthusiasts[6]. As this car will bring the best in 3 worlds. The basic purpose of this car is to solve the problems pertaining to traffic congestion on roads, where we find many people getting stuck every day in this traffic which not only damages their health and also waste lot of time on travelling. The concept here is to see that this car not only allow people to travel on road but also to fly in the sky depending on the requirement and distance to travel, apart from swimming in the water[2]. The car will have to cater to different needs of the people and will help the future generation to travel in the manner they prefer. The designing of this car will have multiple obstacles

as it has to satisfy the regulation of the 3 different worlds. Tech titans like Uber, Amazon, and Google have all laid out ambitious plans for filling the skies with autonomous aircraft[4]. Uber wants to move people

around with flying taxis, and Airbus is committed to producing this kind of vehicle.



Meanwhile Google and Amazon are hoping to deliver packages with much smaller drones. All see the potential for fleets of unmanned aerial vehicles that can pilot themselves[5]. But to make that vision a reality, we're going to need a new breed of sense and avoid technology. Echodyne, a Bellevue, Washington-based startup, believes it has the answer. The company announced preliminary test results from field trials of its MESA-DAA radar system today[2]. It says the device, which is barely larger than smartphone, is capable of detecting even small aircraft at a distance 1.8 miles in varying weather conditions. The company says this breakthrough is driven by the use of meta materials, which allow the radar to eliminate moving parts, making the hardware smaller and more battery efficient without sacrificing range[1].A lot of modern

automobiles are now equipped with radar systems, in fact Tesla recently announced that it would be focusing on radar as the core technology in its autonomous driving system[1]. But even long-range automotive radar from the likes of Bosch and Delphi only claim a range of a few hundred meters. They also don't typically have a very wide field of view. Echodyne's technology claims to be able to have a

120 degree field of view in azimuth (horizontal) and 80 degrees in elevation (vertical). Founder and CEO Eben Frankenberg also says his tech was designed to track a Cessna-sized object, which is much smaller than a car in radar cross section[5].

#### **DESIGN CONSIDERATION:**

Some of the specifications to be determined by the design approach are:

- Range of run
- Endurance
- Rate of climb
- Cruise speed in air
- Cruise speed in land
- Cruise speed in Water
- Airworthiness standards
- Automobile safety and emissions.

#### **Challenge:**

A practical flying car should be capable of safely taking off, flying and landing throughout heavily populated urban environments. However, to date, no vertical takeoff and landing (VTOL) vehicle has ever demonstrated such capabilities[3].

Driving a flying car would require a pilot's certificate and also an initial training of 18- 40 hours and foremost, along with a driver's license for a flying car. The flying car would require intensive maintenance for keeping it in perfect workable condition matching initial technical standards and rules for government regulatory requirements. Some of the major challenges in the flying car technology is the VTOL capability, the powering-system for the vehicle and also many safety issues. It seems the technologists have really been able to find viable solutions to these problems and the name 'flying cars' would actually be replaced by something more technical. And very soon the flying cars will be taking to the skies[1].

#### **Safety:**

Although statistically commercial flying is much safer than driving, unlike commercial planes, personal flying cars might not have as many safety checks and their pilots would not be as well trained[2]. Humans already have problems with the aspect of driving in two dimensions (forward and backwards, side to side),

adding in the up and down aspect would make "driving" or flying as it would be, much more difficult; however, this problem might be solved via the sole use of self-flying and self-driving cars.

#### **Economies:**

In addition, the flying car's energy efficiency would be much lower compared to conventional cars and other aircraft; optimal fuel efficiency for airplanes is at high speeds and high altitudes, while flying cars would be used for shorter distances, at higher frequency, lower speeds and lower altitude. For both environmental and economic reasons, flying cars would be a tremendous waste of resources[4].

**Existing System:** In the existing aviation industry, much of the mechanics of flying is automated[1]. Given the challenges of a person flying compared to driving a car, and the efforts to reduce human error in aviation, there is even more likelihood of flying cars becoming automated so that no human pilot is needed. But there will be differences between existing aviation practice and flying cars. Passenger jet air travel owes much of its impressive safety record to improvements in aircraft maintenance procedures and our understanding of failures.

It is unlikely that the business case for small flying cars will allow for such rigorous practices. Instead, flying cars will be less complex than modern jets, and the latest demonstrators show exactly that. The use of large numbers of small electric motors, such as in the Lilium all-electric aircraft, reduces the maintenance complexity drastically. It also provides an inbuilt measure of redundancy in case one motor fails[4].

**How flying car works:** we're currently in the midst of a new round of flying-car hype. Uber is even having some big flying car event in Texas this week. Historically, every bit of flying-car hype proves to be bullshit. But it may not have to be; I think I have an idea about how flying cars could make sense, even it's not exactly how Uber is imagining it.



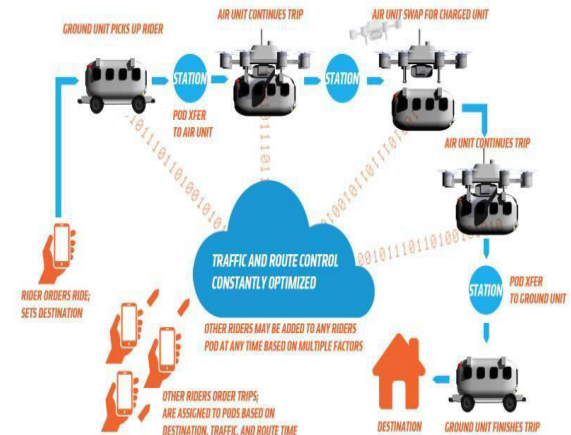
Uber's plan, as it seems to stand now, relies on the use of small, vertical-takeoff-and-land (VTOL) aircraft. Basically, just using tiny planes inside a city.

I think this approach is too simplistic, and won't be able to scale in any way that makes sense. I think I have a better idea.

Now, your gut reaction may be to not like it because the fundamental, romantic appeal of flying cars has always been the incredible independence and freedom of them. Flying cars, going all the way back to the first experiments in the late '20s and through the era of flying Pintos and into today, have always conjured up images of living on inaccessible islands and flying into work every day, and jaunting off to wherever you'd like, looking contemptuously at the ground and those miserable bastards stuck in traffic[2].

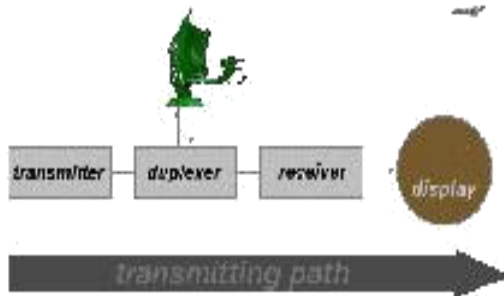
The truth is, though, that's just not going to work. If we want flying cars to happen, we have to get rid of the concept of just translating personal, private cars into the air. The logistics and traffic management of large numbers of independent flying cars is just too inefficient and difficult to manage, and the issues arising from mechanical failure or driver/pilot error are too unforgiving. We need a new way to conceptualize the whole idea, and the way that works isn't exactly sexy: flying (and driving) buses. Buses, like most public transportation, isn't nearly as exciting as your own private car, but it makes more sense in this context. Plus flying has a huge advantage over every other public transportation network out there, and that advantage is the key to why flying buses actually make sense: the routes can be dynamic, optimized, and changeable on demand.

There's no massive, city-wide infrastructure cost for flying vehicles. There's no roads to build, no tracks to lay, no tunnels to dig. I'm imagining a system of flying/driving vehicles that use a resource we have plenty of computing power to constantly adapt and change to meet demand and traffic. First, we have the basic vehicles, which consist of three main components: a passenger pod, a VTOL air unit, and a street-going wheeled ground unit. The passenger pod can dock with either the air unit or the ground unit as needed.



The air unit and ground unit are autonomously controlled, in constant communication with a central traffic controlling and route managing system. A human pilot/driver can be on board to act as an emergency safety backup, but would not be able to fly the vehicle in normal operation[1].

Oh, and the airspace for the flying would be in a zone just above skyscraper-level in a given city, well below the altitude of commercial (and private) aircraft. The FAA can figure out the parameters of that. These vehicles can be powered by combustion motors or can be electric; since there seems to be a desire for systems like these to be electric (reduced pollution, noise, efficiency advantages) let's consider them to be electric for now. Both vehicles can drive or fly with or without the passenger pod. Moving without the passenger pod will only need to happen for purposes of moving resources/vehicles from station to station or for recharging. There are two basic types of stations in this transportation system, rooftop and ground-level. Rooftop stations will be more common in large, dense cities, and will only have air units available. The distance between stations and number of stations will be dependent on the range of the air units[2].



All targets produce a diffuse reflection i.e. it is reflected in a wide number of directions. The reflected signal is also called scattering. Backscatter is the term given to reflections in the opposite direction to the incident rays[4].

Radar signals can be displayed on the traditional plan position indicator (PPI) or other more advanced radar display systems. A PPI has a rotating vector with the radar at the origin, which indicates the pointing direction of the antenna and hence the bearing of targets[2].

#### Transmitter

The radar transmitter produces the short duration high-power rf pulses of energy that are into space by the antenna[3].

#### Duplexer

The duplexer alternately switches the antenna between the transmitter and receiver so that only one antenna need be used. This switching is necessary because the high-power pulses of the transmitter would destroy the receiver if energy were allowed to enter the receiver.

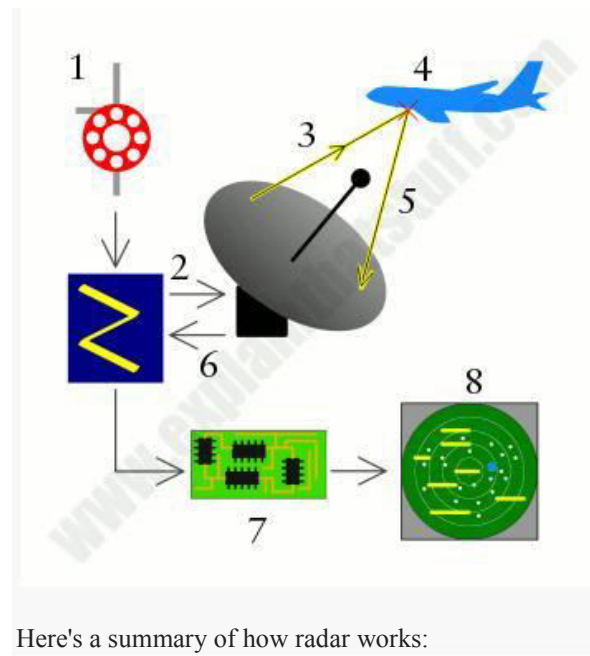
#### Receiver

The receivers amplify and demodulate the received RF-signals. The receiver provides video signals on the output. Radar Antenna The Antenna transfers the transmitter energy to signals in space with the required distribution and efficiency. This process is applied in an identical way on reception.

#### Indicator

The indicator should present to the observer a continuous, easily understandable, graphic picture of the relative position of radar targets.

The radar screen (in this case a PPI-scope) displays the produced from the echo signals bright blips. The longer the pulses were delayed by the runtime, the further away from the center of this radar scope they are displayed. The direction of the deflection on this screen is that in which the antenna is currently pointing.



Here's a summary of how radar works:

1. Magnetron generates high-frequency radio waves.
2. Duplexer switches magnetron through to antenna.
3. Antenna acts as transmitter, sending narrow beam of radio waves through the air.
4. Radio waves hit enemy airplane and reflect back.
5. Antenna picks up reflected waves during a break between transmissions. Note that the same antenna acts as both transmitter and receiver, alternately sending out radio waves and receiving them.
6. Duplexer switches antenna through to receiver unit.
7. Computer in receiver unit processes reflected waves and draws them on a TV screen.
8. Enemy plane shows up on TV radar display with any other nearby targets.

**Conclusion:** As flying car companies innovate their business models, an array of new business services are expected, such as aerial sightseeing, air surveillance-as-a-service, aerial critical aid delivery, air taxi pay-per-ride, and flying car corporate lease. Various flying car market participants have adopted different strategies for growth and expansion: Ehang is developing a flying drone with VTOL and autonomous flying capabilities · Toyota has acquired a patent for [Aerocar](#), a shape-shifting flying car and also invested in Cartivator, a Japanese flying car start-up· Airbus, [Carplane](#), and Lillium are expected

to release flying cars in the next five years[2]. Pre-selling of PAL-V's Liberty Pioneer flying car has begun, with delivery expected by 2018. Airbus self-flying aircraft Vahana is scheduled for production by 2021. Kitty Hawk is developing a flying car with investment from Google. Flying car prototypes are being developed by AeroMobil and Terrafugia. Airbus, in collaboration with Italdesign, is developing autonomous systems for its Pop.Up flying car. VTOL capabilities, autonomous flying technologies and the development of fail-safe features, will be imperative to inspire confidence in potential customers and overall acceptance of flying cars as vehicles for urban mobility," noted the analyst. "Makers of flying cars must work with regulators to ensure that clearly defined and industry-friendly rules for flying car operations are passed."

#### References:

1. B. Karthik IOSR Journal of Mechanical and Civil Engineering (IOSR-JMCE) e-ISSN: 2278-1684, p-ISSN: 2320-334X
2. Klein, S. and Smrcek, L. (2009) Flying car design and testing. In: CEAS 2009 European Air and Space Conference, 26-29 Oct 2009, Manchester, UK.
3. EduRef (2002). "Integrating Expertise into the NSDL: Putting a Human Face on the Digital Library." [Online]  
<http://www.eduref.org/eduref/Default.htm>
4. Gross, M. (2001). "What About the User?" Quality Study Bulletin 080102. Available from the Information Institute of Syracuse [online]
5. Janes, Joseph, 1962-; McClure, Charles R. Source: Public Libraries v. 38 no1 (Jan./Feb. 1999) p. 30-3 Libraries: 854
6. Lamolinara, G. and Grünke, R. (1998). "Reference Service in a Digital Age." LC INFORMATION BULLETIN August 1998. [Online]  
<http://www.loc.gov/loc/lcib/9808/ref.html>
7. Lankes, R. D. and Kasowitz, A. and Collins, J. (2000). Digital Reference: Models for the New Millennium. Lankes, R. David, Collins, J. & Kasowitz, A. S. (Eds.). New York: Neal-Schuman.
8. Wurman, R. S. (1989). Information Anxiety. New York : Doubleday

# Optimal Jamming Attack Detection in WSN

<sup>1</sup>P.Poovizhi, <sup>2</sup>S.Yamuna <sup>3</sup>M.Brindha, <sup>4</sup>V.Poorani

<sup>1</sup>Assistant Professor, Department of Computer Science and Engineering, SNS College of Engineering, Coimbatore

<sup>2</sup>Assistant Professor, Department of Computer Science and Engineering, SNS College of Engineering, Coimbatore

<sup>3</sup>Assistant Professor, Department of Electronics and Instrumentation Engineering, SNS College of Technology, Coimbatore

<sup>4</sup>Assistant Professor, Department of Electronics and Instrumentation Engineering, SNS College of Technology, Coimbatore

Email id:poovizhiponnusamy27@gmail.com, yamuna205@gmail.com,  
mbrindhaie@gmail.com, poorani28@gmail.com

**Abstract—** Consider a scenario where a sophisticated jammer jams an area in which a single-channel random-access-based wireless sensor network operates. The jammer controls the probability of jamming and the transmission range in order to cause maximal damage to the network in terms of corrupted communication links. The jammer action ceases when it is detected by the network (namely by a monitoring node), and a notification message is transferred out of the jammed region. The jammer is detected by employing an optimal detection test based on the percentage of incurred collisions. On the other hand, the network defends itself by computing the channel access probability to minimize the jamming detection plus notification time. The necessary knowledge of the jammer in order to optimize its benefit consists of knowledge about the network channel access probability and the number of neighbors of the monitor node. Accordingly, the network needs to know the jamming probability of the jammer. The latter is captured by formulating and solving optimization problems where the attacker and the network respond optimally to the worst-case or the average-case strategies of the other party. We extend the problem to the case of multiple observers and adaptable jamming transmission range and propose a meaningful heuristic algorithm for an efficient jamming strategy.

## 1. INTRODUCTION

The fundamental characteristic of wireless networks that renders them more vulnerable to attacks than their wire line counterparts is the open, shared nature of their medium. This exposes them to two fundamentally different attacks: passive and active attacks. In the former ones, the malicious

entity does not take any action apart from passively observing the ongoing communication that is, eavesdropping with the intention to intervene with the privacy of network entities involved in the transaction. On the other hand, in active attacks the attacker is involved in transmission as well. Depending on attacker objectives, different terminology is used. If the attacker abuses a protocol with the primary goal to obtain performance benefits itself, the attack is referred to as misbehavior. If the attacker does not directly manipulate protocol parameters but exploits protocol semantics and aims at indirect benefits by unconditionally disrupting network operation, the attack is termed jamming or Denial-of-Service (DoS), depending on whether one looks at the cause or the consequences of it. Misbehavior in wireless networks stems from the selfish inclination of wireless network entities to improve their own derived utility at the expense of other nodes' performance deterioration, by deviating from legitimate protocol operation at various layers. The utility is expressed in terms of consumed energy or achievable throughput on a per link or end-to-end basis. The first case arises if a node denies to forward messages from other nodes so as to preserve battery. The latter case occurs when a node prevents other nodes from accessing the channel or from routing their messages to destinations by selfish manipulation of the access control and routing protocol, respectively. The work in focuses on optimal detection of access layer misbehavior in terms of number of required observation samples to derive a decision. The worst-case attack is found out of the class of most significant attacks in terms of incurred

performance losses. The framework captures uncertainty of attacks and the case of intelligent attacker that can adapt its policy to delay its detection. Jamming can disrupt wireless transmission and occur either unintentionally in the form of interference, noise, or collision at the receiver, or in the context of an attack.

A jamming attack is particularly effective from the attacker's point of view since 1) the adversary does not need special hardware to launch it, 2) the attack can be implemented by simply listening to the open medium and broadcasting in the same frequency band as the network uses, and 3) if launched wisely, it can lead to significant benefits with small incurred cost for the attacker. With regard to the machinery and impact of jamming attacks, they usually aim at the physical layer in the sense that they are realized by means of a high transmission power signal that corrupts a communication link or an entire area. Conventional defense techniques against physical layer jamming rely on spread spectrum which can be too energy consuming for resource-constrained sensors. Jamming attacks also occur at the access layer, whereby an adversary either corrupts control packets or reserves the channel for the maximum allowable number of slots, so that other nodes experience lower throughput by not being able to access the channel. The work in studies the problem of a legitimate node and a jammer transmitting to a common receiver in an on-off mode in a game-theoretic framework. Other jamming instances can have impact on the network layer by malicious packet injection along certain routes or at the transport layer by SYN message flooding for instance. The work in presents attack detection in computer networks based on observing the IP port scanning profile prior to an attack and using sequential detection techniques. The work uses controlled authentication to detect spam message attacks in wireless sensor networks launched by a set of malicious nodes and addresses the tradeoff between resilience to attacks and computational cost. Sensor networks are susceptible to jamming attacks since they rely on deployed miniature energy-constrained devices to perform a certain task without a central powerful monitoring point.

Wood and Stankovic provide a taxonomy of DoS attacks launched against sensor networks from the physical up to the transport layer. Law et al. present attacks aimed at sensor network protocols that are based on learning protocol semantics such as temporal arrangement of packets, time slot size, or packet preamble size. In one paper, low-energy attacks are analyzed, which corrupt a packet by

jamming only a few bits such that the code error correction capability is exceeded. Low Density Parity Check (LDPC) codes are proposed as a method to defend against these attacks. The work in considers passing attack notification messages out of a jammed region by creation of wormhole links between sensors, one of which resides out of the jammed area. The links are created through frequency hopping over a channel set either in a predetermined or in an ad hoc fashion. In one paper, a physical layer jammer termed constant jammer, and three types of link layer jammer termed deceptive, random, and reactive jammer are studied. The reactive jammer is the most sophisticated one as it launches its attack after sensing ongoing transmission. The authors propose empirical methods based on signal strength and packet delivery ratio measurements to detect jamming. In paper, Channel surfing involves on-demand frequency hopping as a countermeasure against jamming is studied. The case of an attacker that corrupts broadcasts from a base station (BS) to a sensor network is considered. The interaction between the attacker and the BS is modeled as a zero-sum game with a long-term payoff for the attacker. The attacker selects the number of sensors it will jam and the BS chooses the probability with which it will sample sensor status with regard to message reception.

In this paper, we study controllable jamming attacks that are easy to launch but are difficult to detect and confront, since they differ from brute force attacks. The jammer controls the probability of jamming and the transmission range in order to cause maximal damage to the network in terms of corrupted communication links. We assume that the effect of jammer action ceases when it is detected by one or more monitoring nodes, and a notification message is transferred out of the jamming region. Following this notification message, drastic actions are presumably taken by the network in order to isolate, penalize, localize, and even physically capture the attacker. These actions are, however, not addressed further in this work. The fundamental trade-off faced by the attacker is the following: a more aggressive attack, either in terms of higher jamming probability or larger transmission range increases the instantaneous payoff but exposes the attacker to the network and facilitates its detection and, later on, its isolation. In an effort to withstand the attack, alleviate the attacker benefit, and expose the attacker to the detection system, the network controls the channel access probability of the employed random access protocol. The necessary knowledge of the jammer in order to optimize its benefit consists of knowledge about the network channel access probability and the number

of neighbors of the monitor node. Accordingly, the network needs to know the jamming probability.

## 2. MODELING ASSUMPTIONS

### 2.1 Sensor Network Model

We consider a wireless sensor network deployed over a large area and operating under a single-carrier slotted Aloha type random access protocol. We assume symmetric transmission and reception in the sense that a node  $i$  can receive a signal from node  $j$  if and only if node  $j$  can receive a signal from  $i$ . Time is divided into time slots and the slot size is equal to the size of a packet. All nodes are assumed to be synchronized when transmitting with respect to time slot boundaries.

Each node transmits at a fixed power level  $P$  with an omnidirectional antenna and its transmission range  $R$  and sensing range  $R_s$  are circular with sharp boundary. Transmission and sensing ranges are defined by two thresholds of received signal strength. A node within transmission range of node  $i$  can correctly decode transmitted messages from  $i$ , while a node within sensing range can just sense activity due to higher signal strength than noise, but cannot decode the transmitted message. Typically,  $R_s$  is a small multiple of  $R$ , ranging from 2 to 3. A node within distance  $R$  of a node  $i$  (excluding node  $i$  itself) is called a neighbor of  $i$ . The neighborhood of  $i$ ,  $N_i$  is the set of all neighbors of  $i$  with  $n_i = |N_i|$  being the size of  $i$ 's neighborhood. Transmissions from node  $i$  are received by all its neighbors. The sensor network is represented by an undirected graph  $G = (S, E)$  where  $S$  is the set of sensor nodes and  $E$  is the set of edges where edge  $(i, j)$  denotes that sensor  $i$  and  $j$  are within transmission range of each other. Sensor nodes are uniformly distributed in an area, with spatial density  $\lambda$  nodes per unit area and the topology is static, i.e., we assume no mobility. Each node has an initial amount of energy  $E$ . We do not consider the energy consumed in reception.

Each node is equipped with a single transceiver, so that it cannot transmit and receive simultaneously. All nodes are assumed to be continuously backlogged, so that there are always packets in each node's buffer in each slot. Packets can be generated by higher layers of a node, or they may come from other nodes and need to be forwarded or they may be previously sent and collided packets to be retransmitted. A transmission on edge  $\delta i; j$  is successful if and only if no node in  $N_j$  transmits during that transmission. In this

work, we consider the class of slotted Aloha type random access protocols that are characterized by a common channel access probability  $\gamma$  for all network nodes in each slot. This provides us with a straightforward means to quantify the network effort to withstand and confront the attack by regulating the amount of transmitted traffic and essentially exposing the attacker to the detection system, as will become clear in the sequel. Provided that it remains silent in a slot, a receiver node  $j$  experiences collision if at least two nodes in its neighborhood transmit simultaneously, regardless of whether the transmitted packets are destined for node  $j$  or for other nodes. Thus, the probability of collision at node  $j$  in a slot is

$$\theta_0 = 1 - (1 - \gamma)^{n_j} - n_j \gamma (1 - \gamma)^{n_j - 1}$$

If node  $j$  attempts to transmit at a slot while it receives a message, a collision occurs as well. In that case, the receiver is not in position to tell whether the collision is due to its own transmission or whether it would occur anyway. In the sequel, we will term collision an event addressing the case of multiple simultaneous transmissions received by (not necessarily intended to) a node and no transmission attempt by that node.

### 2.2 Attacker Model

We consider one attacker, the jammer, in the sensor network area. The jammer is neither authenticated nor associated with the network. The objective of the jammer is to corrupt legitimate transmissions of sensor nodes by causing intentional packet collisions at receivers. Intentional collision leads to retransmission, which is translated into additional energy consumption for a certain amount of attainable throughput or equivalently reduced throughput for a given amount of consumed energy. In this paper, we do not consider the attacker that is capable of node capture. The jammer may use its sensing ability in order to sense ongoing activity in the network. Clearly, sensing ongoing network activity prior to jamming is beneficial for the attacker in the sense that its energy resources are not aimlessly consumed and the jammer is not needlessly exposed to the network. The jammer transmits a small packet which collides with legitimate transmitted packets at their intended receivers. A beacon packet of a few bits suffices to disrupt a transmitted packet in the network. The jammer is assumed to have energy resources denoted by  $E_m$ , yet the corresponding energy constraint in the optimization problems of the next section may be

redundant if the jammer adheres to the policy above. The jammer uses an omnidirectional antenna with circular sensing range  $R_{ms}$  and adaptable transmission range  $R_m$  that is realized by controlling transmission power  $P_m$  as illustrated in Fig. 1.

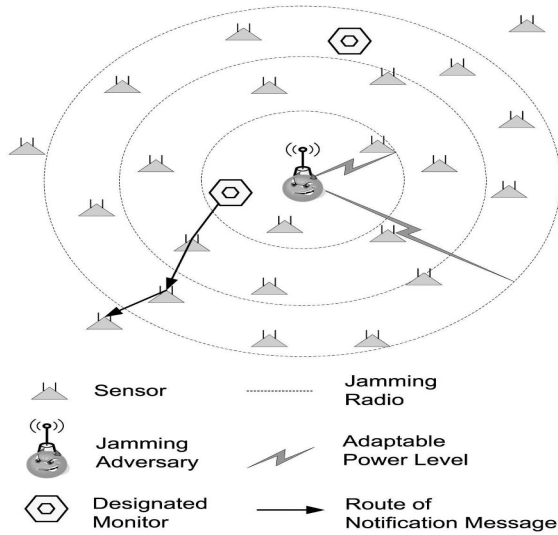


Fig. 1. Illustration of jamming attack. The jamming adversary can jam with different transmit power levels to disrupt network operation while avoiding detection. We assume that there exist designated monitoring nodes for detecting the jamming attack. Upon detection, a notification is routed out of the jammed region in a multihop fashion. The jammer can prolong the transfer of such a notification message by continuing jamming after detection.

### 2.3 Attack Detection Model

The network employs a mechanism for monitoring network status and detecting potential malicious activity. The monitoring mechanism consists of: 1) determination of a subset of nodes that act as monitors, and 2) employment of a detection algorithm at each monitor node. The assignment of the role of monitor to a node is affected by potential existing energy consumption and node computational complexity limitations, and by detection performance specifications. In this work, we consider a fixed set  $M$ , and formulate optimization problems for one or several monitor nodes. We fix attention to a specific monitor node and the detection scheme that it

employs. First, we need to define the quantity to be observed at each monitor. In our case, the readily available metric is the probability of collision that a monitor node experiences, namely the percentage of packets that are erroneously received. During normal network operation and in the absence of a jammer, we consider a large enough training period in which the monitor node learns the percentage of collisions it experiences as the long-term limit of the ratio of number of slots where there was collision over total number of slots of the training period. Now let the network operate in the open after the training period has elapsed and fix attention to a time window much smaller than the training period. An increased percentage of collisions in the time window compared to the learned long-term ratio may be an indication of an ongoing jamming attack that causes additional collisions. However, it may happen as well that the network operates normally and there is just a temporary irregular increase in the percentage of collisions compared to the learned ratio for that specific interval. A detection algorithm is part of the detection module at a monitor node; it takes as input observation samples obtained by the monitor node (i.e., collision/not collision) and decides whether there is an attack or not.

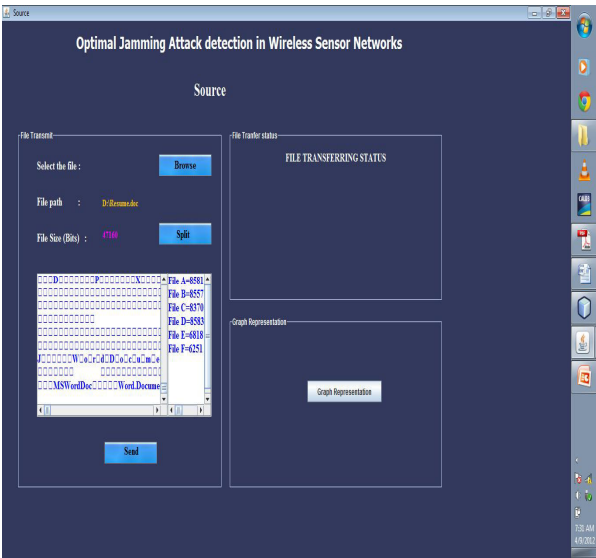
The sequential nature of observations at consecutive time slots motivates the use of sequential detection techniques. A sequential decision rule consists of: 1) a stopping time, indicating when to stop taking observations, and 2) a final decision rule that decides between the two hypotheses (i.e., occurrence or not of jamming). A sequential decision rule is efficient if it can provide reliable decision as fast as possible. The probability of false alarm PFA and probability of missed detection PM constitute inherent trade-offs in a detection scheme in the sense that a faster decision unavoidably leads to higher values of these probabilities while lower values are attained at the expense of detection delay. For given values of PFA and PM, the detection test that minimizes the average number of required observations (and thus average delay) to reach a decision among all sequential and nonsequential tests for which PFA and PM do not exceed the predefined values above is Wald's Sequential Probability Ratio Test (SPRT). When SPRT is used for sequential testing between two hypotheses concerning two probability distributions, SPRT is optimal in that sense as well.

### 3. OPTIMAL JAMMING ATTACK AND DEFENSE POLICIES AS SOLUTIONS TO OPTIMIZATION PROBLEMS

The objective of the adversary is to increase the total number of corrupted links before the attack is detected and the notification alarm is propagated. Following detection, a notification message needs to be passed out of the jammed area and, hence, the damage caused to the network is ceased. An aggressive attack, namely one with large  $q$  has a potential to corrupt more links in successive time slots. Nevertheless, this will be detected relatively quickly due to the large percentage of incurred collisions compared to the nominal one. On the other hand, a milder attack, namely one with smaller  $q$  may turn out to be more beneficial for the attacker. A significant incentive for the attacker to expedite link jamming is when jamming time-critical information. This situation is captured by the weighted cumulative payoff. As a first line of defense, the network selects the access probability  $p$  to control the number of successful transmissions given its energy limitations and at the same time expose the attacker by reducing the number of required samples to obtain a decision. Another useful network constraint is to attempt to maintain a certain minimum level of throughput in the presence of an attack.

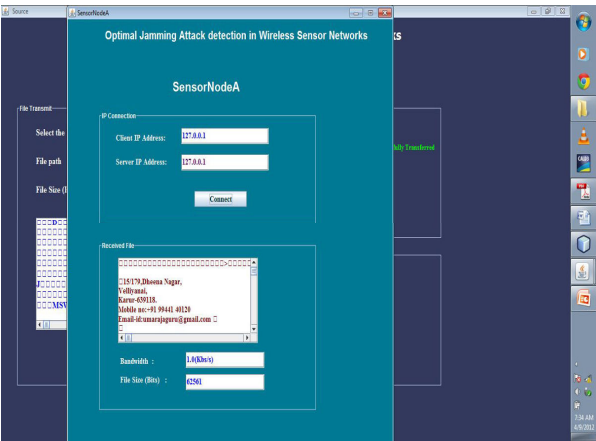
4. ACHIEVED OUTPUT

(i)Splitting of Files:

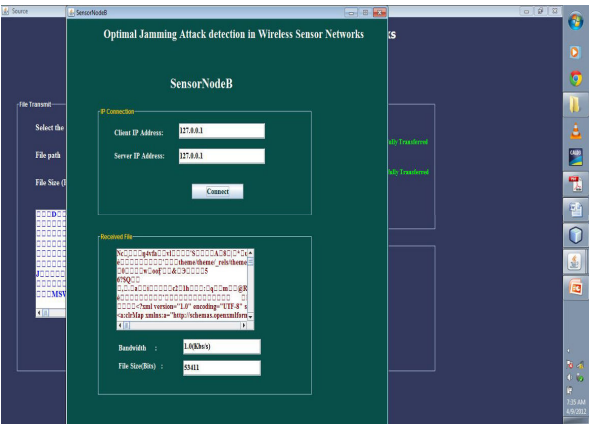


In this module the file which is to be sent is selected & split as 6 Packets based on the length of the file.  
Those 6 packets are sent through the 3 nodes between the source and destination.

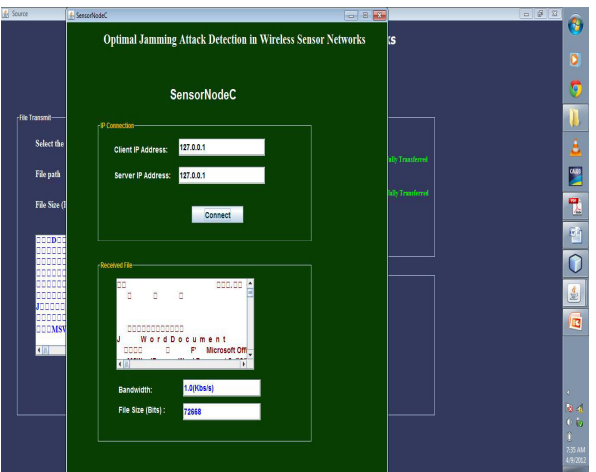
(ii)Node A Screen:



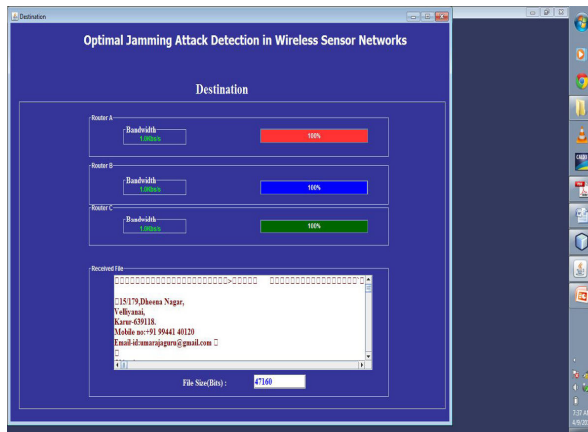
(iii) Node B Screen:



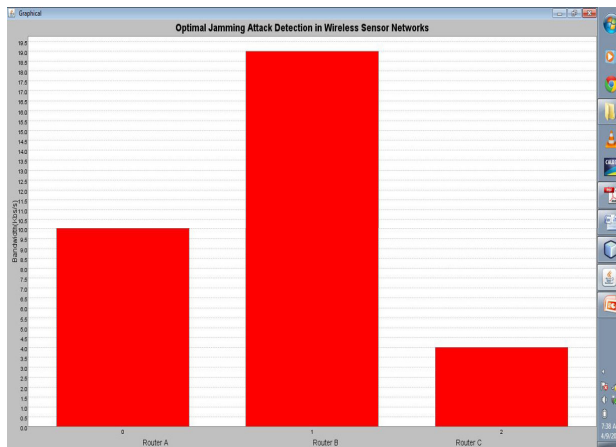
(iv) Node C Screen:



(v) Destination:



(vi) Performance Analysis:



#### 4. CONCLUSION

In this work, we studied controllable jamming attacks against wireless sensor networks, and derived the optimal solutions that dictate optimal jamming attack and network defense strategies. On one hand, the attacker attempts to find an optimal trade-off between the severity of the attack and the extent to which it becomes detectable. On the other hand, the network aims at alleviating the effect of the attack and exposing the attacker to detection. Without loss of generality, we considered an Aloha type of protocol characterized by a common access probability for all sensor nodes. The reason for adhering to this admittedly simple protocol is to abstract out the protocol specifics and focus on the collective impact of network defense (captured through a single parameter) when confronting the attack. It is understood that a similar approach can be applied when the network operates

under other channel access protocols such as CSMA that leverage more composite mechanisms such as back-off and contention window adaptation to regulate the amount of transmitted traffic. Jamming and defending strategies under these composite channel access protocols are left as a future research direction.

Finally, mobility is a dimension that gives an interesting twist in the problem and has a direct impact on network performance. In a network of mobile nodes, one would expect the detection performance to deteriorate since potential attackers move in and out of range of an observer node with a detection system, hence the sequence of observations is intermittent. In that case, interesting topics to consider would be the impact of specific mobility patterns on detection performance and how to engineer mobility patterns of defender nodes in order to alleviate the impact of attacks.

#### REFERENCES

- 1.R. Mallik, R. Scholtz, and G. Papavassilopoulos, "Analysis of an On-Off Jamming Situation as a Dynamic Game," *IEEE Trans. Comm.*, vol. 48, no. 8, pp. 1360-1373, Aug. 2000.
- 2.V. Coskun, E. Cayirci, A. Levi, and S. Sancak, "Quarantine Region Scheme to Mitigate Spam Attacks in Wireless Sensor Networks," *IEEE Trans. Mobile Computing*, vol. 5, no. 8, pp. 1074-1086, Aug. 2006.
- 3.Y.W. Law, L. van Hoesel, J. Doumen, P. Hartel, and P. Havinga, "Energy-Efficient Link-Layer Jamming Attacks Against Wireless Sensor Network MAC Protocols," *ACM Trans. Sensor Networks*, vol. 5, no. 1, pp. 1-38, Feb. 2009.
4. C.W. Helstrom, *Elements of Signal Detection and Estimation*. Prentice-Hall, 1995.
- 5.M. Cagalj, S. Capkun, and J.-P. Hubaux, "Wormhole-Based Anti-Jamming Techniques in Sensor Networks," *IEEE Trans. Mobile Computing*, vol. 6, no. 1, pp. 1-15, Jan. 2007.

# VIRTUALIZATION SECURITY FOR CLOUD COMPUTING

**G.Mamatha<sup>1</sup>**

**Assistant Professor**

**Malla Reddy Institute of Technology & Science  
Hyderabad,Telangana.**

**K.HimaBindu<sup>2</sup>**

**Assistant Professor**

**Malla Reddy Institute of Technology & Science  
Hyderabad,Telangana.**

**Abstract**— presently distributed computing is the overwhelming innovation fit for tending to the present business requests. Be that as it may, similar to all advances it has dangers and vulnerabilities, either from the cloud innovation itself or from its empowering advances. Virtualization is a primary driving innovation for distributed computing. In this administration conveyance demonstrate, the client has the capacity to adaptably make a virtual machine with the coveted necessities of working framework and programming and lets it out to the cloud supplier arrange and to the Internet. Virtualization like some other advancement has a few vulnerabilities and dangers. This paper studies virtualization dangers and vulnerabilities and a portion of the countermeasures for them.

**Keywords** — Cloud Computing, Virtualization, Virtualization Security.

## **I. CLOUD COMPUTING**

### **A. Drivers for cloud computing**

In today's business, the market change rapidly, new technologies arise and people demands' change, which demands the business to change rapidly to be able to compete in this market. This leads to several challenges on the IT to make it more agile and to keep up with the changes and demands of the business. Cloud computing address these challenges, and enable organizations and individuals to obtain and provision IT resources as a service.

Factors affecting cloud domination in this era:

- i) Now a day the hardware cost decreased while on the other hand the computing power and storage capacity are increasing.
- ii) The rapid growth in the data size in science, internet publication and archiving.

The increase in using Services Computing and Web applications. [1]

### **B. Definition**

As defined "Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction. This cloud model is composed of five essential characteristics, three service models, and four deployment models." [2]. "A large-scale distributed computing paradigm that is driven by economies of scale, in which a pool of abstracted virtualized, dynamically-scalable, managed computing power, storage, platforms, and services are delivered on demand to external customers over the Internet." [1]

### **C. Cloud Enabling technologies**

Cloud computing is not a new technology, but it is a module composed and overlaps with several other technologies to enable it.

1) **Grid computing:** Is a kind of distributed computing that allows different heterogeneous machines in a network to work together on the same task simultaneously. Grid computing enables parallel computing and is best for large workloads.

2) **Utility computing:** It is allowing different computer resources to be available for customers to use as needed and the charge depends on the usage, similar to electricity where charges are based on consumption.

3) Virtualization: It isolates the physical resources from the resources that the customer actually see they are using. It enables the resources to be viewed and managed as a pool and lets users create virtual resources from the pool. Virtualization provides better flexibility for provisioning of IT resources compared to provisioning in a non-virtualized environment. It helps optimizing resource utilization and delivering resources more efficiently.

4) Service oriented architecture (SOA): SOA provides a set of services that can communicate with each other. These services work together to perform some activity or simply passes data among services. [3]

#### D. Service Models

1) Software as a Service (SaaS): In this model the customer can use the applications provided by the cloud provider. The customer can access the applications through a web browser or program interface. Here the customer has no control on the infrastructure such as servers, operation system or storage, but only has the right to use the applications running at the provider's side with some limited application configuration the customer can change.

2) Platform as a Service (PaaS): Here the customer is provided with the platform from libraries, services and programming language tools that allow the deployment of its own application either created by itself or acquired. In this model also the customer has no control on the infrastructure such as servers, operating system or storage, but it is different from SaaS that the customer has full control on the deployed application configuration and settings and some setting of the deploying environment. [2].

PaaS is also used as an environment to develop applications, offered by the cloud service provider. The customer may use these environments to first code their applications, and then deploy it on the cloud. Because the workload to the deployed applications varies, the elasticity of the cloud resources guarantees the ability to scale in and out

transparently. [3]

3) Infrastructure as a Service (IaaS): In this model the customer is provided with the fundamental computing resources such a processing, storage and network, with which the customer is free to deploy any desired application or even the desired operating system. Here the customer has more control on the environment as he can control up to the operation system level, its storage and applications, but he still can't control the underlying infrastructure of the whole environment only that of the resources provided to him.

#### E. Deployment Models

1) Private cloud: is provisioned by a single organization, but that doesn't necessarily mean that the organization own, manage and operate the infrastructure, it can be done by a third party or a combination. That allow the infrastructure to exist off premises not only on premises.

2) Community cloud: is provisioned by a certain community of customers that consist of organizations with shared interest such as security requirements, policy or compliance considerations. One of the organizations can own, manage and operate the infrastructure or a third party or some combination. It can be on or off premises. [2]

An example in which a community cloud could be used is government agencies. If they operate under similar guidelines, they could create a community cloud to share the same infrastructure and lower their individual agency's investment.

3) Public cloud: Here where the cloud provider own, manage and operate the infrastructure, and where the infrastructure is on its premises. The provider can be business, academic or government organization or a combination. Public cloud is open for anyone to request the desired resources. [2]

Customers use the cloud services offered by the providers via the Internet and pay metered usage charges or subscription fees. An advantage of the public cloud is its low capital cost with enormous scalability. However, for customers, it has its own risks such as having no control over the resources, the security of confidential data and it being off premises, network performance, and interoperability issues.

4) Hybrid cloud: is a combination of two or more of the previous deployment models, but each remain unique entity, all tied together by standard or proprietary technology that allow data and application portability. The hybrid model allows an organization to deploy less critical applications and data to the public cloud, leveraging the scalability and cost-effectiveness of the public cloud. The organization's mission-critical applications and data remain on the private cloud that provides greater security. [3]

#### F. Cloud Computing Security

As any other technology, cloud computing face many challenges specially regarding security of the cloud. That is a major reason against adoption of cloud by many potential users. When talking about security we need to be concerned with both threats and vulnerabilities. While vulnerabilities are the faults in the system that allows attacks, threats are the potential attacks that may occur and lead to abuse of information or system resources. [4] As virtualization is a key cloud enabling technology and while it can be used as a security component assuring a higher degree of security and leveraging it, virtualization itself impose some security threats, as any security threat on virtualization, is therefore a security threat on cloud computing.

## II. VIRTUALIZATION

Virtualization is the main enabling technology for cloud computing. It means the abstraction of computer resources, and decoupling the software layer from the hardware layer, and isolating running application from used hardware. Recently virtualization started spreading on all levels (system, storage and network) as it can increase system reliability and availability. Virtualization vendors use lots of fancy names for the features of their technology, but behind all the technobabble are a number of revolutionary concepts. Take "fault tolerance" for example. When you use virtualization to pool multiple servers in such a way that they can be used as a single supercomputer, you can drastically increase uptime. If one of those servers goes down, the others continue working uninterrupted.

#### Better disaster recovery:

Data backups are much simpler in a virtualized environment. In a traditional system, you could create an "image" backup of your server complete with operating system, applications and system settings. But it could be restored to a computer only with the exact same hardware specifications.

With virtualization, images of your servers and workstations are much more uniform and can be restored to a wider array of computer hardware setups. This is far more convenient and much faster to restore compared to more traditional backups. Virtualization is done by adding a virtualization layer that consists of a hypervisor and a virtual machine monitor that is placed between the hardware and the operating system using it. Using this layer for system virtualization, a single computer machine can run multiple operating systems simultaneously within virtual machines, where the computer resources such as CPU, memory and storage are dynamically partitioned and shared between the different virtual machines. The VMM's hypervisor is responsible for implementing the virtual machine hardware and running guest OS.

#### A. Virtualization techniques

1) Full virtualization using binary translation: In full virtualization – Figure 1 - depending on the type of instruction that should run, the decision is made on how to handle it. If it is a kernel level instruction the nonvirtualizable instruction are replaced with new ones that gives the same result on the virtual hardware, while user level instructions are executed on the processor directly to help improve the performance.

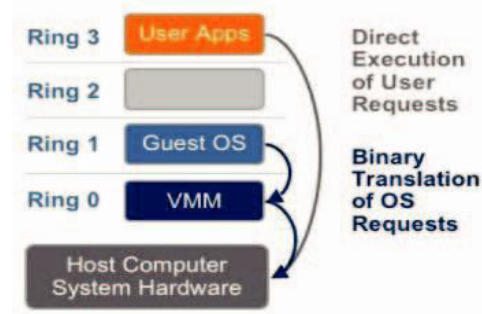


Fig. 1. Full virtualization using binary translation

This translation of kernel level instructions is done by the hypervisor on the fly with suitable instructions to run on virtual hardware, and the hypervisor cache the translation result so it can use it

later, with no hardware or operating system intervention needed. The x86 architecture has 4 levels of privileges for running different instructions, operating system and user applications each need different privilege to run that is divided into 4 categories known as Rings 0, 1, 2 and 3. In a non-virtualized environment operating system runs in ring 0, while user applications run in ring 3. In virtualized environment this change according to the technique used, that is shown in the different figures. It has high level of security and isolation for the virtual machine that allow it to be portable and flexible to run on virtual hardware or directly on physical hardware as there is no changes done to the guest OS. [5] [6]

2) OS assisted virtualization or paravirtualization: While full virtualization doesn't modify the guest OS, paravirtualization – Figure 2- include modifying the kernel of the guest OS to substitute kernel level instructions with hypercalls that is convey directly to the hypervisor without further need to translate this instruction to run it. Hypercall interfaces are provided by the hypervisor for other kernel operations like time keeping, memory management and interrupt handling. [5] [6]

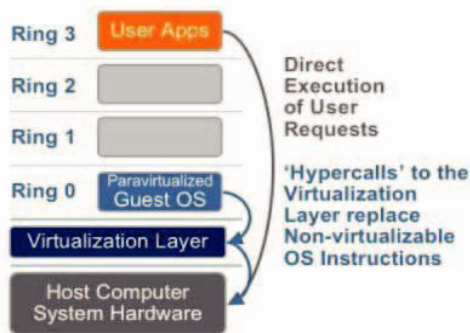


Fig. 2. OS assisted virtualization.

### 3) Hardware assisted virtualization (first generation):

Hardware vendors started enhancing the hardware introducing Intel Virtualization Technology (VT-X) and AMD-V from AMD to add a new execution mode in the CPU with privileged instructions that permit the VMM to run in a mode below ring 0 as shown in figure 3. These privileged instructions are automatically trapped to the hypervisor with no need for translation or changing the OS. [5] [6]

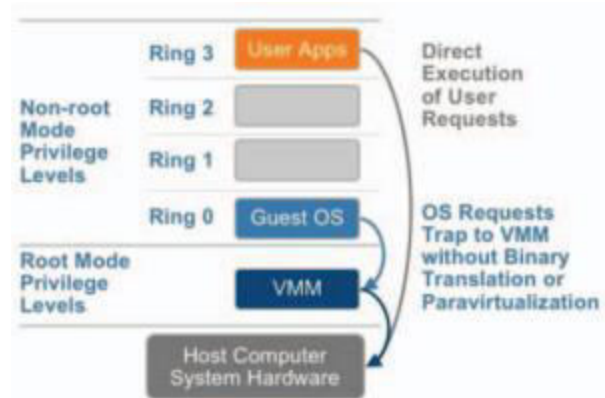


Fig. 3. Hardware assisted virtualization.[5]

## B. VM Services

Virtualization provides some tools to facilitate and help in the administration, and provide some services. One of the services is migration that allows VM to be transferred between physical machines, in case one physical machine

became down, this will allow the work on the VM to continue without waiting for the physical machine. Using virtualization, new VMs can easily be created via images file format, which is a package template. A useful feature of virtualization is clipboard that permits data transfer between guest and host machines. Great care should be taken when using these tools and services, that they don't impose security vulnerability.

## C. Virtualization vulnerabilities and Threats

### 1) VM Escape:

#### • Problem definition:

As illustrated in figure 4 VM escape is where an application running on a VM can directly have access to the host machine by bypassing the hypervisor, being the root of the system it makes this application escape the VM privilege and gain the root privilege. That means it can access resources and gain control over it and all the VMs in the system and the hypervisor, can change assigned resources to VMs (CPU, memory, and disk) or it can shut down or restart a VM or even the hypervisor. Normally VMs are isolated from the host, it can only share resources. A program can't monitor, communicate or affect other VMs or the host. But

due to some organizations configuring flexible isolation meet its needs, or software bugs that may compromise isolation. [4] [8] [9]

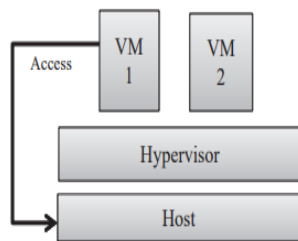


Fig. 4. VM Escape.

- Countermeasures:

That can be solved by properly configuring host and guest interaction rules with minimum flexibility. [8] [9]

Other countermeasures include HyperSafe which is an approach that provides hypervisor control-flow integrity. Trusted cloud computing platform (TCCP) allows providers to propose closed box execution environments, and allow customers to decide before launching VMs if the environment is secure. Trusted virtual datacenter (TVD) guarantee separation and integrity in the cloud[4]are stored and gives a chance to the host to monitor the logs of encrypted terminal connections inside the VMs. [9]

2)Denial of Service:

- Problem definition:

Denial of service attack is an effort to make resources unavailable to the users. It is mainly used with computer networks, but it is not limited to it, it is also used in reference to other resources such as CPU resource management. This can be done by flooding the machine with communication requests which consume its resources that leads it to being unable to answer rightful traffic or answer slowly to be rendered unavailable. In virtualization due to sharing physical resources such as CPU, memory and network between VMs and host, a VM can impose the denial of service attack on other VMs on the same physical machine by taking and consuming all the possible resources of the system, making them unavailable to other VMs. [4] [9]

- Countermeasures:

This can be prevented by allocating limited resources to

each VM, as virtualization technologies provide mechanisms to limit the allocation of resources to individual VM proper configuration for limiting the resource allocation must be considered. [4] [9]

3) Attack on VM at migration:

Problem definition:

Even though features are supposed to help facilitate the work, sometimes it may impose security risk itself.VM migration is one of the virtualization features that allow VMs to be transferred from one physical machine to another in a standby or active state. Through migration VMs' data need to be protected from unsecured network, or attack. For example a malicious user can start or redirect the migration process to a different network in which he has access or untrusted host, or it can just be copied and used elsewhere, which compromise the VM with the passwords, credentials on it and in case of coping it makes it difficult to trace the attacker. [4] [8]

4)Attack on host machine:

Problem definition: The host machine is the control point and root of the virtual system; it can monitor and communicate with the applications running on the VMs. Depending on the virtualization technology, it has different implications in the way of interaction between host machine and the VMs. Possible ways for that are:

- The host can start, shutdown, restart or pause a VM.
- The host is able to monitor and modify the resources gives to the VM.
- With enough rights given it can monitor the application running on a VM.
- View, copy and modify the data stored on the virtual disk assigned to a VM.

All network traffic of the VMs passes through the host, which enable it to monitor each VM network traffic. In some VM technologies keystrokes and screen updates can be logged by the VMM, giving that the host have the necessary permissions, these

- Countermeasures:

The host machine is the control point of the virtual system; therefor it must be strictly protected more than the VMs itself. Isolation of the host should be

provided in order not to allow it to be the gateway for attacks on the VMs. [7] [10]

#### 5)VM Sprawl:

- **Problem definition:** Sprawl problems is not new in the IT world with virtualization, but as creating something becomes easier the more number you get from it and the harder to clean it up, such as email, messages and files. [11] VM sprawl happens when resources are being consumed by unused VMs such as CPU, memory and storage. This can happen either because the VM was provisioned unnecessarily without proper justifications and approvals, or it was ordered, used extensively and then left inactive for prolonged periods, or it was over-provisioned from the start (too much CPU, memory, or disk). This is due to the ease of requesting a new VM, the speed and ease of creating new VM without the complication of the physical word,from limitation at the creation time of a VM. [11] [13] requesting, approvals, ordering machines, arrival, getting it racked and operational which is time consuming that may take weeks or even months.in virtualization this process is reduced to hours or even minutes. [11] [12] [13]

- **Countermeasures:** There are solutions for sprawl some are used and some are suggested and under test. Mainly policies and governance that allow administrators to control what and how many resources can be consumed, and adding optional approval workflow to further ensure the need to provision a machine with the correct resources. Also policies are needed to define what to do after a VM has expired or no longer needed, at this point VMs won't be consuming CPU or memory, but they will need to be archived for a period of time and that will consume storage, the police will govern for how long to keep this archive and when it can be cleaned up, same for snapshots that by time grow and consume more storage. Also automation can clean up and recycle by identifying resources and reuse it, and reduce the number of unnecessarily provisioned or over-provisioned machines.another solution is specifying time.

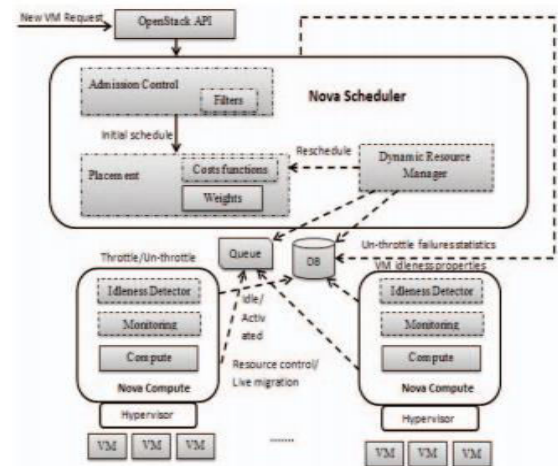


Fig. 5. Utilization Accelerator (PULSAR).

One of the presented solutions is adaptive Utilization Accelerator (PULSAR) - Figure 5 - which is a resource scheduler that spread out OpenStack Nova filter scheduler. It focuses on CPU over-commit by changing the over-commit ratio (OCR) dynamically. While increasing OCR is a way to mitigate sprawl it may have a disadvantage as it doesn't consider possible performance degradation that it might cause if it is static unlike PULSAR. [12]

- Patching vulnerability:
- Problem definition:

With the increased of VMs in different locations, it becomes difficult to manage the OS and application security patches for all VMs. And it becomes more difficult if a VM rolled back to an earlier snapshot, which may leave the system with vulnerabilities that has known solutions so system administrators needs to identify the patches and apply them again. Also for ideal VMs during patching, the latest patches on them need to be checked once they become active and needed patched identified and applied. [10] [14]

- Countermeasures:

There are different patch management systems available. A strong patching management system (example Shavlik) that distributes the patches to the different machines, monitors all the patching level, manages information about installed patches, and controls the patch dependency of the

VMs is needed. Figure 6 shows the schematic view of a patch management framework. [15]

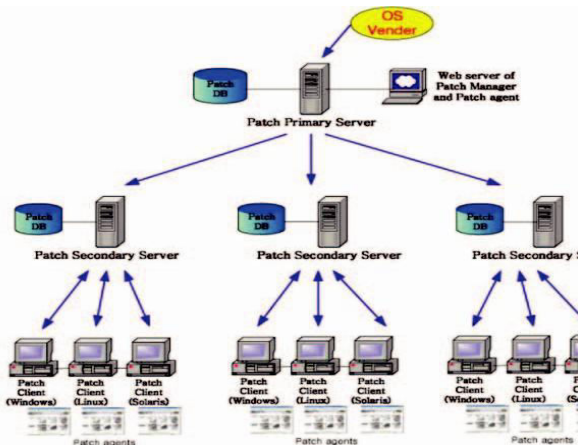
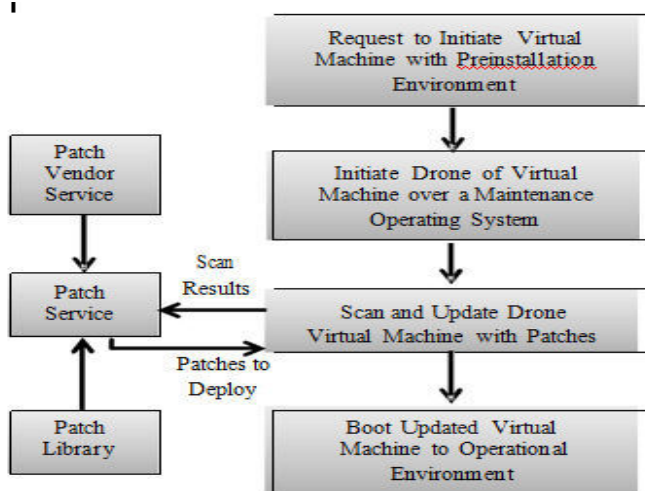


Fig. 6. Patch management framework

In addition to using patch management systems, checking the patching compliance for drone machines, validation of successful updating and no missing patches is needed. This can be done when initiating a drone machine to be in a maintenance environment first to check its patch compliance, then patching is done if it is not according to standard and changes are saved, then the machine can be initiated in the operation environment. Figure 7 shows checking patch



flowchart. [16]

Fig 7.Patch checking flowchart

#### 6) Infected VM Images:

- Problem definition:

VM image is a type of file format used to create VMs, a prepackaged software template that contains the configuration files for a VM, which make its confidentiality and integrity an important issue, as a vulnerable image run by the user may compromise the whole system security. The security risks of VM images has three perspectives, first the publisher's risk by exposing sensitive information on the image like passwords on configured applications, and wanting to share the image with only a limited number of users. Second is the retriever's risk of running a VM image which may contain malicious software, or illegal and unlicensed software. Third is the repository administrator's risk, who is most of the time also the cloud administrator that is the security of the dormant VM image, as vulnerabilities are exposed over time. And the risk of hosting a distributing images that are malicious.[14] [17]

- Countermeasures:

Mirage image management system is suggested with four components to address each of the security problems.

- Access control framework to control who can use this image.
- Filters to remove publisher sensitive information.
- Track the history of an image using provenance tracking mechanism.
- Repository maintenance service, like periodic virus scanning. [17]

#### 7) VM Hopping:

- Problem definition:

If a VM gains access over another VM running on the same host, for example by exploiting

some hypervisor vulnerability, it can control the other VM as illustrated in figure 8. Also as mentioned before in the migration even though features are supposed to help facilitate the work, sometimes it may impose security risk itself, here clipboard in VM is a feature that allows data transfer between VMs and host. As it may be useful but it can be a gateway to transfer data between malicious VMs. [4] [8]

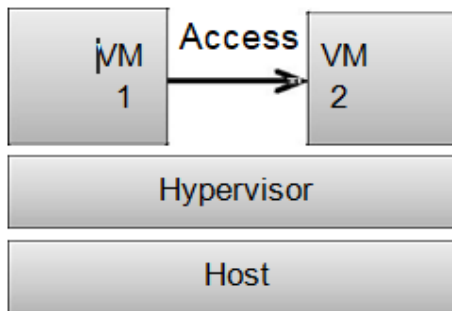


Fig. 8. VM Hopping

Also if an attacker could co-reside his VM, by considering the number of physical machines and the probability of co-resident according to the allocation policy, the attacker starts a minimum number of VMs that gives high probability of residing on the same physical machine with the VM targeted, it can obtain sensitive data using side channels. [18]. The idea of side channel is that although virtualization can be configured to ensure strong isolation between different VMs, this isolation is logical and some hardware components are still share, like using memory caches (L2) as side channel. [19]

- Countermeasures:

This could be avoided if the attacker couldn't co-reside on the same host as the VM that he needs to attack, this needs the allocations policies to be considered by adding a modification that when a user request a new VM, the machines that already holds that user VMs are considered first. This will make a user VMs' together, which will decrease a chance for an attacker VM to co-reside, even if he created multiple VMs they will be on the same machine. [18] Side channels itself can be a detection tool used for defense to detect undesired co-residency, by coordinating with other friendly VMs on the same machine to silence their activity in a certain cache

region for specific time period, and during that time the cache usage is measured, and the result is analyzed for any unexpected behavior that may indicate that there is an unfriendly VM on the same machine. [19]

### III. SECURITY SERVICES

How different virtualization vulnerabilities affect security services and the estimate attack source are illustrated in table 1.

TABLE I. SECURITY SERVICES

|                           | Security services    |                |              |                | Attack source | Solution                                                                                                                 |
|---------------------------|----------------------|----------------|--------------|----------------|---------------|--------------------------------------------------------------------------------------------------------------------------|
|                           | Data Confidentiality | Data Integrity | Availability | Access Control |               |                                                                                                                          |
| VM escape                 |                      |                | ✓            | ✓              | U             | <ul style="list-style-type: none"> <li>• Proper configuration</li> <li>• Safe hypervisor</li> </ul>                      |
| Denial of service         |                      |                | ✓            |                | U/P           | <ul style="list-style-type: none"> <li>• Apply security policy</li> </ul>                                                |
| Attack on VM at migration | ✓                    | ✓              |              |                | U/P           | <ul style="list-style-type: none"> <li>• Consider network security</li> </ul>                                            |
| Attack on host machine    | ✓                    | ✓              | ✓            |                | U             | <ul style="list-style-type: none"> <li>• Proper configuration to isolate host and VMs</li> </ul>                         |
| VM sprawl                 |                      |                | ✓            |                | U             | <ul style="list-style-type: none"> <li>• Creating VM policies</li> <li>• New algorithm for resource scheduler</li> </ul> |
| Patching vulnerability    | ✓                    | ✓              | ✓            | ✓              | P             | <ul style="list-style-type: none"> <li>• Patch management system</li> </ul>                                              |
| Infected VM images        | ✓                    | ✓              |              |                | U             | <ul style="list-style-type: none"> <li>• Image management system</li> </ul>                                              |
| VM hopping                | ✓                    | ✓              |              |                | U             | <ul style="list-style-type: none"> <li>• Modification to allocation policies</li> </ul>                                  |

\*U: attack source is the cloud service user

P: attack source is the cloud service provider

### IV. CONCLUSION

In this paper an overview of cloud computing module was given, focusing on virtualization which is one of its enabling technologies, surveying its vulnerabilities and threats, that may be an obstacle in front of further adopting cloud computing, and illustrating some of the countermeasures that are used against these threats and vulnerabilities. Some

of these countermeasures are already adopted and some are still under testing and development.

These vulnerabilities, countermeasure and future work are:

- VM escape: where VM can completely bypass the hypervisor. Proper configuration of the host and guest interaction rules can solve it, or using a safe hypervisor like HyperSafe, or a closed box execution environments like TCCP.
- Denial of service: similar to network DOS making computer resources unavailable can be prevented by enforcing policies to allocate limited resources to each VM.
- Attack on VM at migration: VMs with their data may be attacked while migrating from one physical machine to another. It needs Network security mechanism and policies should be taken into account, like communication channel security, encryption.
- Attack on host machine: strictly protecting the host as it can control the VMs. Also sufficient isolation should be provided in order to not allow the host to be the gateway for attacks on the VMs.
- VM sprawl: resources are continuously consumed by unused VMs. Mainly policies and governance that allow administrators to control the resource usage are used as countermeasure. Other algorithms are present too like Utilization Accelerator (PULSAR) which is a resource scheduler that spreads Open Stack Nova filter scheduler.
- Patching vulnerability: managing the OS and application security patches for all VMs. A patch management system is necessary, along with an algorithm to check the patching compliance for drone machines.
- Infected VM images: Controlling images as there may be some from malicious user where malicious code may be present or unlicensed software or publisher sensitive information

exposed. Mirage image management system is a suggested solution that uses tracking mechanism, filters, access control and maintenance services.

- VM hopping: where a VM gains access over another VM running on the same host. Modification to the allocations policies are introduced as a countermeasure.

Although some of these are already applied methods and some are still new algorithms, all can be further enhanced to address these vulnerabilities and new challenges that may arise.

## REFERENCES

- [1] EMC, Information Storage and Management. EMC, 2012, pp. 404-429.
- [2] Ian Foster, Yan Zhao, Ioan Raicu, and Shiyong Lu, "Cloud Computing and Grid Computing 360-Degree Compared," in Grid Computing Environments Workshop, 2008.
- [3] Peter Mell and Timothy Grance, "The NIST definition of cloud computing," 800-145, 2011.
- [4] Keiko and Rosado, David G and Fernandez-Medina, Eduardo and Fernandez, Eduardo B Hashizume, "An analysis of security issues for cloud computing," Journal of Internet Services and Applications, vol. 4, 2013.
- [5] IBM, IBM Systems Virtualization: Servers, Storage, and Software.: IBM, 2008.
- [6] VMware Inc, "Understanding full virtualization, paravirtualization, and hardware assist," 2007.
- [7] Sahoo, Jyotiprakash and Mohapatra, Subasish and Lath, Radha, "Virtualization: A survey on concepts, taxonomy and associated security issues," 2010.
- [8] Amarnath and Shah, Payal and Nagaraj, Rajeev and Pendse, Ravi Jasti, "Security in multi-tenancy cloud," in Security Technology (ICCST), 2010 IEEE International Carnahan Conference on, San Jose, CA, 2010, pp. 35-41.
- [9] Joel Kirch, "Virtual machine security guidelines," 2007.
- [10] Shengmei and Lin, Zhaoji and Chen, Xiaohua and Yang, Zhuolin and Chen, Jianyong Luo, "Virtualization security for cloud computing

- service," 2011.
- [11] VMware Inc, "Controlling Virtual Machine Sprawl, How to Better Utilize Virtual Infrastructure," White Paper 2012.
- [12] Breitgand, David and Dubitzky, Zvi and Epstein, Amir and Feder, Oshrit and Glikson, Alex and Shapira, Inbar and Toffetti, Giovanni, "An adaptive utilization accelerator for virtualized environments," in Cloud Engineering (IC2E), 2014 IEEE International Conference on, Boston, 2014, pp. 165 - 174.
- [13] Maik and McDonald, Fiona and McLarnon, Barry and Robinson, Philip Lindner, "Towards automated business-driven indication and mitigation of VM sprawl in Cloud supply chains," in Integrated Network Management (IM), 2011 IFIP/IEEE International Symposium on, Dublin, 2011, pp. 1062 - 1065.
- [14] Qian and Mehrotra, Rajat and Dubey, A and Abdelwahed, Sherif and Rowland, Krisa Chen, "On state of the art in virtual machine security," in Southeastcon, 2012 Proceedings of IEEE, Orlando, FL, 2012.
- [15] Jung-Taek and Choi, Dae-Sik and Park, Eung-Ki and Shon, Tae-Shik and Moon, Jongsub Seo, "Patch management system for multi-platform environment," in Parallel and Distributed Computing: Applications and Technologies.: Springer Berlin Heidelberg, 2005, pp. 654--661.
- [16] Campbell McNeill, "Virtual Machine Asynchronous Patch Management," 14/162,202, January 23, 2014.
- [17] Jinpeng and Zhang, Xiaolan and Ammons, Glenn and Bala, Vasanth and Ning, Peng Wei, "Managing security of virtual machine images in a cloud environment," in Proceedings of the 2009 ACM workshop on Cloud computing security, 2009, pp. 91-96.
- [18] Yi and Chan, Jeffrey and Alpcan, Tansu and Leckie, Christopher Han, "Virtual machine allocation policies against co-resident attacks in cloud computing," in Communications (ICC), 2014 IEEE International Conference on, Sydney, NSW, 2014, pp. 786-792.
- [19] Yinqian and Juels, Ari and Oprea, Alina and Reiter, Michael K Zhang, "Homealone: Co-residency detection in the cloud via side-channel analysis," in Security and Privacy (SP), 2011 IEEE Symposium on, Berkeley, CA, 2011.

# ***DISTRIBUTED TRACKING SYESTEM***

**Ajju P. Benny**  
CSE Department  
Aju5101@gmail.com  
Malla Reddy College of Engineering

***Dr.V.Bhoopathy***  
*Professor, CSE Department*  
v.bhoopathy@gmail.com  
Malla Reddy Engineering College

## **ABSTRACT:**

The Project basically deals with the idea to track the goods, which are sent by different people to their choice of destination. So as to have a control and clear cut information system to the companies involved in this kind of business as well as to the customers those who make use of these services. The companies like DLF and Safex are involved in logistic Business, where they are involved in thorough dispatch of goods to the correct destination with in a stipulated time. These companies can make use of this software or it can be used by any smaller companies, which will have a direct impact on the profit of business happening.

The companies can use this to track and check whether the goods are in process of dispatch or got strucked at place. The customer gets an idea about how and when the packet he/she will receive.

## **1. Introduction:**

These companies can make use of this software or it can be used by any

smaller companies, which will have a direct impact on the profit of business happening. The companies can use this to track and check whether the goods are in process of dispatch or got stroked at place. The customer gets an idea about how and when the packet he/she will receive.

Functional Requirements:

The various aspects associated are

- Design of user interface/Login form
- Validating user data
- Storing the same in database
- Providing unique number to customer.

## **2. Related to work:**

### **2.1 Existing System**

In the existing system the problem faced by the staff and customers are a lot. In this system the tracking of the product and the responsibility of the products received and sent are not fixed and can lead to loss of products and goods sent. The billing and the timeliness is also not accurate. The

customers has to suffer because not getting on-line details or information with respect to the goods sent by him.

## 2.2 Proposed System

So in order to see that the hassles are reduced almost to zero the software module is designed . This module helps people to maintain the details of goods sent and received accurately and perfectly. The responsibility can also be fixed here. The customer can lead to happy life once the goods are handed over to the distribution people. The customer can keep track of the product sent by him by using the internet. Where he can get details like where the product is , when it will reach the destination mention .The customer can hook on to internet and can get details by using the order no given to him by the company people.

## 2.3 The modules of the project

The first module deals with the front-end design part.

The Second module deals with the process of data taken from module one and the same has to be stored in database.

The third module deals with the reporting system where different type of reports can be generated.

- User-Interface module/Login Module.
- Validations check module/Business Logic.
- Data services module.

## 2. System Architecture:

Applications are developed to support companies in their business operations. Applications take data as input, process the data based on business rules, and provide data or information as output. Based on this fact, all applications will have 3 elements

- The user interface or the presentation element, through which data is taken as input.
- The application logic or the business rule element, which helps in implementing the operations to be performed on the input data.

## 3.1 Java Architecture:

### 3.1.1 Two-Tier Architecture:

In a traditional 2-tiered application, the processing load is given to the client PC while the server simply acts as a traffic controller between the application and the data. As a result, not only does the application performance suffer due to the limited resources of the PC, but the network traffic tends to increase as well. When the

entire application is processed on a PC, the flexibility to the design of the application is forced to make multiple application. Multiple user interface requests for data before even presenting can be built and deployed without ever anything to the user. These multiple changing the application logic, database requests can heavily tax the provided the application logic presents a clearly defined interface to the network.



Fig 3.1.1:two-tier architecture

### 3.1.2 Three -tier Architecture:

The first tier is referred to as the presentation layer and typically consists of a graphical user interface of some kind. The middle tier, or business layer, consist of the application or business logic, and the three tier – the data layer – contains the data that is needed for the application.

The middle tier (application logic) is basically the code that the user calls upon (through the presentation layer) to retrieve the desired data. The presentation layer then receives the data and formats it for display. This separation of application logic from the user interface adds enormous

The third tier contains the data that is needed for the application.



fig 3.1.2:three-tier architecture

### 3.1.3N-tierarchitecture breaks down like this:

- ❖ A user interface that handles the users interaction with the application – this can be a web browser running through a firewall, a heavier desktop application, or even a wireless device.

- ❖ Presentation logic that defines what the user interface displays and how a users requests are handled. Depending on what user interfaces are supported, you may need to have slightly different versions of the presentation logic to handle the client appropriately.
- ❖ Business logic that models the application's business rules, often through interaction with the application's data.
- ❖ Infrastructure services that provide additional functionality required by the application components, such as messaging, transactional support.
- ❖ The data layer where the enterprise's data resides.

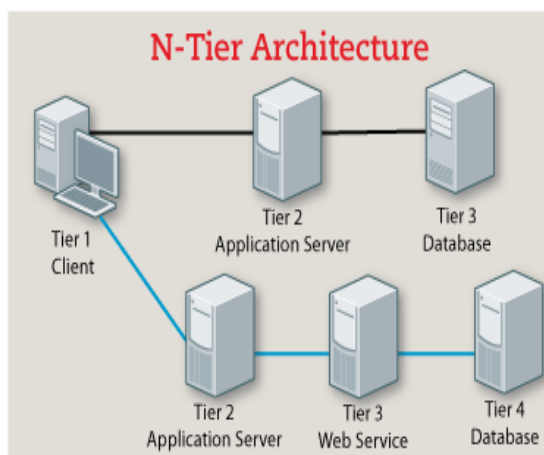


fig 3.1.3: n-tier architecture

## 4. Implementation:

A crucial phase in the system life cycle is the successful implementation of the new system design. Implementation includes all those activities that take place to convert from the old system to the new one. The new system has been implemented and many users have come forward to learn the operation of the new system. Many users have the system and found it to be very useful and efficient in all respects except some.

Implementation the process of having systems personnel checkout and put new equipment into use,train users,install the new application and construct any files of data needed to use it. This phase is less creative than system design. Depending on the size of the organization that will be involved in using the application and the risk involved in its use,system developers may choose to test the operation in only one area with only one or two persons.

Implementation becomes necessary so as to provide a reliable system based on the requirements of the organization. Successful implementation may not guarantee improvement in the organization using the new system,but improper

installation will prevent it. It has been observed that even the best system cannot show good result if the analysts managing the implementation do not attend to every important details. This is an area where the users need to work with utmost care.

Even well designed system can succeed or fail because of the way they are operated and used. Therefore, the quality of training received by the users involved with the system in various capacities helps and may even prevent the successful implementation of the management information system. Those who are directly or indirectly related with the system development work must know in detail what their roles will be, how they can make efficient use of the system and what the system will or will not do for them.

## 6.Conclusion:

The efficiency of any system designed to suit an organization depends on cooperation during the implementation stage and also flexibility of the system to adopt itself to the organization.

With the help of “Distribution Tracking system” it becomes easy for the company and the staff along with

the customers to have the control on the system which is executed on day today basis. The software designed will help them to fix responsibilities on the staff related and can also provide good service to the customer from time to time. Initially there will be a little amount of hesitancy from the staff to make use of the newly designed software, but as time progresses and once they are habituated to the system, then they will appreciate the need for this kind of softwares. It is also very important to see that the staff are trained properly on the newly designed software and some trials can be conducted before we ask to execute it on-line.

## 7.References:

Complete reference java, 2<sup>nd</sup> edition by Herbert Schlitz.

SERVLET programming, by Our'Reilly

Servlet programming, 2<sup>nd</sup> edition by Karl moss

HTML Black book by Steven Holzner

Core Java Foundation Class by Kim Topley

The JDK 1.4 tutorial by Greg Travis

[www.sunmicrosystem.com](http://www.sunmicrosystem.com)

## ***A Technique in Communication with cloud using RPC***

***K.Bharath,***  
Korrollobharath22@gmail.com  
Malla Reddy College of Engineering

***Ch.Vijaya Kumari***  
Malla Reddy College of Engineering,  
hodcse@mrce.in

### **ABSTRACT**

Cloud computing emerging area for distributing applications and accessing application through remote procedure calls are helpful to request data and sending data to clouds. In cloud computing there is lot of mechanisms to transfer data like message queues but by using remote procedure calls can access information remotely without having independent failures. Data available at the source and the consumer can interact with the system. However security is one of the criteria to apply lot of security algorithms for data exchange in between source and the user. The cloud server can search the user data by request through remote procedure call and find out the information by using index and corresponding files. The functions contained within RPC are accessible by any program that must communicate using a client/server methodology.

### **KEYWORDS**

Remote Procedure calls, Data security, Network protocols, Cloud Platforms.

### **1. INTRODUCTION**

Cloud computing is the advanced technology for data storage and describes the web as a platform for accessing information. To users cloud is the pay per use on demand to use resources through internet. Cloud models are helpful to exploring the view of the end user.

**Public Cloud:** is a type of cloud hosting in which the cloud services are delivered over a network which is open for public usage. This model is a true representation of cloud hosting; in this the service provider renders services and infrastructure to various clients. The customers do not have any distinguish ability and control over the location of the infrastructure.

**Private Cloud:** is also known as internal cloud; the platform for cloud computing is Implemented on a cloud-based secure environment that is safeguarded by a firewall which is under the governance of the IT department that belongs to the particular corporate. Private cloud as it permits only the authorized users.

**Hybrid Cloud:** is a type of cloud computing, which is integrated. It can be an arrangement of two or more cloud servers, i.e. private, public or community cloud that is bound together but remain individual entities. Benefits of the multiple deployment models are available in a hybrid cloud hosting. A hybrid cloud can cross isolation and overcome boundaries by the provider; hence, it cannot be simply categorized into public, private or

community cloud. It permits the user to increase the capacity or the capability by aggregation, assimilation or customization with another cloud package / service.

In cloud computing, cloud services are available such as software as a service, platform as a service and infrastructure as a service are basic level services useful in user environments. Software as a service is to allow the user to access software applications. Platform as a service is to provide the deployment tools to the end user. Infrastructure as a service is to access the resources such as virtual storage.

Benefits of cloud computing are cost effective, on demand service, resource available on network, online deployment models, high efficiency and high reliability

## 2. REVIEW ON CLOUD COMPUTING

Cloud computing is a computing based on the internet. In cloud computing resources are available on the internet side it allows the user to access the applications. The journey to the cloud is for accessing the application and running it for flexibility, Automatic software updates, document control and the security.

Cloud platform: Elastic Cloud Computing Platform it provides a programmable virtual cloud infrastructure automated design, deployment and management of virtual applications in the cloud and configure the cloud capacity in an easy to use.

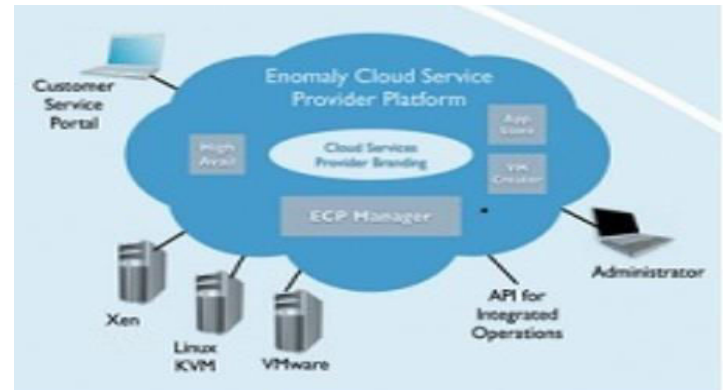


Figure 2: Elastic Cloud Computing Platform

### Why cloud Computing

The cloud computing has leading improvements in the storage and accessing resources remotely. Immediate updates with new features and functionality of software enhancements are available in cloud computing environments. To reduce the size of the own data centers and reducing the servers count without impacting the capabilities for accessing information cloud computing is useful. In the traditional computing requires buying capacity it is sufficient to the end user in cloud environments.

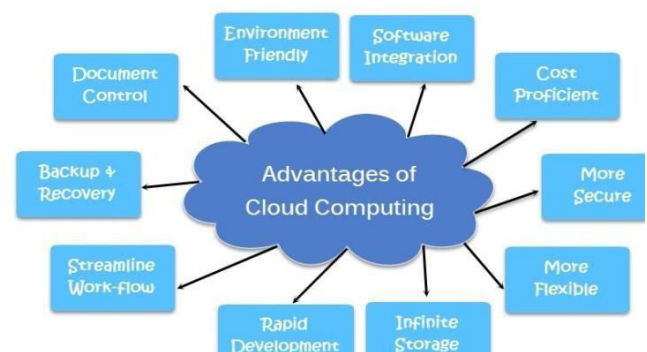


Figure 3: Advantages of Cloud Computing

3. AN EVOLUTIONARY APPROACH IN RPC

RPC Processes and Interactions in the cloud

The RPC components make it easy for clients to call a procedure located in a remote server program. The client and server each have their own address spaces; that is, each has its own memory resource allocated to data used by the procedure. The following figure shows the RPC process.

RPC Process

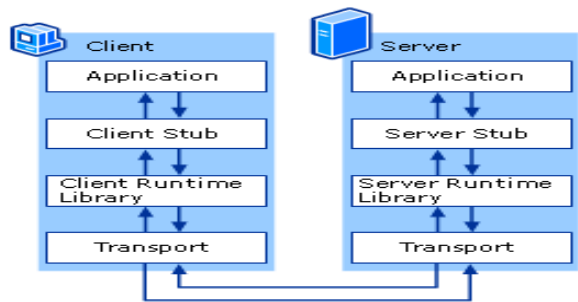


Figure 4: RPC Process between Client and Server

The RPC process starts on the client side. The client application calls a local stub procedure instead of code implementing the procedure. Stubs are compiled and linked with the client application during development. Instead of containing code that implements the remote procedure, the client stub code retrieves the required parameters from the client address space and delivers them to the client runtime library. The client runtime library then translates the parameters as needed into a standard Network Data Representation (NDR) format for transmission to the server.

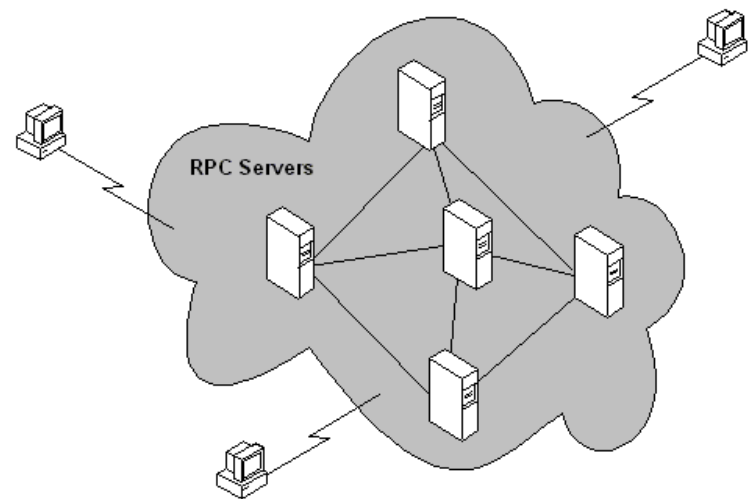


Figure 5: RPC SERVERS

**RPC-Supported Network Protocols in the Cloud:** In the Cloud Computing environment protocols are supported to provide the services to request and response in between the end user and the cloud. Some of the protocols are listed on bellowed table.

| Types of protocols                 |          | RPC Type            |
|------------------------------------|----------|---------------------|
| Transmission Control Protocol(TCP) | Control  | Connection-oriented |
| Sequenced Exchange (SPX)           | Packet   | Connection-oriented |
| Named Pipe                         |          | Connection-oriented |
| HTTP                               |          | Connection-oriented |
| User Datagram Protocol (UDP)       |          | Connectionless      |
| Cluster Datagram Protocol (CDP)    | Datagram | Connectionless      |

#### 4. CONCLUSION

In this paper we proposed the RPC mechanism in the cloud environments and process between the client and the server communication process. RPC is an evolutionary approach to access the information remotely by using the clouds environment. The cloud engaged with the requests and the process for delivering the information by searching on the storage by using index mechanisms. The distributed link and tracking client provide reliable connection between the client and the server.

#### 5. REFERENCES

- [1] Ren, Yulong, and Wen Tang. "A SERVICE INTEGRITY ASSURANCE FRAMEWORK FOR CLOUD COMPUTING BASED ON MAPREDUCE." *Proceedings of IEEE CCIS2012*. Hangzhou: 2012, pp 240 – 244, Oct. 30 2012-Nov. 1 2012
- [2] N, Gonzalez, Miers C, Redigolo F, Carvalho T, Simplicio M, de Sousa G.T, and Pourzandi M. "A Quantitative Analysis of Current Security Concerns and Solutions for Cloud Computing.". Athens: 2011., pp 231 – 238, Nov. 29 2011- Dec. 1 2011
- [3] Hao, Chen, and Ying Qiao. "Research of Cloud Computing based on the Hadoop platform.". Chengdu, China: 2011, pp. 181 – 184, 21-23 Oct 2011.
- [4] Y, Amanatullah, Ipung H.P., Juliandri A, and Lim C. "Toward cloud computing reference architecture: Cloud service management perspective.". Jakarta: 2013, pp. 1-4, 13-14 Jun. 2013.
- [5] Devarakonda Krishna "A Safety Summons on Cloud Data" IETE- Chandigarh, PP:340-343, 17-18 29 July. 2017
- [6] BIRRELL, A. D., AND NELSON, B. J. Implementing remote procedure calls. *ACM Trans. Comput. Syst.* 2, 1 (Feb. 1984), 39-59.

## ***BLUE DOG***

**D. Namratha**

CSE Department

Namrathadonti95@gmail.com

Malla Reddy College of Engineering

***Dr.Chandra Shekar***

*Professor, CSE Department*

drchandru86@gmail.com

Malla Reddy College of Engineering & Technology

### **ABSTRACT:**

Nowadays our beloved ones security issues are getting increased, as they are moving from one location to another. In order to resolve the problem we have developed an app, when ever the person whose mobile number is registered with our application not answering our phone call, it will help to send the person location information to us within 5 minutes and he can view the exact location where he is with the help of Google maps. You need to add phone numbers of your guardian in this app. By default if any one of a guardian gives you missed calls your location is sent to that guardian only.

### **1 Introduction:**

This is an application which is used for security purpose. This is an android application. Here we will be having guardians and users, where users add other people as his guardians. When the user does not lift the phone, the details of his location are sent to the guardian. Here the address will be sent using which we can

locate it in the map. And This application should be installed on both the guardian's and user's phone. So, only the guardian will receive the location details of the user if he does not lift the phone but not when the user calls and when the guardian does not lift the phone

#### **1.1 Objective:**

- Here the main objective is to provide security for people.
- The main purpose is to send the location of user to the guardian.
- We can get the location of the person when the person does not lift the phone.
- We can even see the map regarding the location.

#### **1.2 Scope:**

- Internet connection is compulsory for this application to send the location of the person.
- Once a guardian calls a person and if the person does not lift the phone then automatically the person's

location details will be sent to the guardian.

- The application should be installed on both the guardian's and user's phone. So, only the guardian will receive the location details of the user if he does not lift the phone but not when the user calls and when the guardian does not lift the phone.
- In another method, both the user and guardian will be having application and both of them will add each other as a guardian and the location details will be sent to both of them in case where either of them does not attend the call.

## **2. Related to work:**

- Current problems people face in society.
- Now a days security is becoming a major issue..
- When the user does not lift the phone he dear ones of that user may get worried.

### **2.1 Existing System:**

1. At times, people might not attend the call due to some reason and the people who care about them might get tensed.

2. We are not having any application for security reasons like this. Hence, we are proposing the following solution.

### **2.2 Proposed System:**

1. Here the location is sent in the form of latitude and longitude.

2. The location details will be converted into address internally in the application and the address is sent to the guardian.

3. We even provide a special feature of seeing that location in the map using this application.

## **3. Technologies Used:**

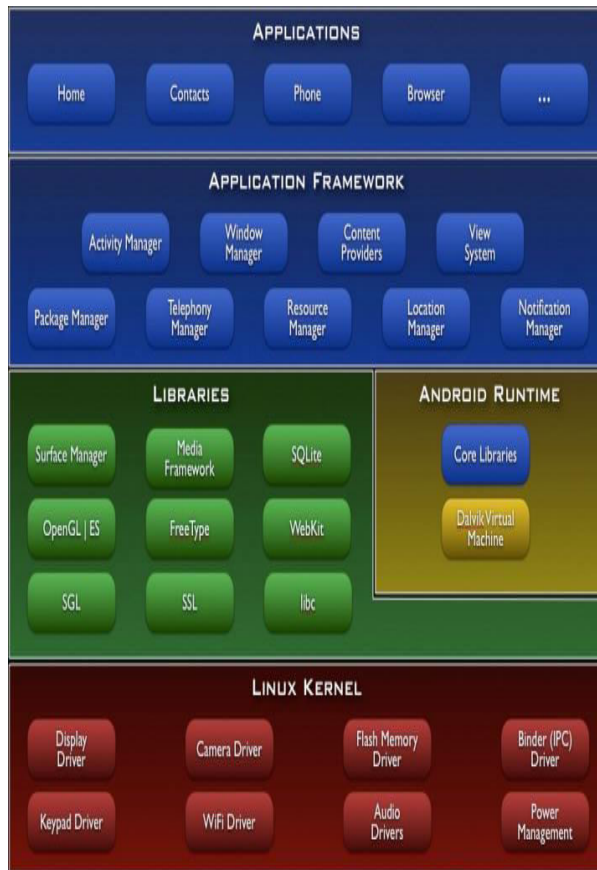
### **3.1 Android:**

#### **What is Android?**

Android is a software stack for mobile devices that includes an operating system, middle ware and key applications. The Android SDK provides the tools and API is necessary to begin developing applications on the Android platform using the Java programming language.

### **3.2 Android Architecture:**

The following diagram shows the major components of the Android operating system.



**Fig:3.1 Android Architecture**

Android is a mobile operating system initially developed by Android Inc. Android was purchased by Google in 2005. Android is based upon a modified version of the Linux kernel. Google and other members of the Open Handset Alliance collaborated to develop and release Android to the world. The Android Open Source Project (AOSP) is tasked with the maintenance and further development of Android. Unit sales for Android OS Smart phone ranked first among all Smartphone OS handsets sold in the U.S. in the second and third quarters of 2010, with a third quarter market share of 43.6%.

Android has a large community of developers writing application programs ("apps") that extend the functionality of the devices. There are currently over 100,000 apps available for Android. Android Market is the online app store run by Google, though apps can be downloaded from third party sites (except on AT&T, which disallows this). Developers write in the Java language, controlling the device via Google-developed Java libraries.

The unveiling of the Android distribution on 5 November 2007 was announced with the founding of the Open Handset Alliance, a consortium of 79 hardware, software, and telecom companies devoted to advancing open standards for mobile devices. Google released most of the Android code under the Apache License, a free software and open source license.

The Android operating system software stack consists of Java applications running on a Java based object oriented application framework on top of Java core libraries running on a Dalvik virtual machine featuring JIT compilation. Libraries written in C include the surface manager, Open Core media framework, SQLite relational database management system, OpenGL ES 2.0 3D graphics API, WebKit layout engine, SGL graphics engine, SSL, and Bionic libc. The Android operating

system consists of 12 million lines of code including 3 million lines of XML, 2.8 million lines of C, 2.1 million lines of Java, and 1.75 million lines of C++.

#### **4. Conclusion:**

This application is mainly used for security purpose. It gives us the details of the location of a person when he does not lift the phone using the GPS system. This application has to be installed in both the phones i.e. the user as well as the guardian. At first the user has to add at least one guardian to his application so that the details will be sent to him. Now, if the guardian calls and the user does not lift the phone the latitude and longitude of the user are calculated based on the GPS

system enabled in his mobile and this sms is sent to the guardians. Then internally these latitude and longitude values are taken and the address is calculated by the application and the address will be displayed in the received sms list. Then by clicking on this address we can see the map using Google maps.

#### **4.1 Future Scope:**

Moving to the future scope, we can provide additional features like sending the location messages free of cost and other features like in which there will be a flexibility of using the application only on one mobile instead of installing on every mobile.

# A Survey of Nature Inspired Load Balancing Algorithms in a Cloud Computing Environment

Priyanka manikonda<sup>1</sup>, Veerender A<sup>2</sup>

Department of computer science and engineering

Mallareddy College Of Engineering

Kompally, Hyderabad, Telangana, India

[manikonda.priyanka@gmail.com](mailto:manikonda.priyanka@gmail.com), [veerender57@gmail.com](mailto:veerender57@gmail.com)

**Abstract-** One of the primary issues in cloud computing is implementation of a novel load balancing approach. The demanding thirst for optimal performance of the system is creating research interest in this area. Many Load Balancing algorithms that aim to enhance the overall system performance have been proposed. In this paper, we survey a special group of Load balancing algorithms that have taken inspiration from nature. We provide an overview of the current trends in the field by discussing and comparing these algorithms.

**Keywords-** Cloud computing, Load balancing, swarm Intelligence.

## I. INTRODUCTION

The exponential growth of cloud computing in the recent years has attracted research and academia to this field. Load Balancing is a primary issue that needs to be taken care of. Several, Load Balancing algorithms have been proposed and investigated; however, there are issues yet to be addressed. Load balancing is “the process of distributing the work load among various nodes of a cloud based system to improve both resource utilization and job response time while also avoiding a situation where some of the nodes are heavily loaded while other nodes are idle or doing very little work.”

Load balancing algorithms are divided as static and dynamic based upon the working environment [1], centralized and distributed based upon the control strategy [2]. Static

algorithms are effective in stable and homogenous environments where as Dynamic algorithms are effective in dynamic and heterogeneous environments. The centralized strategy requires an arbiter or control node to perform the load balancing act whereas in distributed strategy load balancing is performed by all the nodes of the system. Many Load Balancing algorithms have been proposed in the recent past. In this paper, we present a survey of Nature Inspired Load Balancing Algorithms that have been specifically developed for hosted environments. These algorithms use the concept of Swarm intelligence for Load Balancing [3]. We consider some of the potentially viable nature inspired algorithms for load balancing in large scale cloud environments. There are two popular classes of Nature Inspired Algorithms available in the literature namely Ant Colony and Honey Bee Colony. We give an overview of these

algorithms, discuss their pros and cons and analyze their properties.

The rest of this paper is organized as follows. We discuss the related work in Section II. Then, in Section III we discuss the various challenges and issues of Nature Inspired load balancing in cloud computing environment. Afterwards in Section IV, we review and compare Nature Inspired Load Balancing Algorithms that are available currently in the literature. Section V concludes the paper and highlights future enhancements that can be done in Cloud Load Balancing.

## II. RELATED WORK

Klaithem Al Nuaimi, Nader Mohamed, Mariam Al Nuaimi and Jameela Al-jaroodi have presented a survey of Load Balancing in Cloud Computing. The paper gives an overview of Load Balancing Algorithms like INS, ESWLC, CLBDM, Ants Colony, Mapreduce, VM Mapping and DDFTP. It compares the algorithms based upon on certain parameters[1]. Martin Randles, David Lamb,A. and Taleb-Bendiab have presented a comparative study of three distributed load balancing algorithms namely Honeybee based load balancing, Biased Random Sampling and Active Clustering. The paper describes and compares the algorithms by performing experiments using simulations set up in Repast.NET[4].Rich Lee and Bingchiang Jeng in their work made a comparitative analysis of Round-Robin, Weighted Round-Robin,Least Connection, Shortest Expected Delay, Resource Best and Resource Fit algorithms using a simulation program based on GNU R[5]. V. Sesum-Cavic and E. Kuhn have presented the advantages of using swarm intelligence in load balancing in their works [6].

## III. ANALYSIS OF ISSUES RELATED TO NATURE INSPIRED CLOUD LOAD BALANCING ALGORITHMS

In this section we give an introduction to the major challenges a Nature Inspired Cloud Load Balancing Algorithm must address before it is implemented in the system. These challenges if not addressed properly may affect the performance of the algorithm. These challenges are summarized as follows.

### A. *Nature of the Cloud*

Static Nature Inspired Cloud Load Balancing Algorithms are designed to work with static clouds. The performance of these algorithms is satisfactory in stable cloud environments. However, it is a challenge to design a dynamic Nature Inspired Cloud Load Balancing Algorithm that is flexible and adapts to the dynamic changes in the attributes of the system.

### B. *Control Mode*

Nature Inspired Cloud Load Balancing Algorithms are mostly centralized. Having a single point of control in large scale cloud computing environments makes load balancing a daunting task. Moreover, if the arbiter goes down it brings the whole system to a halt which is not desirable. Hence, distributed or even hybrid algorithms are required. Designing a distributed Nature Inspired Cloud Load Balancing Algorithm that gives optimal performance is a challenge.

### C. *Resource Awareness*

Nature Inspired Cloud Load Balancing Algorithms are designed to work with homogenous resources. But in most cases, a cloud computing environment is a collection of heterogeneous resources. If the algorithm

being implemented is not resource aware then its performance will be adversely affected. Hence, developing a Resource aware Nature Inspired Cloud Load Balancing Algorithm is a challenge.

#### *D. Geographical separation of the cloud nodes*

Algorithms that perform well with cloud nodes distributed over geography are yet to be devised. There are some additional parameters that must be considered by the algorithm during balancing. The parameters include network speed, distance between processing node and clients, distance between nodes. Hence algorithms that work well with geographically separated cloud nodes are required [7].

#### *E. Replication Model*

In order to guarantee the SLA's Nature Inspired Cloud Load Balancing Algorithms must support data replication. This can be either complete or partial replication. A fully replicated algorithm comes with additional costs as data on replication nodes must be maintained. It doesn't use the available storage efficiently. On the other hand even though partial replication algorithms utilize the storage resources efficiently, they are very complex to design [8].

#### *F. Ease of implementation*

It is desirable that Nature Inspired Cloud Load Balancing Algorithms are easily implemented and operated. Complexity in algorithm's implementation will raise performance issues. Therefore, Simple algorithms are required [9].

#### *G. Network Overhead*

Nature Inspired Cloud Load Balancing Algorithms have to deal with overhead of the network as the so called software ants honey bees continuously traverse through the network to gather information of the cloud nodes. The presence of large number of ants or bees may sometimes degrade the system performance [10].

#### *H. Scalability*

The scale of a cloud computing platform can be very large. Nature Inspired Cloud Load Balancing Algorithms must be devised to take in to consideration the scalability factor. They must be flexible enough to work in a situation where a resource can be randomly added. The algorithm must quickly take into consideration the newly added resource and perform load balancing with little effect on system performance.

#### *I. Synchronization*

Nature Inspired Cloud Load Balancing Algorithms perform the operation of load balancing with the help of agents like artificial ants or honey bees. The algorithm generates a large number of agents. These agents continuously traverse through the cloud and monitor it. Achieving synchronization among the agents is a challenging task for the algorithms.

## **IV. REVIEW OF NATURE INSPIRED CLOUD LOAD BALANCING ALGORITHMS**

In this section we give a review of Nature inspired cloud load balancing algorithms that are currently available in the literature. Nature inspired cloud load balancing algorithms can be categorized in to two classes namely, Ant Colony Inspired Algorithms and Honey Bee Inspired Algorithms. We first discuss the Ant colony Inspired Algorithms that have been developed for cloud load balancing. Then

later discuss the Honey Bee Inspired cloud load balancing algorithms. TABLE I shows a comparison of Nature Inspired Load Balancing Algorithms that are reviewed in this paper.

#### *A. Ant Colony Inspired Load Balancing Algorithms*

All Ant Colony Inspired Load Balancing Algorithms are based on ACO algorithm [11, 12]. Individual ants are quite simple insects but colony of ants collectively perform a variety of tasks such as building anthills, foraging for food with great reliability and consistency[13,14]. This social behavior of ants has inspired researchers to solve many computational problems that includes even the problem of Load Balancing a Cloud. We discuss three popular ant colony inspired load balancing algorithms.

Zehua Zhang and Xuejie Zhang have proposed a Load Balancing Mechanism based on Ant Colony and Complex Network Theory (LBMACCN) in open cloud computing federation (OCCF)[15]. The mechanism aims to solve the issue of complex and dynamic load balancing in OCCF. Underload load balancing, Overload balancing, Pheromone updating, and Network Evolution are the major modules of LBMACCN. The operation of the algorithm suits for heterogenic cloud environments. The algorithm is designed to work with distributed clouds. The fault tolerance and the scalability factor of the mechanism are excellent. The algorithm requires a full replication model. It has to deal with network overhead problem because ants keep on traversing across the network which sometimes may lead to delays. Synchronization of the ants and Static nature of the algorithm are issues to be considered.

Kumar Nishant et al. - have proposed an algorithm that performs the task of load balancing of nodes in cloud using the Ant colony optimization [10]. The algorithm is an improved version of algorithm presented in [15]. It aims at efficient load distribution among the cloud nodes such that the ants never come across a dead end during network traversal for building an optimum solution set. The algorithm suits for heterogenic and distributed clouds. The fault tolerance and the scalability factor are excellent. The issue of synchronization between ants is solved. The algorithm requires a full replication model and has to deal with network overhead problem. It is a Static Load Balancing Algorithm.

Both the above mentioned algorithms work in the following manner, ants and pheromones are generated from a head node once a request is generated. In order to gather node information for task scheduling ant behavior is used. The ants start their forward route from the 'head'. A forward movement means that the ant is searching for overloaded nodes whereas reverse movement indicates it is searching for an underloaded node. The authors of [10] have introduced a new feature termed 'suicide' to the ants. Once the target node is found the ant will be terminated preventing unnecessary backward movement.

Kun Li et al. - have proposed LBACO (Load Balancing Ant Colony Optimization) algorithm for load balancing in cloud based systems [16]. The algorithm is simulated using CloudSim version 2.1 toolkit package. The algorithm is based on ACO algorithm [11, 12]. The algorithm considers the past task scheduling time to carry out new scheduling strategy. It reserves the current optimal solution and uses it to make a decision in future Load balancing scenarios. The experiment discussed in [16] proves that

LBACO is more effective when compared with ACO algorithm. The algorithm suits for dynamic, distributed and heterogeneous clouds. The algorithm assumes that the tasks are mutually independent, preemptive and computationally intensive.

### *B. Honey Bee Inspired Load Balancing Algorithms*

Honey Bee inspired technique is used as a search technique in many computing applications where the system is highly scalable and dynamic [17]. A bee hive works as follows. There are two roles a bee can play namely forager bee and follower bee. Forager bees search for a suitable source of food, when found, they return back to the colony and advertise the same using a “waggle dance”. The quality, distance and quantity of the food are judged based upon the intensity of the dance. Follower bees are sent out to harvest the discovered food. The waggle dance is used to determine whether more bees are required to harvest or the exploited source to be abandoned. The same phenomenon is being applied in load balancing of nodes in cloud computing environment. We discuss 2 popular bee inspired Load Balancing Algorithms that are available in the literature.

Martin Randles, David Lamb and A. Taleb-Bendiab have proposed a Honey bee inspired Load Balancing algorithm [4]. The working of the algorithm can be explained as follows. The servers in the cloud take the role of either foragers or harvesters. An advert board that mimics the waggle dance of the bees is maintained. A server successfully executing a request will post on the board. A server that reads the board follows the chosen advert, and then serves the request; thus mimicking the harvest bee. The server not reading the advert board serves a random virtual server's queue request; thus mimicking the forager bee. The

total colony profit is calculated based upon just-serviced virtual server profit. If the just serviced virtual server's profit is high then the forager server will keep on posting an advert for it until the calculated profit becomes low. The algorithm is suitable for large scale, dynamic and heterogeneous cloud environments. The algorithm suffers with Network Overhead and replication issues. Given an optimum profit calculation method, this algorithm provides a good distributed solution to the problem of Load balancing in cloud computing environment.

Jing Yao and Ju-hou He, have presented a Load Balancing Strategy for cloud computing based on Artificial Bee algorithm [18]. The algorithm is an improved version of the algorithm proposed in [4]. The previous mechanisms were considering only lightly loaded nodes for balancing where the technique in [18] considers all resources in the cloud iteratively. The algorithm includes some extra operations that improve the quality of service. Experimental results show that in a system with a fixed number of systems and increasing number of requests the algorithm generates an improved throughput whereas a change in number of servers and fixed number of requests the original ABC algorithm performs well. The algorithm is dynamic, distributed and suitable for heterogeneous cloud systems; whereas, it suffers in a highly scalable system. Network Overhead and replication model should also be considered before implementing this algorithm in a cloud based system.

## V. CONCLUSIONS AND FUTURE WORK

We present a survey of Nature Inspired Load Balancing Algorithms that are available in the literature. We have discussed the challenges that these algorithms must address. The pros and cons of the algorithms have also been depicted. We have drawn out comparisons between the algorithms by taking into account the necessary parameters. All the algorithms discussed in the paper suffer with one or the other issue. Hence there is scope for improvement in the algorithms. Therefore, as part of future work we plan to propose our own Nature Inspired Load Balancing Algorithm for cloud computing that addresses the issues discussed earlier.

**TABLE I Comparison of Nature Inspired Load Balancing Algorithms**

|                                        | LBMA<br>CCN | LBMA<br>CCN in<br>[17] | LBACO  | AB<br>C in<br>[20] | ABC in<br>[21] |
|----------------------------------------|-------------|------------------------|--------|--------------------|----------------|
| Dynamic                                | NO          | NO                     | YES    | YES                | YES            |
| Distributed                            | YES         | YES                    | YES    | YES                | YES            |
| Heterogeneous                          | YES         | YES                    | YES    | YES                | NO             |
| Geographical<br>Separation<br>of nodes | NO          | NO                     | NO     | NO                 | NO             |
| Complexity                             | EASY        | EASY                   | MEDIUM | EASY               | MEDIUM         |
| Network<br>overhead                    | YES         | LESS                   | YES    | YES                | YES            |
| Scalability                            | NO          | NO                     | YES    | NO                 | NO             |
| Synchronization                        | NO          | YES                    | YES    | NO                 | YES            |

## REFERENCES

- [1] Klaithem Al Nuaimi, Nader Mohamed, Mariam Al Nuaimi and Jameela Al-Jaroodi "A Survey of Load Balancing in Cloud Computing:Challenges and Algorithms" IEEE Second Symposium on Network Cloud Computing and Applications. 2012
- [2] Yilin Lu,Shaochun Wu,Jian Zhang and Shujuan Zhang "A Hybrid Dynamic Load Balancing Approach for Cloud Storage" International Conference on Industrial Control and Electronics Engineering 2012
- [3] Vesna Sesum-Cavic and Eva Kühn "Applying swarm intelligence algorithms for dynamic load balancing to a Cloud Based Call Center" Fourth IEEE International Conference on Self-Adaptive and Self-Organizing Systems 2010.
- [4] Martin Randles, David Lamb and A. Taleb-Bendiab "A Comparative Study into Distributed Load Balancing Algorithms for Cloud Computing" IEEE 24th International Conference on Advanced Information Networking and Applications Workshops 2010.
- [5] Rich Lee and Bingchiang Jeng "Load-Balancing Tactics in Cloud" International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery 2011.
- [6] V. Sesum-Cavic and E. Kühn, "Comparing configurable parameters of swarm intelligent algorithms for dynamic load balancing", submitted for publication.
- [7] Buyya R., R. Ranjan and RN. Calheiros, "InterCloud: Utility-oriented federation of cloud computing environments for scaling of application services," in proc. 10th International Conference on Algorithms and Architectures for Parallel Processing (ICA3PP), Busan, South Korea, 2010.
- [8] Foster, I., Y. Zhao, I. Raicu and S. Lu, "Cloud Computing and Grid Computing 360-degree compared," in proc. Grid Computing Environments Workshop, pp: 99-106, 2008.
- [9] Grosu, D., A.T. Chronopoulos and M. Leung, "Cooperative load balancing in distributed systems," in Concurrency and Computation:Practice and Experience, Vol. 20, No. 16, pp: 1953-1976, 2008.

[10] Kumar Nishant, Pratik Sharma, Vishal Krishna, Chhavi Gupta, Kuwar Pratap Singh, Nitin and Ravi Rastogi “ Load Balancing of Nodes in Cloud Using Ant Colony Optimization” 14th International Conference on Modelling and Simulation 2012

[11] M. Dorigo, C. Blum, “Ant colony optimization theory: A survey” in Theoretical Computer Science 344 (2–3) (2005), DOI: 10.1016/j.tcs.2005.05.020, pp.243–278, 2005.

[12] M. Dorigo, M. Birattari, T. Stutzel, “Ant colony optimization”, in IEEE Computational Intelligence Magazine, DOI: 10.1109/MCI.2006.329691, pp.28-39, 2006.

[13] E.Bonabeau, MDorigo, and G.Theraulaz, "Inspiration for optimization from social insect behavior," Nature, vol.406, pp.39-42, July2000.

[14] MDorigo, G.D.Caro, and L.M.Gambardella, "Ant algorithms for discrete optimization," ArtifLife, vol.5, no.2, pp.137-172, 1999.

[15] Zehua Zhang and Xuejie Zhang “A Load Balancing Mechanism Based on Ant Colony and Complex Network Theory in Open Cloud Computing Federation” 2nd International Conference on Industrial Mechatronics and Automation 2010.

[16] Kun Li, Gaochao Xu, Guangyu Zhao, Yushuang Dong and Dan Wang “Cloud Task scheduling based on Load Balancing Ant Colony Optimization” Sixth Annual ChinaGrid Conference 2011.

[17] S. Nakrani and C. Tovey, On Honey Bees and Dynamic Server Allocation in Internet Hosting Centers. Adaptive Behavior 12, pp:223-240 (2004).

[18] Jing Yao and Ju-hou He Load Balancing Strategy of Cloud Computing based on Artificial Bee Algorithm 8<sup>th</sup> International conference on computing technology and Information Management, vol1 (2012)

# Need for Various Energy Efficient Mechanisms in the Wireless Sensor Networks

Puladas Sandhya Priyanka<sup>1</sup>

<sup>1</sup>Assistant Professor, CSE Department,

Malla Reddy College of Engineering, Dhulapally,  
Kompally, Secunderabad

Rashmitha<sup>2</sup>

<sup>2</sup>Assistant Professor, CSE Department,

Malla Reddy College of Engineering, Dhulapally,  
Kompally, Secunderabad

**Abstract** – Sensor nodes in the wireless sensor networks (WSNs) are of battery made-up. Hence they are provided with limited amount of energy. The battery of a sensor node is unchargeable and also it is difficult to recharge them manually. Therefore they are to be used in more efficient manner for longer life time. It is difficult to increase the life time of the battery by manually but by implementing some protocols or methods of energy efficiency that can be achieved. Many research works and various protocols have been developed for the past few years in the WSNs. Hence there are various schemes or mechanisms designed for the efficiency of the wireless sensor networks. These things majorly concentrate on the improvement of the energy efficiency of the nodes in the networks where energy is drained out by performing several actions like communicating to the neighbor nodes for the transmission of the packets, by sensing the medium, during the reception of the acknowledgements, during the neighbor node discovery etc. Some protocols have been developed in the regard of providing energy efficiency.

**Key Words:** Wireless Sensor Networks, Sensor Nodes, Energy – Efficiency.

## I. INTRODUCTION

Wireless Sensor Networks (WSNs) are the networks which are prominently used in the information technology. They consist of the tiny devices with sensors named as Sensor Nodes. The Sensors perform the tasks like – sensing, minor computations and forward the sensed data to the sink or the base station. Mobile Sensor Networks are those which are a class of the WSN and are growing rapidly in the field of various applications. Hence these are having many challenges like – to provide connectivity among the sensor nodes as well as at the same time to maintain the energy consumption at minimum in the sensor nodes as they are made up with batteries

where energy cannot be recharged further. Hence these issues made researchers to think and develop various mechanisms in order to provide energy efficiency in the network. Few applications of the WSN are: Home and Office, Control and Automation, Environment Monitoring, Health care, Military, Security and Surveillance etc.

The following lines give an idea about the sensor node – the block diagram, characteristics etc.

### 1.1 Characteristics of Sensor Node:

The following is the list of characteristics of a sensor node

- They are tiny in size,
- Have limited memory size to store the values or queue them before transmitting,
- Short range of transmission,
- Low energy capacity,
- Limited lifetime.

The Fig 1, illustrates the block diagram of the sensor node. It shows how a sensor node is and gives a clear picture of the above listed characteristics of the sensor node.

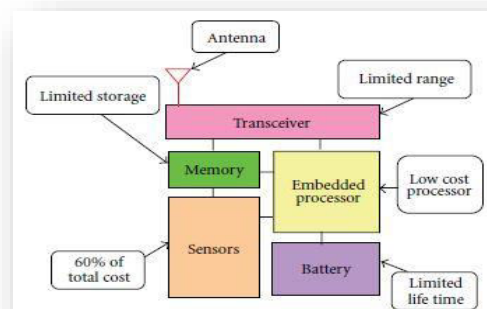


Fig 1: Block Diagram of Sensor Node

## 1.2 Conceptual Analysis of the Sensor Node:

The sensor nodes are made of batteries which have limited amount of energy and this energy can be neither recharged nor replaced. Therefore it has to be more efficiently used. This is because the most of the energy of sensor node is consumed during the process of Communication with other sensor nodes in the network. Hence the key point is to minimize the energy consumption of the sensor nodes and improve the performance of the network.

## II. ROLE OF THE VARIOUS ENERGY MECHANISMS

The energy consumption is the key factor in wireless sensor networks, which is leading the researchers to perform various alternatives in this area to improve the energy efficiency by making the consumption to minimum and utilize only when needed. The following figure illustrates the various energy mechanisms in the wireless sensor networks.

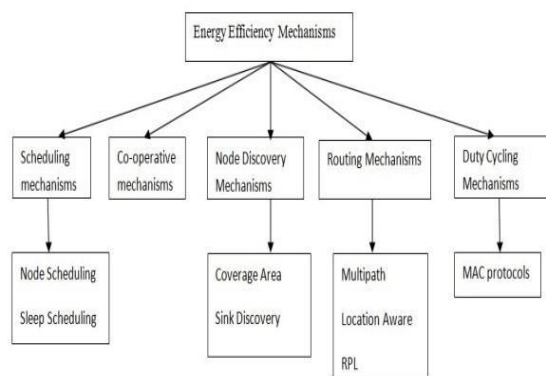


Fig 2: Hierarchical Representation of the Energy Efficient Mechanisms

### 2.1 Role of Co-Operative Mechanisms

**Need of Co-operative Mechanisms:** In order to have proper transportation of data among the sensor nodes in the network, there must be cooperation among the nodes. Cooperation mechanism leads to establish sending of data where no data is lost. The following lines make us clear in knowing how the cooperative mechanism is carried out in the wireless sensor network.

**Various Co-operative Mechanisms:** The energy efficiency is carried out by the coalition-based data transport mechanism where neighboring nodes are organized into groups to form coalitions and sensor nodes within one coalition carry out cooperative communications [20]. The other scheme is that the

nodes choose the efficient path selection by relay nodes and then transport the data to the next node [29]. If the medium is of broadcast nature then the network adapts two phase model where in one phase the data transport is sent to the receiver node and if the data is not received by that node then in second phase it is sent through the relay node [24]. Also the cooperation can be based on the two protocols called DRP (Data Reservation Protocol) and PCA (Prioritized Channel Access) where the slot reservation take place and hence the collisions are avoided [15].

### 2.2 Role of Node Discovery Mechanisms

**Need of Node Discovery Mechanisms:** One of the important design issues of the wireless sensor networks is to have knowledge about the neighbor nodes and the proper connectivity among them in the network. The following lines describe how the node discovery mechanism is being implemented so far.

**Various Node Discovery Mechanisms:** By sending the beacon messages the nodes present in the node can be known [13]. Also target discovery plays a major role as the node need to know about the nearby node to which it has to send the data. Thereby, the nodes can be classified into disjoint sets so as to perform the tasks in an efficient way [12]. By adapting the selective communication mechanism the node can know about the efficient path in which nodes are capable of handling the data when routed and by having prioritization of the data high priority data can be sent fast whereas low priority can be discarded [23].

### 2.3 Role of Routing Mechanisms

**Need of Routing Mechanisms:** To send the data from one node to another series of instructions to be followed, so that there exists some path between the nodes to transfer the data to reach the destination. There by, nodes do spend their energy in locating the path of the nodes to route the packets and hence many routing mechanisms have been proposed by many researchers. The following lines describe about the various mechanisms implemented so far in the wsn.

**Various Routing Mechanisms:** In routing the packet, the location of a node to which the sender node has to be sent must be known and thereby, location aware comes into existence since all the nodes in the network will be in mobile in nature [2]. When the network is considered with the cluster formation for performing the tasks in more efficient manner, Table maintenance and the cluster head maintenance comes into the picture in order to maintain the consistency of the nodes and their

energy information [16]. A New Multipath Routing Approach for the energy efficiency is also proposed in which Route Request control messages can be sent by and controlled [21].

## 2.4 Role of Scheduling Mechanisms

**Need of Scheduling Mechanisms:** By scheduling the nodes in the network the amount of energy consumption can be reduced. Scheduling is to be used because there will be division of work among the nodes, also depending upon the scheduling times nodes can be made awake either to receive the data or to send the data.

**Various Scheduling Mechanisms:** Some of the schedulings are: Wake Scheduling, in which the edge nodes when arranged in grid form are taken into account and are scheduled accordingly [19]. Sleep Scheduling tells that when the nodes are to be in sleep state by sending or making the node to hear the alarm message [8]. Also by coordination of the nodes through implementation of the eligibility of the nodes and by the maintaining the two back-off schemes so that no two nodes remain off at the same time, the energy efficiency can be maintained [6]. By knowing the coverage area of particular number of nodes till where they can sense the data, which can be applicable for large dense wireless sensor networks, can also be considered for the energy consumption [3]. Also to avoid blind points when any two nodes switch to off state scheduling algorithms have been proposed [25].

## 2.5 Need of Duty Cycling MAC Protocols:

**Various Duty Cycling MAC Protocols:** In the aspect of the providence of the energy efficiency different MAC protocols have been designed namely: S-MAC [4], Q-MAC[17], CF-MAC and H-MAC [5] Asynchronous MAC for QoS[9], T-MAC[27], EQ-MAC for energy efficient and quality of service providence [1]. P-MAC where patterns are being exchanged between the nodes for communication initiation[6], and also the which tasks can a sensor node handle in a network when the sensor network is designed for multi application scenarios which is being called as task aware protocol, shortly TA-MAC [22]. Instead of beacons, sensor nodes can also make use of the preambles sending, so that other nodes can know when they have to be in sleep state and when they have to be awake [11].

## 2.6 Importance of Priority Mechanisms:

There are many ways of assigning the priorities to the nodes [30]. Some of them are:

- ADISK – Local disk currently available in the batch jobs
- AMEM – Real memory currently available
- JOBCOUNT – Number of jobs currently running on a node
- POWER – Depending on the power of the node
- SPEED – Depending on the processing speed it could perform a job

Apart from these, there are different methods of how the priority can be allocated. Priority Hybrid model, which consists of two buffers like event-driven and clock-driven buffers. The emergency packets when received from the other nodes are stored in the event-driven model saying that some event has occurred and the transmission of packets from one to other in a normal way, that are stored in the regular clock-driven model [8]. The priority is assigned in such a way that if there is any real time data then that is send to the highest priority queue and if not then the data in the normal priority is being processed. Also, it is said that the real time data can preempt the data at other queues [10]. An algorithm where there exists three different levels of queues with priorities given. All the real time data packets go to the highest level priority and the non-real time data which has been received from the other nodes go into the second level priority queue and the non-real time data which is at the local node is sent to the lowest priority queue[18]. In addition to, the assigning of the priority by implementing the fuzzy logic where, the data is initially turned into the linguistic values and then stored fuzzy rule base and the decision making like human is done in the inference engine whether to be given highest priority or not [28].

## III. CONCLUSION

As the technology has grown into the automatic mechanism, the sensors are part of it. Thereby, several applications of the sensors have been taking place in various fields like – hospitals, crime investigation by vehicle tracking, military, civil engineering, etc. In order to increase the application of sensors widely, the major thing is to design the sensor networks in a more efficient and effective way. In this regard, the major problem faced by the sensor networks is the energy degradation of the nodes which are battery made up and have limited amount of energy which is also not be able to recharged further. All these mechanisms indicate how to maintain the efficiency and the performance of the network by reducing the energy consumption of the nodes eventually of the entire network. These help the researchers to perform research and increase the performance levels of the wireless sensor networks which would be the most benefitted part in the recent technology development where sensor

world would be seen in the very few years of the present generation.

## REFERENCES

- [1] Bashir Yahya and Jalel Ben-Othman "An Energy Efficient Hybrid Medium Access Protocol Scheme for Wireless Sensor Networks With Quality of Service", in IEEE"GLOBECOM" 2008 proceedings.
- [2] C.Shanti and D.Sharmila "A self-organized location aware energy efficient protocol for wireless sensor networks" in Elsevier Computers and Electrical Engineering,2014.
- [3] Di Tian and Nicolas D.Georganas "A coverage-Preserving Node Scheduling Scheme for Large Wireless Sensor Networks" in WSNA'02, Atlanta, September 2002.
- [4] D Saha, M R Yousuf and M A Matin "Energy Efficient Scheduling Algorithm For S-MAC protocol in Wireless Sensor Networks" in International Journal Wireless Mobile Networks, 2011.
- [5] G.Boggia, P.Camarda, O.Fiume and L.A. Grieco "CF-MAC and H-MAC protocols for Energy Saving in Wireless Adhoc Networks", in IEEE (0-7803-8887-9), 2005.
- [6] Hend Alqamzi and Jing Li "An Efficient Energy-balanced Coordinated Node Scheduling (ECONS) Protocol for Dense Wireless Sensor Networks" in IEEE GLOBECOM proceedings 2006.
- [7] Hsu-Jung Liu, Mei-Wen Huang, Wen-Shyong Hsieh, Chenhuan Jack Jan, "Priority-based Hybrid Protocol in Wireless Sensor Networks", 11<sup>th</sup> IEEE International Conference on High Performance Computing and Communications, 2009.
- [8] Jae Hyun and H. Jin kim "Predictive Target Detection and Sleep Scheduling for Wireless Sensor Networks" in IEEE International Conference on Systems, Man, and Cybernetics, 2013.
- [9] Kien Nguyen and Yusheng Ji "Asynchronous MAC protocol with QoS awareness in Wireless Sensor Networks" in IEEE (978-1-4673-0921-9), 2012.
- [10] Lutful Karim, Nidal Nasser, Tarik Taleb, Abdullah Alqallaf, "An Efficient Priority Packet Scheduling Algorithm for Wireless Sensor Network", IEEE ICC, Adhoc Network Symposium, 2012.
- [11] Michael Buettner, Gary V. Yee and Eric Anderson , Richard Han "X-MAC: A short Preamble MAC protocol for Duty-cycled Wireless Sensor Networks" in SenSys'06, November 2006.
- [12] Mihaela Cardei, My T Thai Yingshu Li and Weili Wu "Energy-Efficient Target Coverage in Wireless Sensor Networks" in IEEE (0-7803-8968-9), 2005.
- [13] Mikko Kohvakka, Jukka Suhonen, Mauri Kuorilehto, Aille Kaseva, Marko HAnnikainen and Timo D.Hamalainen " Energy Efficient neighbor Discovery protocol for mobile sensor networks", in Elsevier Ad hoc Networks, 2009.
- [14] Moshadique Al Ameen, S.M. Riazual and Kyungsup Kwak " Energy Saving Mechanisms for MAC protocols in Wireless Sensor Networks", in Hindwai Publications, International Journal of Distributed Sensor Networks, 2010.
- [15] M.Rengasamy, E.Dutkiewicz and M.Hedley "MAC Design and Analysis for Wireless Sensor Networks" in International Symposium on Communication Information Technologies, 2007.
- [16] Mudasser Iqbal, Iqbal Gondal and Laurence Dooley "A cross-layer Dissemination protocol for Energy Efficient Sink Discovery in wireless sensor networks" in IEEE for publication in ICC proceedings 2007.
- [17] N.A Vasanthi and S.Annadurai " Energy Efficient Sleep Schedule for Achieving Minimum Latency in Query based Sensor Networks" in Proceeding of IEEE International Conference on Sensor Networks Ubiquitous and Trustworthy Computing, 2006.
- [18] Nidil Nasser, Lutful Karim, Tarik Taleb, "Dynamic Multilevel Priority Packet Scheduling Scheme for Wireless Sensor Network", IEEE transactions on Wireless Communications, Volume-12, No. 4, 2013.
- [19] Niki Trigoni,Yong Yao,Alan Demers,Johannes Gehrke and Rajmohan Rajaraman "Wave Scheduling: Energy-Efficient Data Dissemination for Sensor Networks" in proceedings of Data Management for Sensor Networks, 2004.
- [20] Qinghal Gao, Junshan Zang, Xuemin and Bryan Larish "A Cross-Layer optimization Approach for Energy-Efficient Wireless Sensor Networks: Coalition-Aided Data Aggregation, Co-operative Communication, and Energy Balancing" in Hindwai Publication, 2007.
- [21] Saira Banu and R.Dhanasekaran "A New Multipath Routing Approach for Energy Efficiency in Wireless Sensor Networks" in International Journal of Computing Applications (0975-8887), 2012.
- [22] Sangheon Pack and Jaeyoung Choi,Taekyoung Kwon and Yanghee Choi " TA-MAC: Task Aware MAC protocol for Wireless Sensor Networks", in IEEE 2006.
- [23] Sara Pino-Povedano, Rocio Arroyo-Valles and Jesus Cid-Sueiro "Selective Forwarding for energy-efficient target tracking in Sensor Networks" in Elsevier Signal Processing, 2014.
- [24] Shaoqing Wang and Jingnam Nie " Energy Efficiency Optimization of Co-operative Communication in Wireless Sensor Networks", in Hindwai Publication 2010.
- [25] Singaram M, Finney Daniel Shadrach S, Sathish Kumar N and Chandraprasad V " Energy Efficient Self-Scheduling Algorithm For Wireless Sensor Networks" in International Journal of Scientific Technology Research, 2013.
- [26] Tao Zheng, Sridhar Radha Krishnan and Venkatesh Sarangan "PMAC: An Adaptive energy-efficient MAC protocol for Wireless Sensor Networks" in proceedings of 19<sup>th</sup> IEEE International Parallel and Distributed Processing Symposium , 2005.
- [27] Tjis van Dam and Koen Langendoen "An Adaptive Energy-Efficient MAC Protocol for Wireless Sensor Networks" in SenSys'03,November 2003,California.
- [28] Varsha Jain, Shwetha Agarwal, Kuldeep Goswami., "Priority Based Fuzzy Decision Packet Scheduling Algorithm for QoS in Wireless Sensor Network", International Journal of Computer Applications (0975-8887), Volume-97, 2013.
- [29] Weiwei Fang, Feng Liu, Fangnan Yang, Lei Shu and Shojiro Nishio, "Energy – Efficient Cooperative Communication for Data Transmission in Wireless Sensor Networks", in IEEE 2010, (0098 3063), Page No: 2185-2192.
- [30] [http://docs.adaptivecomputing.com/mwm/Content/topics/prio\\_res/nodeallocation.html](http://docs.adaptivecomputing.com/mwm/Content/topics/prio_res/nodeallocation.html)

# **CAPTCHA BREAKING USING IMAGE TEMPLATE MATCHING AND MACHINE LEARNING ALGORITHMS**

**G.MADHURI**  
Assistant Professor

**V.Sandhya**  
Assistant Professor

## **Abstract**

CAPTCHAs (Completely Automated Public Turing test to Tell Computers and Humans Apart) have become a very popular security mechanism used to prevent automated abuse of online services intended for humans. Different flavors of CAPTCHA can be seen on Internet. However, a wide variety of CAPTCHAs have been successfully attacked by automated programs. This has made CAPTCHA design an interesting area for research. Among various flavors of CAPTCHA text based are most preferable because of its low implementation cost but this category of CAPTCHA can be attacked with very high success rate. One of the problems with respect to CAPTCHA is “hard to separate text from background”. The underlying security assumption is that if an automated program cannot locate text in CAPTCHA then it cannot do further processing to identify the characters of it.

This project proposing two novel approaches for breaking CAPTCHAs. The first proposed approach aim is, segmenting denoised characters from the background using morphological operations and recognition of the characters using histogram of white pixels. Many researchers are concluded that a significant effort is needed in segmentation before individual character classification can be attempted in CAPTCHAs breaking. This proposed method successfully segments the CAPTCHA from the background and segment it in terms of individual characters. And second one is, single character classification using of machine learning techniques. Very few researchers worked on machine learning techniques for character recognition. Hence my project gives a new direction in single character classification. The performance of this method is measurable in character classification using machine learning classifiers like Libsvm, Multilayer perceptron, Naviebayes, J48 and IbK classifiers. This method performs extremely well for breaking easily segmentable CAPTCHAs, with a robust recognition with 90% recognition rate.

## INTRODUCTION

### 1.1 MOTIVATION

Digital image processing is a rapidly flourishing field of computer science. The recent developments in digital image acquisition, microprocessors triggered the growth of digital image processing. Since digital systems offer better flexibility and affordability, fields which are using analog imaging are shifting to digital systems. Medicine, video production, remote sensing, photography, security monitoring are the best examples of fields which adopted digital imaging. These and added origins bring forth large amounts of digital image information day-to-day.

Mathematical morphology (MM) is rich in theoretical framework for solving image analysis and processing problems [32]. MM is a part of natural science that deal with image topology, structure, shape, connectivity etc., and MM is basically derived from algebra of non-linear operators. The linear approaches of segmentation fail to find solution for the problems which involve geometrical aspects of the image. Therefore non-linear approaches are required. Mathematical morphology turns to be the most powerful non-linear methodology which can solve the problems involved with geometrical shapes in image segmentation. Most of the texture images tend to be nonlinear in nature because Image is not linear combination of gray levels in the neighborhood. Hence morphological methods are well suited for segmentation of image images. MM transformations represent images in easy way, quantify the images and retain the fundamental shape properties of objects.

MM is quite often applied in problems where object's shape and speed is important. MM is highly fascinating for segmentation as it effectively represents geometrical attributes such as shape, size, connectivity or dissimilarity which is taken as segmentation aligned properties. MM tools find importance because they offer better simplicity and efficiency.

#### 1.1.1 Texture

Images can be understood by their intensity differences that typically originate from object surfaces. Image is the most vital visual property in distinguishing uniform regions. Image of an image, expressed in terms of the spatial change in gray values of pixels, is helpful for several real time implementations and is an area of deep survey for researchers. By this one can understand or estimate the quantifying parameters of surfaces like fine, coarse, smooth, or grained. Analysis of image images requires careful design of statistical parameters because the intensity variation for a well-defined image leads to the derivation of many aspects. There are some commonly used methods for image analysis. These methods depend on the intensity variations and new algorithms are coming up for more efficient modelling, segmenting and classifying images [31]. Therefore it is more reliable to consider object surfaces and it's characteristics for various applications like segmentation, pre-processing, noise reduction, image analysis and recognition etc. Image can be defined as the attribute of the surface that gives rise to this reflectance change locally. In many instances, this attribute arises because of surface

roughness, which leads to scattering of light randomly, thereby reducing or enhancing local reflectance in the viewing direction.

With over a half a century of digital image processing, image is still a very active field of research fascinating various mathematical techniques and algorithms developed for the need of improving object surface characterization.

### **1.1.2 Effect of noise in texture images**

Noise in the image is random change of color information or intensity and usually an aspect of noise is due to electronic parts. It can be originated because of the scanner circuit and sensor or digital camera. Also noise in image can start in ideal photon detector and film grain. Noise in image is an unwanted consequence of image acquisition that adds extraneous and spurious information. The digital images affected with random noise specifically have Gaussian or normal distribution. The magnitude of the error, in the acquired signal because of this kind of noise is independent of the entire signal, similar to quantum mottle. Because of Gaussian noise in the image, light and dark regions are affected to the equal extent. Another name of Gaussian noise is additive noise as the noisy image is equal to adding of noise having Gaussian probability distribution to the original image. Impulse or salt and pepper noise visually appear as random distribution of black and white pixels in an image. The origins of Impulse noise are flaws in sensor components, flaws in storage devices, noise affecting image reconstruction and data transmission. Noise factor affects a lot in accurate, proper and

critical segmentation process. Noise may lead to connected regions, edges and boundaries as non-connected, and also vice-versa and also sometimes leads to over segmentation. Any unwanted further splitting of uniform regions is called over segmentation.

## **1.2 IMAGE SEGMENTATION**

Image segmentation [33, 34, 38] divides an image in the form of groups of different regions depending on image attributes, such that every region is uniform based on specific image properties. Segmentation results can be used for further image analysis and processing like object classification and identification. Image segmentation also consists of obtaining attributes and deducing quantitative parameters to separate images. However, segmentation is usually more complex than classification, since boundaries that segregate different image regions have to be obtained furthermore to recognizing image in each region.

Image segmentation could also be unsupervised or supervised based on if prior knowledge about the image class or image is available. Supervised image segmentation detects and segregates one or more regions that satisfy image properties given in the training images. Unsupervised segmentation has to initially recover various image classes from an image before dividing them into regions. The unsupervised segmentation is more adaptable for real world applications when compared to the supervised case. The unsupervised segmentation is widely used in spite of the

fact that it is generally more computationally expensive.

Dividing an image into uniform regions is very useful in several applications of machine learning and pattern recognition. For instance, in Geographic Information System (GIS) analysis and remote sensing, image segmentation could be useful to identify landscape variation from an aerial photo. The image segmentation is one of the vital tasks in computer vision systems. It plays a vital role in many pattern recognition tasks such as cartography, remote sensing, robot vision, military surveillance, inspection of textile products, and CAPTCHA imaging [35].

Segmentation methods are dependent upon some region or pixel homogeneity parameter which is related to their local neighborhood. Boundary-oriented algorithms find the most different pixels that exhibit discontinuity in the gray values in the image, in contrary region oriented algorithms search for the most likely pixels that can be grouped into regions. These parameters based on similarity in image segmentation techniques employ some textural spectral-spatial-temporal properties like Markov Random Field (MRF) statistics, Gabor features, co-occurrence matrix oriented properties, Local Binary Pattern (LBP), autocorrelation properties and many other properties. Segmentation algorithms can be classified based on different criteria, such as boundary/region oriented methods, clustering algorithms, graph theoretic algorithms, etc.

### **1.2.1 Region growing**

The fundamental principle of region growing technique [36] is to begin from

seed pixels that are supposed to be located within the object of interest which is to be segmented. The adjacent pixels to every seed pixel are estimated whether they belong to the object or not. If the adjacent pixels are identified to be part of the object then they are merged to the object region and the procedure proceeds until any unchosen pixels persist. Region growing methods change based on the similarity property, selection of seed region, connectivity type employed to choose neighboring pixels and the criteria used to traverse adjacent pixels.

Disadvantages of the region growing are

1. Time consuming
2. Sensitive to seed selection

### **1.2.2 Split and merge**

Split and merge methods begin with iterative dividing of image into tiny regions till they fail to obey some uniformity principle. The second step, combines adjacent regions with same properties. Advantage of split and merge is it combines advantages of region splitting and region growing. Disadvantage of split and merge is it produces boundaries similar to square shape.

### **1.2.3 Watershed**

Based on the gradient descent on image properties and survey of less significant pixels on boundaries of region, watershed method combines pixels into regions. Based on satisfactory mapping the image feature space is considered as a topological surface in which larger pixel values signify that boundaries exist in the original image. This method is analogous to landscape basins in which water slowly fills in. As the amount of water increases the sizes of basins become large until the water

overflow from one basin to another. Smaller basins are slowly grouped to form large basins. Based on the local geometric information to combine the features of the image domain with local extremes calculation regions are obtained. The obtained segmentation by watershed is based on the choice of either manual or prior information. These methods are highly applicable for various measurements fusion and they are robust to user specified thresholds. The advantage of watershed method is it produces closed boundaries and the disadvantage is it results in over segmentation.

#### **1.2.4 Level set segmentation**

The prime example of the level set [29] is that it is a numerical technique to find the starting point of surfaces and contours. The contour is implanted as the zero level set of level set function inspite of directly modifying the contour. The level set function is produced by a differential equation that depends on image properties. From the output, zero level-set is extracted and based on this evolving contour can be obtained at any time. Level sets provide path to model randomly topological modifications and complex shapes such as splitting and merging are handled accurately.

#### **1.2.5 Mean shift segmentation**

In [30] Edison image segmentation method using embedded edge information and mean shift is presented. Its starting step uses the mean shift segmenter in the mixed coordinate feature space and color  $L*u*v^*$ . Using edge confidence measure the mean

shift weights are obtained. The next step iteratively combines the regions of modes.

#### **1.2.6 Graph-theoretic segmentation**

The algorithms in [28, 37] employ graphical representation for image regions or pixels. In these methods weighted graph edges are connected mutually small neighborhood pixels. The pair wise elements similarities are represented by these weights.

Processing of CAPTCHA images has experienced dramatic extension, it has been an interdisciplinary field of research attracting expertise from computer sciences, applied mathematics, engineering, statistics, biology, physics and medicine. Computer-aided diagnostic image processing has already become a vital part of clinical routine. Accompanied by a rush of recent development of high technology and use of several imaging modalities, more challenges arise; for example, how to process and analyze a significant number of images so that high quality information can be produced for an effective and accurate disease diagnosis and treatment.

CAPTCHA imaging refers to a number of techniques like computer aided tomography, magnetic resonance imaging, ultrasound imaging, radio isotope imaging, electrical impedance tomography, microscopic imaging that can be used as non-invasive methods to assist diagnosis or treatment of different CAPTCHA conditions. CAPTCHA images are often deteriorated by noise due to various sources of interference and other phenomena that affect the measurement processes in imaging and data acquisition systems.

Recently texture emerged as a powerful visual primitive to succinctly

describe any image content and texture based CAPTCHA image synthesis has become one of the interesting research fields. Texture exists all around us and serves as an important visual cue for the human visual system. Captured within an image, one can identify texture by its recognizable visual pattern. That is the reason texture analysis and characterization is widely studied over the last three decades in a variety of applications, including CAPTCHA imaging, remote sensing, pattern recognition, industrial inspection, texture based image retrieval, human visual perception and provides information for recognition and interpretation.

Images and graphics are among the most important media formats for human communication and they provide a rich amount of information for people to understand the world. With the rapid development of digital imaging techniques and internet, more and more images are available to public. Consequently, there is an increasingly high demand for effective and efficient image indexing and classification methods. That is the reason image classification has become one of the most popular topics in the field of pattern recognition and image mining.

Recently CAPTCHA emerged as a powerful tool for security oriented system. CAPTCHA analysis and characterization was widely studied over the last three decades in a variety of applications, including medical imaging, remote sensing, pattern recognition, industrial inspection CAPTCHA based image retrieval, human visual perception, and provides information for recognition and interpretation. An image CAPTCHA is described by the number and types of its primitives and the spatial organization or layout of its

primitives. The spatial organization may be random, may have a pair-wise dependence of one primitive on a neighboring primitive, or may have a dependence of 'n' primitives at a time. CAPTCHA is undoubtedly one of the main features used in Image processing, pattern recognition and multispectral scanner images obtained from aircraft or satellite platforms for microscopic images of cell cultures or tissue samples. There are several other areas that make extensive use of textural features such as grain shapes, size, and distribution for classifying and analyzing specimens. CAPTCHA is very important in quality control since many inspection decisions are based on the appearance of the CAPTCHA of the material. That's why the research on CAPTCHA analysis has received considerable attention in recent years and it has become a subject of intense study for many researchers.

### 1.3 CAPTCHA AND PATTERN

CAPTCHA and pattern have been recognized as important attributes of image data. They are used extensively in the visual interpretation of image data, in which CAPTCHA is often more important than the other image attributes. The patterns have received somewhat less attention than the spectral characteristics of image data in digital image processing. The reason given for this lower priority has been usually the relatively coarse spatial resolution of sensor image data, while the apparent ease of analysis of the spectral data is probably just as significant. The recent advent of high-resolution image data acquisition systems and the increasing recognition of the difficulties inherent in the analysis of spectral data are providing strong incentives

to consider pattern as an important attribute of CAPTCHA analysis.

There is no universally accepted definition for the CAPTCHA. Some of the following definitions of CAPTCHA s are pattern based and they are given below.

- Faugeras and Pratt [Faugeras-1980] defined CAPTCHA as “the basic pattern and repetition frequency of a CAPTCHA sample could be perceptually invisible, although quantitatively present ... In the deterministic formulation CAPTCHA is considered as a basic local pattern that is periodically or quasi-periodically repeated over some area.”
- Jain and Karu [Jain-1996] defined “CAPTCHA is characterized not only by the grey value at a given pixel, but also by the grey value ‘pattern’ in a neighborhood surrounding the pixel.”

Generally speaking, CAPTCHA s are complex visual patterns composed of entities, or sub patterns that have characteristic brightness, color, slope, size, etc. Thus CAPTCHA can be regarded as a similarity grouping in an image [Rosenfeld-1982]. The local sub pattern properties give rise to the perceived lightness, uniformity, density, roughness, regularity, linearity, frequency, phase, directionality, coarseness, randomness, fineness, smoothness, granulation, etc., of the CAPTCHA as a whole [Levine-1985].

The word CAPTCHA certainly has many interpretations in the graphics community. The word CAPTCHA is used in the sense of a pattern applied to the surface

of an object. Intuitively, one can think of CAPTCHA as visual information which gives clues about the nature of the object, usually expressed at the object’s surface. The difference between a pattern and a CAPTCHA is that a CAPTCHA involves the attachment of the pattern to the surface of an object. Depending on the context, the word pattern has many different interpretations. The biology community seems to use the word pattern without defining it [Stev-74]. The implicit meaning generally brings to mind some kind of repeated arrangement (regular or not) and the term is often defined by examples.

The Oxford English Dictionary [Simp-1989] has 13 entries concerning the substantive “Pattern.”

The present study collected the following definitions of patterns:

- “The term pattern often used teleological to mean a set of objects which can in some (useful) way be treated alike. For each problem area we must ask: ‘what patterns would be useful for a machine working on such problems?’.” [Minsky-1961]
- “Two or more objects where the entire object, the pattern, cannot be predicted or deduced from any thing that can be known about any of its parts taken separately. For example, the words ‘ON’ and ‘NO’ can not be deduced from the individual letters ‘N’ and ‘O’. A (Statistical) interaction, a whole, a shape, a form, a gestalt.”[Uhr-1973]
- “A pattern is a distinctive combination of qualities, acts or tendencies”... “for a machine a

pattern consists of all data configurations which, by the application of the rules of transformation and matching, lead to a known standard representation and meaning.” ... “for a human it is probably not possible to give a better definition than: a pattern is something which somebody recognizes as a pattern.” [Giuliano-1967]

- “In their widest sense, patterns are the means by which we interpret the world.”... “Some patterns have a physical manifestation ... other patterns have only an abstract existence e.g. patterns in social or economic data.” [Meisel-1972]
- “A pattern is an ordering among elements which is such that (1) when the arrangement of a subset of these elements is determined, the range of possible arrangements of the remainder may be subdivided into two (not necessarily exhaustive) ranges of possibilities, any member of the first of which is more probable than any member of the second and (90) in general, as the subset of already determined elements increases with sufficiently large increments, the range of probable arrangements of the remainder decreases and the range of improbable arrangements increases.” [Sayre-1965]
- When one tries to apply mathematical methods in relations with patterns a much more formal and exact definition is necessary.

Grenander in his “foundations of patterns analysis” gives great attention to the definition of a pattern: “Our goal is to develop a model flexible enough to make it possible for us to discuss patterns in general: to give us a precise language in terms of which we shall be able to analyze and describe patterns.”

From the above the present study concludes that no abstract definition of pattern is given, but a certain distinct problem is solved, often without an explicit definition of the pattern used.

Study of patterns on CAPTCHA is recognized as an important step in characterization and classification of CAPTCHA. Various approaches are existing to investigate the textural and spatial structural characteristics of image data, including measures of CAPTCHA, Fourier analysis, fractal dimension, variograms and local variance measures. Fourier analysis is found as the most useful when dealing with regular patterns within image data. It has been used to filter out speckle in radar data and to remove the effects of regular agricultural patterns in image data. Study of regular patterns based on fundamentals of local variance was carried out recently.

That is the reason, the study of patterns still plays a significant area of research in classification and characterization of CAPTCHA s. That’s why the present thesis investigates how the frequency of occurrences of patterns varies after applying the preprocessing steps on the original CAPTCHA d image.

The present study assumes CAPTCHA is characterized not only by the grey value at

a given pixel, but also by the grey value pattern in a neighborhood surrounding the pixel. The ability to efficiently analyze and describe CAPTCHA patterns is thus of fundamental importance. A simple pattern of a neighborhood can be considered as one of the CAPTCHA primitive feature. Textural patterns can often be used to recognize familiar objects in an image or retrieve images with similar CAPTCHA from a database.

#### **1.4 CAPTCHA ANALYSIS AND CLASSIFICATION**

The analysis of CAPTCHA in images provides an important cue to the recognition of objects. It has been recently observed that different image objects are best characterized by different CAPTCHA methods. Since there are a lot of variations among natural CAPTCHA, to achieve the best performance for CAPTCHA analysis or retrieval, different features should be chosen according to the characteristics of CAPTCHA images. CAPTCHA, as a measure of the variation of the intensity of a surface, quantifies properties such as smoothness, coarseness and regularity. Image analysis usually refers to processing of images by computer with the goal of finding what objects are presented in the image. One immediate application of CAPTCHA is the recognition of image regions using CAPTCHA properties. In the past decades, numerous algorithms for CAPTCHA feature extraction have been proposed, many of which focus on extracting CAPTCHA features that are robust to noises, rotation and illumination variants. CAPTCHA classification system is

a computer system for browsing, searching, recognizing, comparing and classifying images from a large volume of digital images. The goal of CAPTCHA classification then is to produce a classification map of the input image where each uniform CAPTCHA region is identified by the CAPTCHA class to which it belongs. In CAPTCHA classification the first and most important task is to extract some features which efficiently embody information about the characteristics of the original image. These features can then be used for the classification of different images.

#### **Literature Survey**

Since its first appearance in 2000, a safety mechanism based on Completely Automated Public Turing Test to Tell Computers and Humans Apart (CAPTCHA) has been subjected to multiple attacks that seek to compromise their efficiency [1–5]. Therefore various security verification methods have been proposed covering a broad spectrum of options for generation of robust CAPTCHAs able to resist attacks by malicious programs [6–8]. Recently used methods are based on solutions of CAPTCHA imagerecognition tasks; text- or voice-processing; logical and mathematical puzzles, that besides offering a recognition challenge make the user apply some additional knowledge; even more, complex approaches that analyze patterns of clicks or face recognition [2,5,6,9]. However, text-based systems appear to be the most popular due to their easy implementation and usability. For this reason a set of design rules has been proposed to increase CAPTCHA security without compromising

the user experience [2,9–11]. This has significantly reduced the number of CAPTCHA APIs allowing only the most mature remain in preference of Web security providers. Those systems, which properly follow the CAPTCHA design guidelines, are currently used as safety mechanisms in high-traffic sites such as Facebook, Ticket Master, Gmail, Livenation, Uploading, CNN, YouTube and others [4,11–13].

Recently, two basic concepts are assumed to address automatic recognition of CAPTCHA: to break anti-segmentation techniques used to protect regions corresponding to characters and to overcome anti-recognition techniques for each character. While anti-recognition mechanisms alter individual letter features such as font size, type and count, distortion, blurring and independent rotation of each character, the main secure mechanism to avoid breaking CAPTCHAs relies on anti-segmentation techniques that guarantee their robustness [9,11,13]. Some of the principal anti-segmentation techniques used in new versions of CAPTCHA and reCAPTCHA are the variable orientation of characters in word, the collapse between letters in a word, the addition of random dots and lines of different sizes, cluttered backgrounds and similar foreground/background colors [2,3,9,10]. For example, reCAPTCHA test proposes to recognize two out-of-context words, where waviness and horizontal stroke were added to increase the difficulty of breaking the CAPTCHA by a computer program. According to Bursztein, the unpredictable collapse in CAPTCHA is the best option to avoid segmentation of characters now widely used by various sites

like Google, Facebook, Twitter and others [9,13]. Developed techniques for breaking CAPTCHAs are also used in pattern recognition applications particularly, for handwritten text interpretation or for optical character recognition (OCR) during automatic degraded text scanning. For example, the proposed approach provides simple and fast character recognition mechanism for scanning books in large scale such as Google Books and News Archive Search and their conversion to plain text [14,15].

Thus, the main purpose of this project is to reduce vulnerability of CAPTCHAs from frauds and to protect users against cyber-criminal activities as well as to introduce a novel approach for recognizing either handwritten or damaged texts in ancient books, manuscripts and newspapers.

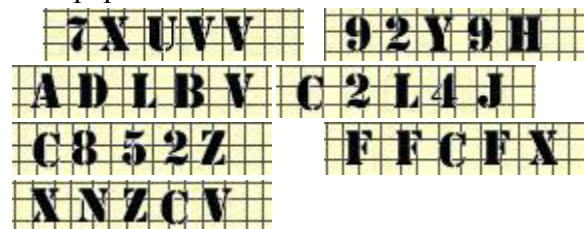


Fig. 2.1. Examples of security mechanisms in CAPTCHAs

Currently, there are numerous techniques breaking CAPTCHAs. The most complete analysis of CAPTCHA beating mechanisms provided by Bursztein presents various systems able to recognize CAPTCHAs of some popular Internet sites, which include Wikipedia, eBay, CNN, Baidu, Megaupload and others with accuracy rates ranging from 40% to 95% [9,11,13]. However, these systems do not recognize CAPTCHAs provided by sites like Google or

reCAPTCHA of new versions [10,13]. In Ref. [16] Yan presented attack on previous 2010 version of Google CAPTCHA, where character segmentation is based on analysis of patterns grouped in following categories: (1) point-shaped patterns (letter such as i or j); (2) cycle shaped patterns (letters a, b, d, etc.); (3) cross-shaped patterns (letters t and f) and (4) pattern, which juxta-poses three vertical lines to form the character (m or w). Although Google recently revamped its reCAPTCHA system nevertheless, Google's reCAPTCHA now is vulnerable once again after newly launched reCAPTCHA-solving/breaking service [1].

Several well-known approaches have broken CAPTCHAs such as Yahoo early CAPTCHAs [17], the CAPTCHAs used by PayPal site [18], Windows Hotmail and Gmail free e-mail providers [19], LiveJournal, php BB, e-banking CAPTCHAs used by a lot of financial institutions and other services [20,21]. After that, when the Newcastle Uni-versity research team has broken segmentation of Microsoft CAPTCHA with 90% of success rate, Microsoft uses improved CAPTCHA Control ASP.NET 2.0 [22]. This approach creates histogram of black pixels found in column assuming that characters are not overlapped and subsequently, defines letter separation point when no pixels are found in column. This segmentation approach fails, when characters are connected at least by a single pixel.

Recently reported in scientific literature approaches apply different algorithms to obtain CAPTCHA imageskeleton for easy manipulation of characters overcoming in this way anti-segmentation mechanisms [9,

11, 23]. The precision of the segmentation step reported by new CAPTCHA beating systems lies between 40% [10,13,16] and 95% [9,10,24]. Interesting approach is presented by Liu [25], which exploits a set of morphological filters that break satisfactorily security mechanism based on asymmetric-ellipses sometimes presented in reCAPTCHA. Another approach presented by Indian research group [26] considers that the pre-processing stage is not necessarily must generate complete letter blobs. It may be used only for fast global feature extraction however the correct segmentation task must be handled by the recognition module, which looks along the CAPTCHA imager to define the character boundaries. The proposed approach achieves recognition accuracy about 72% with response time less than 14.5s per 400 CAPTCHAs.

Although these results could give an idea that the problem is already solved, unfortunately, these reports frequently present theoretical proposal and have not formal evaluation of whole CAPTCHA breaking process. The main security mechanisms implemented in reCAPTCHA are focused on exploiting different font sizes, which suffer from a particular pattern of waving rotation and random collapse overlapping characters in words. As shown in Fig. 1 character-level blurring is also seen in certain areas. That represents a challenge for binarization and correct segmentation of characters. Additionally, some extra security features such as length and text-size randomization, character tilting and waving are used, which may guarantee that CAPTCHA scheme is secure against attacks [9,27].

Another requirement of systems for automatic CAPTCHA beating is providing high-speed recognition useful for real-time applications that not always are reported in well-known approaches. As usually, these CAPTCHA beating schemes apply the following stages: preprocessing for removal of background clutter and noise, segmentation for sub-division of CAPTCHA image into single regions and recognition of characters. The most difficult task is segmentation step although the development of fast and robust classifier is also a challenging task.

In this research we propose to subdivide the CAPTCHA breaking process into the following stages: CAPTCHA image acquisition, preprocessing, segmentation and recognition. The sequence of steps is presented in Fig. 2. In the next sections detailed description of the proposed steps is provided.

### **CAPTCHA Breaking Using Segmentation and Morphological Operations**

#### **3.1 Brief Outline**

Segmentation subdivides an CAPTCHA image into its constituent regions or objects. The point to which the subdivision is carried depends on the problem being solved. That is, segmentation should end when the objects of interest in an application have been isolated. Without a good segmentation algorithm, an object may never be identifiable [41]. Image segmentation continues to be an vital and active research area in image analysis [48]. Many techniques have been proposed to deal with the image segmentation problem. They

can be broadly grouped into the following categories. Histogram-Based Techniques, Edge-Based Techniques, Region-Based Techniques, Hybrid Techniques [48]. The accuracy of segmentation is highly dependent on the success or failure of each computerized analysis procedure. After the segmentation process is over, we should be familiar with, which pixel belongs to which object, the discontinuities where abrupt changes lie, tell us the locations of boundaries of regions. The connectedness of any two pixels is identified when there exists a connected path wholly within the set, where a connected path is a path that always moves between neighboring pixels. Therefore, region is a set of adjacent connected pixels. Extensive researches have been made in designing and creating different segmentation algorithms, however, still no algorithm is found from the researches results that can be accepted and appropriate for all kinds of images, obviously, all segmentation algorithms cannot be equally applicable to a certain application [39].

#### **3.2 Introduction**

The term morphology is generally attributed to the German poet, novelist, playwright, and philosopher Johann Wolfgang von Goethe (1749–1832), who coined it early in the nineteenth century in a biological context. Its etymology is Greek: morph- means ‘shape, form’, and morphology is the study of form or forms. In biology morphology refers to the study of the form and structure of organisms, and in geology it refers to the study of the configuration and evolution of land forms.

In linguistics morphology refers to the mental system involved in word formation or to the branch of linguistics that deals with words, their internal structure, and how they are formed.

*Morphology* is a broad set of image processing operations that process images based on shapes. Morphological operations apply a structuring element to an input image, creating an output image of the same size. In a morphological operation, the value of each pixel in the output image is based on a comparison of the corresponding pixel in the input image with its neighbors. By choosing the size and shape of the neighborhood, we can construct a morphological operation that is sensitive to specific shapes in the input image.

The most basic morphological operations are dilation and erosion. Dilation adds pixels to the boundaries of objects in an image, while erosion removes pixels on object boundaries. The number of pixels added or removed from the objects in an image depends on the size and shape of the *structuring element* used to process the image. In the morphological dilation and erosion operations, the state of any given pixel in the output image is determined by applying a rule to the corresponding pixel and its neighbors in the input image. The rule used to process the pixels defines the operation as dilation or an erosion. This table lists the rules for both dilation and erosion.

One of the most informative visual cues in CAPTCHA images is CAPTCHA . Despite that, CAPTCHA is regarded as a "fuzzy" concept with no mathematical or comprehensive definition agreed upon yet.

This may be due to the vague concept that CAPTCHA may actually hold resulting in the many interpretations related to human perception. The Oxford dictionary defines CAPTCHA with three different meanings: "the way a surface, substrate or fabric looks or feels to the touch, i.e. whether it is rough, smooth, hard, soft, etc; the way food or drink tastes or appears; the way in which a piece of music or literature is constructed, with regard to the way in which its parts are combined". CAPTCHA carries extensive information and plays an important role in our interpretation of a visual scene. A CAPTCHA is a measure of the variation of the surface intensity and quantifying properties such as density, regularity, uniformity, roughness, phase, linearity, directionality, frequency, randomness, smoothness, coarseness, fineness, granulation, etc., as a whole. The CAPTCHA images can be described with terms such as smoothness, grain, regularity or homogeneity. Various CAPTCHA metrics like variance, kurtosis, angular second moment, inertia, entropy etc. are used to quantify CAPTCHA images. Therefore, CAPTCHA images can be treated as CAPTCHA d images. CAPTCHA based CAPTCHA image segmentation offers reliable results and is considered to be very efficient. That's why the research on CAPTCHA analysis has received considerable attention in CAPTCHA image processing in recent years and it has become a subject of intense study for many researchers.

### 3.3 Effect of Noise in CAPTCHA images

In CAPTCHA image processing, CAPTCHA images are corrupted by different types of noises. It is very important to obtain precise images without noise to facilitate accurate observations for the given application. Removing of noise from CAPTCHA images is now a very challenging issue in the field of CAPTCHA image processing and many researchers are working in this area. Most well known noise reduction methods, which are usually based on the local statistics of a CAPTCHA image, are not efficient for CAPTCHA image noise reduction.

Low image quality and even a fractional noise images are obstacles for effective feature extraction, analysis, recognition and quantitative measurements especially in CAPTCHA image processing. Even though sometimes the quality is good, a small noise leads to an improper diagnosis by CAPTCHA expert thus leads to a false treatment, which may affect seriously the patient. Therefore, there is a fundamental need of noise reduction on CAPTCHA images.

### 3.4 CAPTCHA Image Segmentation

CAPTCHA Image processing can be defined as the manipulation of an CAPTCHA image for the purpose of either extracting accurate information from the image or producing an alternative representation of the image. There are numerous specific motivations for CAPTCHA image processing but many fall into the following categories: (i) to remove unwanted signal components that are corrupting the image and (ii) to extract

information by rendering it in a more obvious or more useful form.

CAPTCHA image segmentation is one of the most critical tasks of image analysis because the segmentation results will affect all the subsequent processes of image analysis, such as representation and description, feature measurement and even the following higher level tasks such as classification and interpretation. CAPTCHA image segmentation can be done using color, gray level, depth, CAPTCHA or any other feature of interest according to the specific application.

The segmentation algorithm depends on the envisioned application and on the imaging modality employed. For instance, segmentation of gray and white matter in a cerebral MRI induces vastly different constraints from that of a vertebrae in an X-ray of the vertebral column, in terms of target topology, prior knowledge, choice of target representation, signal to noise ratio, and dimensionality of the input data. The selection of an adequate segmentation paradigm is therefore pivotal as it affects how efficiently the segmentation system can deal with the target organ or structure and conditions its accuracy and robustness. A deeper understanding of both the anatomical characteristics of the tissues and organs of the human body (or, more precisely, of the sub-structures we distinguish within them) and of their inter-relationships is crucial in diagnostic and interventional medicine which is possible only by the effective and accurate segmentation process.

3.5 CAPTCHA SEGMENTATION

The CAPTCHA image segmentation also plays a vital role in various pattern recognition applications such as cartography, remote sensing, robot vision, military surveillance, inspection of textile products and CAPTCHA imaging. CAPTCHA segmentation divides an image into a set of distinct regions based on CAPTCHA properties, so that each region is uniform with respect to certain CAPTCHA characteristics. Results of segmentation can be further used to image processing and analysis, for example, to object recognition. Similar to classification, segmentation of CAPTCHA also involves extracting features and deriving metrics to separate CAPTCHA s. However, segmentation is generally more complex than classification, since boundaries that separate different CAPTCHA regions have to be detected in addition to recognizing CAPTCHA in each region. CAPTCHA segmentation could also be supervised or unsupervised based on if prior knowledge about the image or CAPTCHA class is available. Supervised CAPTCHA segmentation identifies and splits one or more regions that match CAPTCHA properties shown in the training CAPTCHA s. For CAPTCHA images unsupervised segmentation is preferred as the information to be segmented is not known in advance. Unsupervised segmentation has to first recover dissimilar CAPTCHA classes from an image before dividing them into regions. Compared to the supervised case, the unsupervised segmentation is more flexible for real world applications despite that it is generally more computationally expensive.

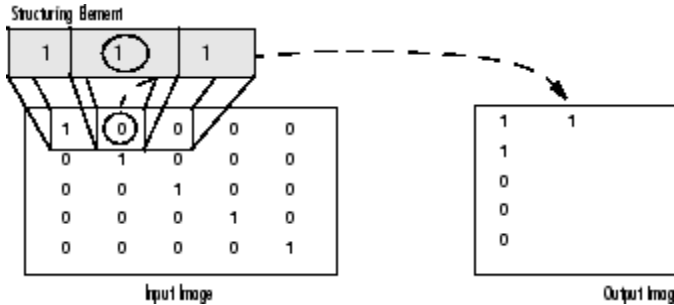
Segmenting an image into uniform regions is very useful in various applications of proper identification of abnormality, shape and volume in CAPTCHA image processing, pattern recognition and machine learning.

3.6 Rules for Dilation and Erosion

| Operation | Rule                                                                                                                                                                                 |
|-----------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Dilation  | The value of the output pixel is 1 if any of the pixels in the input pixel's neighborhood is set to the value 1, the output pixel is set to 1.                                       |
| Erosion   | The value of the output pixel is 1 only if all the pixels in the input pixel's neighborhood are set to 1. In a binary image, if any pixel is set to 0, the output pixel is set to 0. |

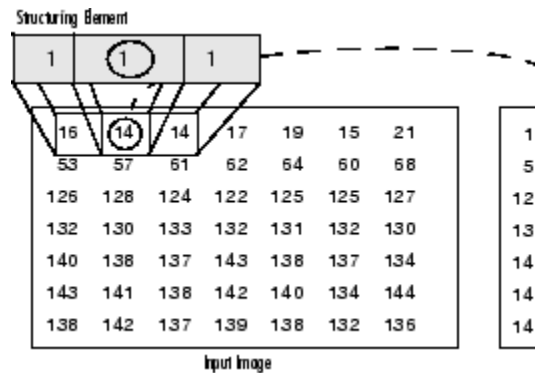
The following figure illustrates the dilation of a binary image. Note how the structuring element defines the neighborhood of the pixel of interest, which is circled. The dilation function applies the appropriate rule to the pixels in the neighborhood and assigns a value to the corresponding pixel in the output image. In the figure, the morphological dilation function sets the value of the output pixel to 1 because one of the elements in the neighborhood defined by the structuring element is on.

Morphological Dilation of a Binary Image



The following figure illustrates this processing for a grayscale image. The figure shows the processing of a particular pixel in the input image. Note how the function applies the rule to the input pixel's neighborhood and uses the highest value of all the pixels in the neighborhood as the value of the corresponding pixel in the output image.

**Morphological Dilation of a Grayscale Image**



**Processing Pixels at Image Borders (Padding Behavior)**

Morphological functions position the origin of the structuring element, its center element, over the pixel of interest in the input image. For pixels at the edge of an image, parts of the neighborhood defined by the structuring element can extend past the border of the image.

To process border pixels, the morphological functions assign a value to these undefined pixels, as if the functions had padded the image with additional rows and columns. The value of these padding pixels varies for dilation and erosion operations. The following table describes the padding rules for dilation and erosion for both binary and grayscale images.

**3.7 Combining Dilation and Erosion**

Dilation and erosion are often used in combination to implement image processing operations. For example, the definition of a morphological *opening* of an image is an erosion followed by a dilation, using the same structuring element for both operations. The related operation, morphological *closing* of an image, is the reverse: it consists of dilation followed by an erosion with the same structuring element.

The following section uses imdilate and imerode to illustrate how to implement a morphological opening. Note, however, that the toolbox already includes the imopen function, which performs this processing. The toolbox includes functions that perform many common morphological operations.

**3.7.1 Morphological Opening**

We can use morphological opening to remove small objects from an image while preserving the shape and size of larger objects in the image. For example, we can use the imopen function to remove all the circuit lines from the original images, creating an output image that contains only the rectangular shapes of the microchips.

Open=erosion followed by dilation

Close=dilation followed by erosion

**3.8 The proposed approach for character segmentation and recognition**

Based on analysis of character morphology in CAPTCHA alphabet, we grouped them by some characteristics in the following categories: Some characters with

circular regions such as a, b, d, e, g, o, p and q. These letters generate regions of 20 pixels wide. Characters with occurrence of more than one pixel per column for letters like c, e, f, k, s, t, z. There usually exist at least two pixels for each column.

u-shape pattern characters such as u, n, h are normally presented as pattern with two narrow sections with more than one pixel in column separated by a wide section of columns with only one pixel (the part of letter that connects vertical segments).

Characters of one pixel per column with slope representing letters like v, x, y with a slope about 451 7101.

Thin characters are letters such as i, j, l; they consist of a small vertical block of approximately 5 pixels wide.

r-shape pattern character is formed by narrow stripe with some pixels in column followed by a much larger section of columns with only one pixel.

Double Characters are letters m and w, which can be commonly confused with letter n or v only, when they are separated by column without black pixels.

The obtained three-color bar may be enhanced using some proposed rules. They must be applied one at a time and in the following order to avoid interference between them:

- 1.Noise reduction: if a bar in generated three-color bar code is black and it is only of one pixel wide, then the bar is replaced by a white bar.

- 2.Slope calculation: calculate the slope of white segments

- 3.m and w pattern matching: the m-type pattern can be wrongly interpreted as two consecutive n characters. This character is detected by a segment of pixels represented by two wide white bars separated by black bar. To find m-type pattern in the three-color bar code we run a template matching algorithm using template image. For each found m-type pattern, the corresponding region is segmented in stringent CAPTCHA image.

- 4.u pattern matching: this pattern represented by two black bars separated by a white bar is found in letters n, u, v, y and h. To find u-type pattern in three color bar code we run a template matching algorithm using the template. Then for each found patterns the corresponding region is segmented in stringent CAPTCHA image.

### 3.9 Proposed Algorithm:

Step1: Read color CAPTCH image

Step 2: Convert color CAPTCH into gray scale image

Step 3: By using threshold value extract the CAPTCH from the back ground of image

Step 4: Apply spatial filter for further cleaning

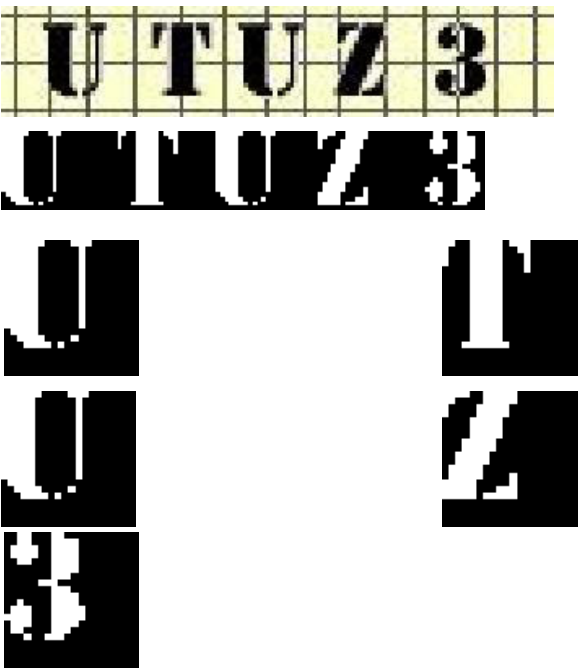
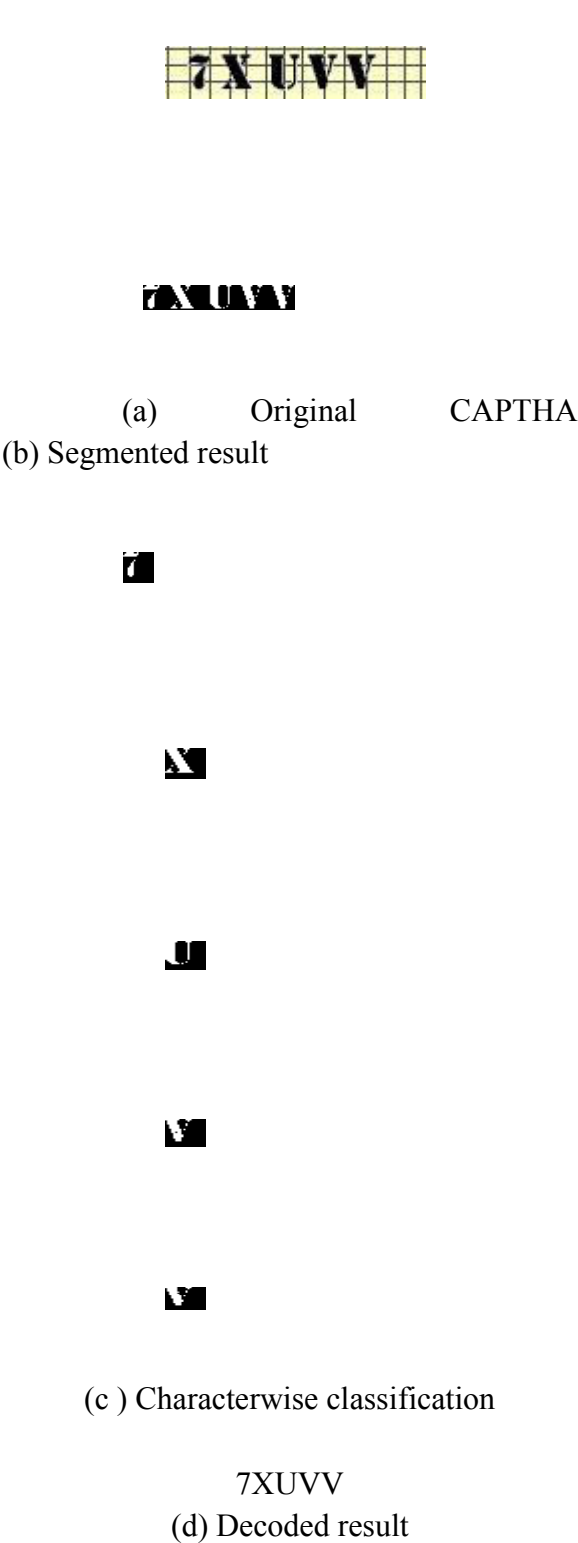
Step 5: Segment the cleaned CAPTCHA into individual characters

Step 6: Recognize the characters by using Template matching

Step 7: Decode the characters print recognized characters and display confidence

Confidence = 0.9798

3.10 Result and Discussions



decoded =UTUZ3  
confidence = 0.9995

4.5 Conclusions

In summary, three designs of text based CAPTCHA are proposed in this PROJECT. This CAPTCHA breaking design follows the principle “hard to separate text from background using segmentation techniques”. The CAPTCHAs are designed considering the techniques and concepts involved in cracking various existing CAPTCHAs. The proposed designs of CAPTCHA are thus too strong to get cracked using template matching and at the same time very user friendly with high confidence rate. The proposed method classification performance as follows

- Pixel Counting: 8% Break Rate

Fig.4.1: CAPTHA breaking system results.

- Vertical Projections: 97% Break Rate
- Horizontal Projections: 100% Break Rate
- Template Correlations: 100% Break Rate

## **Breaking Text-Based CAPTCHAS with Machine Learning Techniques**

### **4.1 Brief outline**

Weka is open source software for data mining under the GNU General public license. This system is developed at the University of Waikato in New Zealand. “Weka” stands for the Waikato Environment for knowledge analysis. Weka is freely available at <http://www.cs.waikato.ac.nz/ml/weka>. The system is written using object oriented language java. Weka provides implementation of state-of-the-art data mining and machine learning algorithm. User can perform association, filtering, classification, clustering, visualization, regression etc. by using weka tool. Each and every organization is accession vast and amplifying amounts of data in different formats and different databases at different platforms. This data provides any meaningful information that can be used to know anything about any object. Information is nothing just data with some meaning or processed data. Information is then converted to knowledge to use with KDD[43].

Data Mining is a non trivial extraction of implicit, previously unknown, and imaginable useful information from data. Data mining finds important information hidden in large volumes of data.

Data mining is the reasoning of data. It is the use of software techniques for finding patterns and consistency in sets of data. Data Mining is an interdisciplinary field involving: Databases, Statistics, and Machine Learning. There are various techniques available for data mining as given below:-

*A. Association Rule Learning:* - This is also called market basket analysis or dependency modelling. It is used to discover relationship and association rules among variables.

*B. Clustering:* - This technique creates and discovers group of similar data items. This is also called unsupervised classification.

*C. Classification:* - This can classify data according to their classes i.e. put data in single group that belongs to a common class. This is also called supervised classification.

*D. Regression:* - It tries to find a function that model the data with least errors.

*E. Summarization:* - It provides easy to understand and analysis facility through visualization, reports etc .

It is possible to mine data with computer that automates this process. Various data mining tools are available in market some are:-

- Environment for DeveLoping KDD-Applications Supported by Index-Structures (ELKI)
- jHepWork
- Konstanz Information Miner (KNIME)
- Orange (software)
- RapidMiner

- Scriptella ETL — ETL (Extract-Transform-Load) and script execution tool
- Weka [42]

## 4.2 Naive Bayes Classifier

In machine learning, naive Bayes classifiers are a family of simple probabilistic classifiers based on applying Bayes' theorem with strong (naive) independence assumptions between the features.

Naive Bayes has been studied extensively since the 1950s. It was introduced under a different name into the text retrieval community in the early 1960s,[1]:488 and remains a popular (baseline) method for text categorization, the problem of judging documents as belonging to one category or the other (such as spam or legitimate, sports or politics, etc.) with word frequencies as the features. With appropriate preprocessing, it is competitive in this domain with more advanced methods including support vector machines.

Naive Bayes classifiers are highly scalable, requiring a number of parameters linear in the number of variables (features/predictors) in a learning problem. Maximum-likelihood training can be done by evaluating a closed-form expression, which takes linear time, rather than by expensive iterative approximation as used for many other types of classifiers.

In the statistics and computer science literature, Naive Bayes models are known under a variety of names, including simple Bayes and independence Bayes. All these names reference the use of Bayes' theorem in the classifier's decision rule, but naive Bayes is not (necessarily) a Bayesian

method Russell and Norvig note that "[naive Bayes] is sometimes called a Bayesian classifier, a somewhat careless usage that has prompted true Bayesians to call it the idiot Bayes model."

The Naive Bayes Classifier technique is based on the so-called Bayesian theorem and is particularly suited when the dimensionality of the inputs is high. Despite its simplicity, Naive Bayes can often outperform more sophisticated classification methods.



Fig.4.1: Classification demo.

To demonstrate the concept of Naïve Bayes Classification, consider the example displayed in the illustration above. As indicated, the objects can be classified as either GREEN or RED. Our task is to classify new cases as they arrive, i.e., decide to which class label they belong, based on the currently existing objects.

Since there are twice as many GREEN objects as RED, it is reasonable to believe that a new case (which hasn't been observed yet) is twice as likely to have membership GREEN rather than RED. In the Bayesian analysis, this belief is known as the prior probability. Prior probabilities are based on previous experience, in this case the percentage of GREEN and RED objects, and often used to predict outcomes before they actually happen.

Prior probability GREEN  $\square$  Number of GREEN objects / Total number of Objects  
 Prior probability for RED  $\square$  Number of RED objects / Total number of Objects

Since there is a total of 60 objects, 40 of which are GREEN and 20 RED, our prior probabilities for class membership are:

Prior probability for GREEN  $\square$  40/60

Prior probability for RED  $\square$  20/60

Having formulated our prior probability, we are now ready to classify a new object (WHITE circle). Since the objects are well clustered, it is reasonable to assume that the more GREEN (or RED) objects in the vicinity of X, the more likely that the new cases belong to that particular color. To measure this likelihood, we draw a circle around X which encompasses a number (to be chosen a priori) of points irrespective of their class labels. Then we calculate the number of points in the circle belonging to each class label.

Although the prior probabilities indicate that X may belong to GREEN (given that there are twice as many GREEN compared to RED) the likelihood indicates otherwise; that the class membership of X is RED (given that there are more RED objects in the vicinity of X than GREEN). In the Bayesian analysis, the final classification is produced by combining both sources of information, i.e., the prior and the likelihood, to form a posterior probability using the so-called Bayes' rule (named after Rev.

Finally, we classify X as RED since its class membership achieves the largest posterior probability.

**Note.** The above probabilities is not normalized. However, this does not affect the classification outcome since their normalizing constants are the same.

## 4.2 LIBLINEAR

It is a linear classifier for data with millions of instances and features. It supports

L2-regularized classifiers

L2-loss linear SVM, L1-loss linear SVM, and logistic regression (LR)

L1-regularized classifiers (after version 1.4)

L2-loss linear SVM and logistic regression (LR)

L2-regularized support vector regression (after version 1.9)

L2-loss linear SVR and L1-loss linear SVR.

Main features of LIBLINEAR include

Same data format as LIBSVM, our general-purpose SVM solver, and also similar usage

Multi-class classification: 1) one-vs-the rest, 2) Crammer & Singer

Cross validation for model selection

Probability estimates (logistic regression only)

Weights for unbalanced data

MATLAB/Octave, Java, Python, Ruby interfaces

## 4.3 LibSVMs Linlinear

In practice the complexity of the SMO algorithm (that works both for kernel and linear SVM) as implemented in libsvm is  $O(n^2)$  or  $O(n^3)$  whereas liblinear is  $O(n)$  but does not support kernel SVMs.  $n$  is

the number of samples in the training dataset.

Hence for medium to large scale forget about kernels and use liblinear (or maybe have a look at approximate kernel SVM solvers such as LaSVM).

#### 4.4 MULTILAYER PERCEPTRON

A multilayer perceptron (MLP) is a feedforward artificial neural network model that maps sets of input data onto a set of appropriate outputs. A MLP consists of multiple layers of nodes in a directed graph, with each layer fully connected to the next one. Except for the input nodes, each node is a neuron (or processing element) with a nonlinear activation function. MLP utilizes a supervised learning technique called backpropagation for training the network. MLP is a modification of the standard linear perceptron and can distinguish data that are not linearly separable.

If a multilayer perceptron has a linear activation function in all neurons, that is, a linear function that maps the weighted inputs to the output of each neuron, then it is easily proved with linear algebra that any number of layers can be reduced to the standard two-layer input-output model (see perceptron). What makes a multilayer perceptron different is that some neurons use a nonlinear activation function which was developed to model the frequency of action potentials, or firing, of biological neurons in the brain. This function is modeled in several ways.

The multilayer perceptron consists of three or more layers (an input and an output layer with one or more hidden layers) of nonlinearly-activating nodes and is thus

considered a deep neural network. Each node in one layer connects with a certain weight  $w_{ij}$  to every node in the following layer. Some people do not include the input layer when counting the number of layers and there is disagreement about whether  $w_{ij}$  should be interpreted as the weight from  $i$  to  $j$  or the other way around.

Learning occurs in the perceptron by changing connection weights after each piece of data is processed, based on the amount of error in the output compared to the expected result. This is an example of supervised learning, and is carried out through back propagation, a generalization of the least mean squares algorithm in the linear perceptron.

The derivative to be calculated depends on the induced local field  $v_j$ , which itself varies. It is easy to prove that for an output node this derivative can be simplified to

This depends on the change in weights of the  $k$ th nodes, which represent the output layer. So to change the hidden layer weights, we must first change the output layer weights according to the derivative of the activation function, and so this algorithm represents a back propagation of the activation function.

The next architecture we are going to present using Theano is the single-hidden-layer Multi-Layer Perceptron (MLP). An MLP can be viewed as a logistic regression classifier where the input is first transformed using a learnt non-linear transformation  $\Phi$ . This transformation projects the input data into a space where it becomes linearly separable. This intermediate layer is referred to as a hidden layer. A single hidden layer is

sufficient to make MLPs a universal approximator. However we will see later on that there are substantial benefits to using many such hidden layers, i.e. the very premise of deep learning. See these course notes for an introduction to MLPs, the back-propagation algorithm, and how to train MLPs.

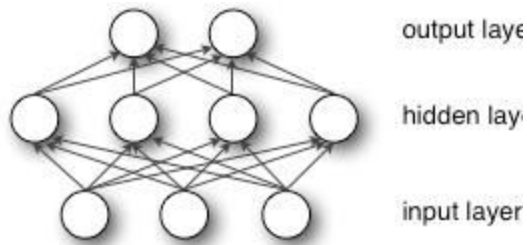


Fig.4.2: Neural network.

The output vector is then obtained as:  $o(x) = G(b^{(2)} + W^{(2)} h(x))$ . The reader should recognize the form we already used for Classifying MNIST digits using Logistic Regression. As before, class-membership probabilities can be obtained by choosing  $G$  as the softmax function (in the case of multi-class classification).

To train an MLP, we learn all parameters of the model, and here we use Stochastic Gradient Descent with mini batches. The set of parameters to learn is the set  $\theta =$

$$\{W^{(2)}, b^{(2)}, W^{(1)}, b^{(1)}\}.$$

Obtaining the gradients  $\frac{\partial \ell}{\partial \theta}$  can be achieved through the backpropagation algorithm (a special case of the chain-rule of derivation).

#### 4.5 IBK

Ibk algorithm, implements the  $k$ -nearest neighbor algorithm. Nearest-neighbor learning is also known as "Instance-based" learning.

K-Nearest Neighbors, or KNN, is a family of simple:

Classification and regression algorithms based on Similarity (Distance) calculation between instances. Nearest Neighbor implements rote learning. It's based on a local average calculation. It's a smoother algorithm.

Some experts have written that  $k$ -nearest neighbours do the best about one third of the time. It's so simple that, in the game of doing classification, we always want to have it in our toolbox.

#### 4.6. Proposed Algorithm

1. Read CAPTHA image and convert it into grayscale image handle in  $I$
2. Apply pre-processing on gray scale image  $I$ 
  - As part of pre processing use global threshold and extract only characters from the back ground.
3. Apply segmentation on resulted image of step 2 and identify the characters region wise.
4. Crop the character along with region with the size 20x20 (Max size to hold the character)
5. Calculate number of ones in each column and considered as feature vector. So that 1x20 feature vector is obtained for each character.
6. Design the feature vector for all alphabets (A-Z and numbers (0-9))

7. Apply the machine learning algorithms for recognizing the character and returns corresponding recognition rate.

4.7 Results and Discussions

| 4.7 Results and Discussions |    |    |    |    |    |    |    |    |     |     |     | 7   | 7   | 8   | 10  | 12  | 11  | 12  | 10  | 9     | 6  | 6  | 0  | 0  |    |   |   |
|-----------------------------|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-------|----|----|----|----|----|---|---|
|                             |    |    |    |    |    |    |    |    |     |     |     | 4   | 4   | 3   | 3   | 16  | 16  | 16  | 16  | 16    | 1  | 1  | 5  | 5  |    |   |   |
| F1                          | F2 | F3 | F4 | F5 | F6 | F7 | F8 | F9 | F10 | F11 | F12 | F13 | F14 | F15 | F16 | F17 | F18 | F19 | F20 | label |    |    |    |    |    |   |   |
| 1                           | 3  | 3  | 3  | 1  | 4  | 8  | 13 | 16 | 16  | 12  | 8   | 14  | 6   | 1   | 10  | 0   | 9   | 0   | 2   | 16    | 0  | 13 | A  | 10 | 2  | 0 |   |
| 3                           | 16 | 16 | 16 | 16 | 0  | 0  | 2  | 6  | 16  | 14  | 13  | 20  | 1   | 2   | 1   | 0   | 2   | 0   | 0   | 150   | 16 | 0  | 16 | B  | 2  | 1 | 0 |
| 8                           | 12 | 14 | 14 | 15 | 1  | 1  | 0  | 0  | 1   | 3   | 5   | 50  | 9   | 0   | 8   | 0   | 3   | 0   | 3   | 110   | 10 | 0  | 7  | C  | 2  | 0 | 0 |
| 1                           | 3  | 16 | 16 | 16 | 16 | 0  | 0  | 1  | 4   | 15  | 13  | 31  | 10  | 9   | 13  | 0   | 15  | 0   | 15  | 30    | 13 | 0  | 13 | D  | 11 | 6 | 0 |
| 1                           | 3  | 16 | 16 | 16 | 16 | 0  | 0  | 2  | 7   | 3   | 7   | 11  | 2   | 0   | 3   | 0   | 8   | 0   | 10  | 30    | 2  | 0  | 2  | E  | 0  | 0 | 0 |
| 1                           | 3  | 16 | 16 | 16 | 16 | 3  | 1  | 1  | 6   | 1   | 3   | 20  | 11  | 0   | 13  | 0   | 15  | 0   | 16  | 30    | 16 | 0  | 14 | F  | 11 | 5 | 0 |
| 7                           | 11 | 13 | 15 | 16 | 2  | 0  | 0  | 3  | 9   | 9   | 10  | 42  | 10  | 0   | 13  | 0   | 13  | 0   | 6   | 30    | 15 | 0  | 13 | G  | 11 | 6 | 0 |
| 2                           | 16 | 16 | 16 | 15 | 16 | 1  | 1  | 2  | 16  | 16  | 16  | 16  | 16  | 1   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | H  |    |    |   |   |
| 4                           | 4  | 3  | 3  | 16 | 16 | 16 | 16 | 16 | 1   | 1   | 5   | 5   | 9   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | I  |    |    |   |   |
| 1                           | 5  | 5  | 5  | 1  | 0  | 16 | 16 | 15 | 15  | 14  | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | J  |    |    |   |   |
| 2                           | 16 | 16 | 16 | 16 | 16 | 1  | 2  | 8  | 12  | 13  | 12  | 7   | 5   | 1   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | K  |    |    |   |   |
| 1                           | 3  | 16 | 16 | 16 | 16 | 1  | 0  | 1  | 2   | 3   | 5   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | L  |    |    |   |   |
| 1                           | 3  | 11 | 7  | 9  | 12 | 15 | 11 | 4  | 1   | 1   | 16  | 16  | 16  | 16  | 16  | 1   | 0   | 0   | 0   | 0     | 0  |    | M  |    |    |   |   |
| 1                           | 2  | 12 | 5  | 6  | 7  | 9  | 10 | 9  | 10  | 10  | 8   | 8   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | N  |    |    |   |   |
| 1                           | 3  | 16 | 16 | 16 | 16 | 0  | 0  | 1  | 4   | 15  | 13  | 11  | 9   | 0   | 0   | 0   | 0   | 0   | 0   | 0     |    | O  |    |    |    |   |   |
| 1                           | 3  | 16 | 16 | 16 | 15 | 3  | 1  | 0  | 3   | 9   | 9   | 7   | 5   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | P  |    |    |   |   |
| 1                           | 3  | 16 | 16 | 16 | 16 | 0  | 0  | 1  | 4   | 15  | 13  | 11  | 9   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | Q  |    |    |   |   |
| 1                           | 3  | 16 | 16 | 16 | 16 | 2  | 1  | 0  | 3   | 14  | 15  | 14  | 11  | 1   | 1   | 0   | 0   | 0   | 0   | 0     | 0  |    | R  |    |    |   |   |
| 9                           | 10 | 10 | 10 | 9  | 6  | 7  | 7  | 10 | 11  | 11  | 3   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | S  |    |    |   |   |
| 4                           | 1  | 0  | 1  | 16 | 16 | 16 | 16 | 16 | 1   | 1   | 2   | 5   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | T  |    |    |   |   |
| 13                          | 14 | 15 | 16 | 16 | 3  | 1  | 0  | 1  | 0   | 2   | 13  | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | U  |    |    |   |   |
| 1                           | 4  | 7  | 10 | 14 | 16 | 11 | 9  | 2  | 1   | 3   | 2   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | V  |    |    |   |   |
| 1                           | 4  | 8  | 11 | 15 | 15 | 12 | 5  | 0  | 7   | 12  | 16  | 14  | 9   | 1   | 2   | 3   | 1   | 0   | 0   | 0     | 0  |    | W  |    |    |   |   |
| 2                           | 5  | 8  | 9  | 9  | 11 | 12 | 12 | 10 | 10  | 9   | 5   | 2   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | X  |    |    |   |   |
| 1                           | 4  | 7  | 10 | 16 | 16 | 12 | 10 | 8  | 2   | 2   | 2   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0     | 0  |    | Y  |    |    |   |   |

**Naviey Bais Output:**  
Kappa statistic  
-0.0294  
Mean absolute error  
0.0571  
Root mean squared error  
0.2336  
Relative absolute error  
101.5298 %  
Root relative squared error  
138.2426 %  
Coverage of cases (0.95 level)  
0 %  
Mean rel. region size (0.95 level)  
3.4286 %  
Total Number of Instances  
35  
Ignored Class Unknown Instances

1

**Multilayer perceptron Output:**  
Kappa statistic  
-0.0294  
Mean absolute error  
0.0569  
Root mean squared error  
0.1923  
Relative absolute error  
101.148 %  
Root relative squared error  
113.8114 %  
Coverage of cases (0.95 level)  
2.8571 %  
Mean rel. region size (0.95 level)  
26.6939 %  
Total Number of Instances  
35  
Ignored Class Unknown Instances

1

**LibLinear Output:**  
Kappa statistic  
-0.0294  
Mean absolute error  
0.0571  
Root mean squared error  
0.239  
Relative absolute error  
101.5298 %  
Root relative squared error  
141.4798 %  
Coverage of cases (0.95 level)  
0 %  
Mean rel. region size (0.95 level)  
2.8571 %  
Total Number of Instances  
35  
Ignored Class Unknown Instances

1

**IBK Output:**  
Kappa statistic  
-0.0294  
Mean absolute error  
0.0563  
Root mean squared error  
0.1861  
Relative absolute error  
100.0987 %  
Root relative squared error  
110.1262 %  
Coverage of cases (0.95 level)  
88.5714 %  
Mean rel. region size (0.95 level)  
90.6122 %  
Total Number of Instances  
35  
Ignored Class Unknown Instances

1

**J48 output:**

|                                    |           |         |       |
|------------------------------------|-----------|---------|-------|
| Kappa                              | statistic |         |       |
| -0.0294                            |           |         |       |
| Mean                               | absolute  | error   |       |
| 0.0571                             |           |         |       |
| Root                               | mean      | squared | error |
| 0.2017                             |           |         |       |
| Relative                           | absolute  | error   |       |
| 101.5298 %                         |           |         |       |
| Root                               | relative  | squared | error |
| 119.3728 %                         |           |         |       |
| Coverage of cases (0.95 level)     |           |         |       |
| 0 %                                |           |         |       |
| Mean rel. region size (0.95 level) |           |         |       |
| 7.102 %                            |           |         |       |
| Total Number of Instances          |           |         |       |
| 35                                 |           |         |       |
| Ignored Class Unknown Instances    |           |         |       |

1

**4.8 Conclusions**

This projects presents discussion about machine learning approaches. Weka is used as data mining tool that provides various algorithms to be applied on data sets. The J48 algorithm is used to implement Univariate Decision Tree approach, while its results are discussed. The Multivariate approach is introduced as the Linear Machine approach that makes the use of the Absolute Error Correction and also the Thermal Perceptron Rules. Decision Tree is a popular technique for supervised classification, especially when the results are interpreted by human. Multivariate Decision Tree uses the concept of attributes correlation and provides the best way to perform conditional tests as compare to Univariate approach. The project concludes that multi layer perceptron approach is far

better than remaining approaches while it allow us dealing with large amount of data.

**FUTURE SCOPE**

This model may extent to solve for different CAPTCHA designs- BarCAPTCHA, TransparentCAPTCHA and ThreadCAPTCHA. The obtained very satisfactory results confirm that the proposed approach may be considered as robust techniques for CAPTCHA improvement as well as a new way of developing systems able to protect users against cyber-criminal activities and Internet threats.

**REFERENCES**

[1] D. Danchev, Google's reCAPTCHA under automatic fire from a newly launched reCAPTCHA-solving/breaking service, Internet Security Threat Updates & Insights, <http://www.webroot.com/blog/2014/01/21/googles-recaptcha-automatic-fire-newly-launched-recaptcha-solving-breaking-service/>, 2014.

[2] C. Obimbo, A. Halligan, P. De Freitas, Captch All: an improvement on the modern textbased CAPTCHA, J. Procedia Comput. Sci. 20 (2013) 496–501.

[3] G. Baxter Bell, Strengthening CAPTCHA-based Web security, First Monday J. 17 (2012) 2 <http://firstmonday.org/ojs/index.php/fm/article/view/3630>.

[4] S. Kulkarni, H.S. Fadewar, CAPTCHA based web security: an

- overview, *Int. J. Adv. Res. Comput. Sci. Softw. Eng.* 3 (11) (2013) 154–158 ([http://www.ijarcsse.com/docs/papers/Volume\\_3/11\\_November2013/V3I110-0379.pdf](http://www.ijarcsse.com/docs/papers/Volume_3/11_November2013/V3I110-0379.pdf)).
- [5] M. Serrao, S. Salunke, A. Mathur, Cracking CAPTCHAs for cash: a review of CAPTCHA crackers, *Int. J. Eng. Res. Technol.* 2 (1) (2013) 1–5.
- [6] G. Goswami, B.M. Powell, M. Vatsa, R. Singh, A. Noore, FaceDCAPTCHA: face detection based color image CAPTCHA, *Futur. Gener. Comput. Syst.* 31 (2014) 59–68.
- [7] L.D. Priya, S. Karthik, Secure captcha input based spam prevention, *Int. J. Emerg. Sci. Eng.* 1 (7) (2013) 9–12.
- [8] S. Azad, K. Jain, CAPTCHA: attacks and weaknesses against OCR technology, *Global J. Comput. Sci. Technol. Neural Artif. Intell.* 13 (3) (2013) 14–18.
- [9] E. Bursztein, A. Moscicki, C. Fabry, S. Bethard, J.C. Mitchell, D. Jurafsky, Easy does it: more usable CAPTCHAs, in: *Proceedings of the 32nd ACM Conference on Human Factors in Computing Systems*, Canada, 2014, pp. 2637–2646, <http://dx.doi.org/10.1145/2556288.2557322>.
- [10] C. Cruz-Perez, O. Starostenko, F. Uceda-Ponga, V. Alarcon-Aquino, L. Reyes-Cabrera, Breaking reCAPTCHAs with unpredictable collapse: heuristic character segmentation and recognition, in: J.A. Carrasco-Ochoa, J.F. Martinez-Trinidad, J.A. Overa Lopez, K. Boyer (Eds.), *LNCS: Pattern Recognition*, 7329, Springer-Verlag, Berlin Heidelberg, 2012, pp. 155–165.
- [11] E. Bursztein, M. Matthieu, Text-based CAPTCHA strengths and weaknesses, in: *Proceedings of the 18th ACM Conference on Computer and Communications Security*, IL, USA, 2011, pp. 125–138, (<http://ly.tl/p22>).
- [12] K. Fang, Z. Bu, Z.Y. Xia, Segmentation of CAPTCHAs based on complex networks, in: J. Lei, F. Lee Wang, H. Deng, D. Miao (Eds.), *LNCS: Artificial Intelligence and Computational Intelligence*, 7530, 2012, pp. 735–743.
- [13] E. Bursztein, H. Paskov, et al., Science of CAPTCHAs: solvability by humans and machines, AFOSR MURI Project, 2013, pp. 1–59.
- [14] Committee on Institutional Cooperation. Google Book Search Project, (<http://www.cic.net/projects/library/book-search/introduction>), 2014.
- [15] C. Lim Tan, X. Zhang, L. Li, Image based retrieval and keyword spotting in documents, in: D. Doermann, K. Tompa (Eds.), *Handbook of Document Image*

- Processing and Recognition, Springer-Verlag, London, 2014, pp. 805–842.
- [16] J. Yan, A. Salah, E. Ahmad, The robustness of a new CAPTCHA, in: 3rd Workshop on System Security, NY, USA, 2010, pp. 36–41, (<http://doi.acm.org/10.1145/1752046.1752052>).
- [17] M. Wehner, Internet advertisers kill text-based CAPTCHA, (<http://news.yahoo.com/internet-advertisers-kill-text-based-captcha-205416291.html>), 2013.
- [18] K. Kluever, R. Zanibbi, Breaking the PayPal CAPTCHA, (<http://www.kloover.com/2008/05/12/breaking-the-paypalcom-captcha/>), 2014 (retrieved on 25th of May).
- [19] K. Dawson, Windows Live Hotmail CAPTCHA Cracked, Exploited, (<http://tech.slashdot.org/article.pl?sid=08/04/15/1941236&from=rss>), and Gmail CAPTCHA Cracked, (<http://it.slashdot.org/article.pl?sid=08/02/27/0045242>), 2014.
- [20] S. Li, A. Syed, et. al., Breaking e-Banking CAPTCHAs, in: Proceedings of the 26th Computer Security Applications Conference, NY, USA, 2010, pp. 171–180, ([http://www.acsac.org/2010/openconf/modules/request.php?module=oc\\_program&action=summary.php&id=53](http://www.acsac.org/2010/openconf/modules/request.php?module=oc_program&action=summary.php&id=53)).
- [21] S. Kruglov, Defeating of weak CAPTCHAs, (<http://www.captcha.ru/en/breakings/>), 2013.
- [22] Microsoft ASP.NET Team, Using a CAPTCHA to Prevent Bots from Using our ASP.NET Web Razor) Site, ([http://www.asp.net/web-pages/tutorials/security/using-a-captcha-to-prevent-automated-programs-\(bots\)-from-using-your-aspnet-web-site](http://www.asp.net/web-pages/tutorials/security/using-a-captcha-to-prevent-automated-programs-(bots)-from-using-your-aspnet-web-site)), 2012.
- [23] S.E. Ahmad, J. Yan, M. Tayara, The Robustness of Google CAPTCHAs, Newcastle University Print, England, 2011 (Technical report).
- [24] A. Baluni, S. Gole, Two-step CAPTCHA: using a simple two step turing test to differentiate between humans and bots, Int. J. Comput. Appl. 81 (16) (2013) 48–51.
- [25] P. Liu, J. Shi, L. Wang, L. Guo, An efficient ellipse-shaped blobs detection algorithm for breaking facebook CAPTCHA, in: Y. Yuan, X. Wu, Y. Lu (Eds.), CCIS: Trustworthy Computing and Services, 320, Springer-Verlag, Berlin Heidelberg, 2013, pp. 420–428.

- [26] D. Kapoor, H. Bangar, A. Chaurasia, A. Sethi, An ingenious technique for symbol identification from high noise CAPTCHA images, in: Proceedings of the Annual IEEE India Conference, 2012, pp. 98–103.
- [27] M. Takaya, H. Kato, T. Komatsubara, Y. Watanabe, A. Yamamura, Recognition of one-stroke symbols by humans and computers, J. Procedia – Soc. Behav. Sci. 97 (6) (2013) 666–674.
- [28] Barbu.A and Zhu.S.C, “Multigrid and multi-level swendsen-wang cuts for hierarchic graph partition”, In Computer Vision and Pattern Recog.(CVPR), vol. 2, pp. 731–738, July 2004.
- [29] Brox.T and Weickert.J, "Level set segmentation with multiple regions", in IEEE Trans. Image Processing, vol.15, no.10, pp. 3213–3218, 2006.
- [30] Christoudias.C, Georgescu.B and Meer.P, "Synergism in low level vision", In Proc. of the 16th Int. Conf. on Pattern Recog., Los Alamitos: IEEE Computer Society, vol. 4, pp. 150–155, 2002.
- [31] Davies.E.R , “Handbook of texture analysis”, Imperial college press, Chapter 1, 2009.
- [32] Jean Serra, Pierre Soille, "Mathematical morphology and its applications to image processing", Springer, 1 edition, 1994.
- [33] Kekre.H.B, Saylee Gharge, "Texture Based Segmentation using Statistical Properties for Mammographic Images”, Int. Journal of Advanced Computer Science and Applications(IJACSA), Vol. 1, No. 5, pp. 102-107, Nov 2010.
- [34] Lutz Goldmann, Tomasz Adamek, Peter Vajda, “Towards Fully Automatic Image Segmentation Evaluation”, in Advanced Concepts for Intelligent Vision Systems Lecture Notes in Comp. Science, vol. 5259, pp. 566–577, 2008.
- [35] Mounir Sayadi, Lotfi Tlig and Farhat Fnaiech, “A new texture segmentation method based on the fuzzy C-mean algorithm and statistical features”, Applied Mathematical Sciences, vol. 1, no. 60, pp. 2999 – 3007, 2007.
- [36] Scarpa.G and Haindl.M, "Unsupervised texture segmentation by spectral spatial independent clustering," In Proc. of the 18th Int. Conf. on Pattern Recog., Los

Alamitos: IEEE Computer Society,  
vol. 2, pp. 151–154, 2006.

[37] Shi.J and Malik.J,  
"Normalized cuts and image  
segmentation," in IEEE Trans.  
Pattern Anal. Mach. Intell., vol. 22,  
no. 8, pp. 888–905, Aug.2000.

[38] Singh.S and Sharma.M,  
"Texture experiments with Meastex  
and Vistex benchmarks," in  
Proc.Inter. Conf. on Advances in  
Pattern recog., Lecture Notes in  
Computer Science, 2013.

[39] Ahmed R. Khalifa et  
al., "Evaluating The Effectiveness Of  
Region Growing And Edge  
Detection Segmentation  
Algorithms",. Journal of American  
Science, 2010;6(10),pp.580-587.

[40] Kostas Haris et al., "Hybrid  
Image Segmentation Using

Watersheds and Fast Region  
Merging", IEEE Transactions On  
Image Processing, Vol. 7, No. 12,  
December 1998,pp.1684-1699.

[41] S.Lakshmi et. al., "A study of  
Edge Detection Techniques for  
Segmentation Computing  
Approaches", IJCA Special Issue on  
"Computer Aided Soft Computing  
Techniques for Imaging and  
Biomedical Applications" CASCT,  
2010, pp.35-41.

[42] "Data Mining" from  
Wikipedia the free Encyclopedia.  
Web.  
<[http://en.wikipedia.org/wiki/Data\\_mining](http://en.wikipedia.org/wiki/Data_mining)>.

[43] Term "INTRODUCTION OF  
DATA MINING", "Data Mining:  
What is Data Mining ", source from  
[http://www.anderson.ucla.edu/](http://www.anderson.ucla.edu/faculty/jason.frand/teacher/technologies/palace/datamining.htm)  
faculty/jason.frand /  
teacher/technologies/palace/datamini  
ng.htm.

## UNSTRUCTURAL DATA USING HADOOP

K.Madan Mohan<sup>1</sup> Asst.Professor,Malla Reddy college of Engineering,Hyderabad

R.Bangari<sup>2</sup> Asst.Professor,Malla Reddy college of Engineering,Hyderabad

### ABSTRACT:

Big data came into existence when the traditional relational database systems were not able to handle the unstructured data (weblogs, videos, photos, social updates, human behavior) generated today by organization, social media, or from any other data generating source. Data is increasing in size day by day and Hadoop is

used to process such large amount of data. In is paper, I made a study of various security issues associated with big data in context with the Hadoop environment and the various solution techniques and technologies involve in securing the big data Hadoop.

---

Keywords: Big Data, SASL, delegation, cell level, variety, unauthorized.

---

### 1.INTRODUCTION

Big data means data which is large in size, volume, variety. Nowadays the size a data is increasing rapidly, use of social media, Smartphone's, online shopping's etc. The volumes of Big data are on a roll, which can be inferred from the fact that as far back in the year 2012, there were a few dozen terabytes of data in a single dataset, which has interestingly been catapulted to many petabytes today. Such large amount of data is used for commercial purpose by enterprise to increase their business profit and many other applications, and therefore there is a need to secure such large amount of data and its processing. Big data has the following characteristics: Volume: In Big data the word big it-self define the size the data. Volume is associated with the size of big data. At• present the data is supposed to be petabytes with could increase to zettabytes in near future. Velocity: Velocity in Big data deals with the speed of the data coming from various sources. Velocity characteristic• is not limited to the speed of incoming data but also speed at which the data flows and aggregated. Variety: Data variety is a

measure of the richness of the data representation – text, images video, audio, etc. the• data processed is not of a single type it consists of semi structured data and unstructured data. Value: Data value measures the usefulness of data for making decisions. The data science is useful in getting to• know the data, but “analytic science” encompasses the predictive power of big data. Various users can run certain queries against the data stored and thus can deduct important results from the filtered data obtained and also rank it according to the dimensions they require. These reports help people to find the business trends according to which they can make change in their strategies. Complexity: Complexity measures the degree of a interdependence in big data structures such that a small change• (or combination of small changes) in one or a few elements can yield very large changes or a small change that ripple across or cascade through the system and substantially affect its behavior, or no change at all (Katal, Wazid, & Goudar, 2013) and interconnectedness (possibly very large).

### III. PROCESSING BIG DATA

Hadoop allows running applications on systems with thousands of nodes with thousands of terabytes of data [2]. The distributed file system supports fast data transfer rates among nodes allowing the system to continue operating uninterrupted at times of node failure. Hadoop has of distributed file system, analytics platforms and data storage and a layer handling parallel computation, rate of flow (workflow) and configuration administration [8].the HDFS runs across the nodes in a Hadoop cluster with together connects the file systems on many input and output data nodes and make one big file system [2]. Hadoop ecosystem have Hadoop kernel, Map- Reduce, the Hadoop distributed file system (HDFS) and a number of related components such as Apache, Oozie, Hive, HBase ,Pig and Zookeeper and these

components that are explained as below[7,8]:  
HDFS: A high faults tolerant distributed file system which is responsible for storing data on the clusters.● MapReduce: highly powerful parallel programming technique used for distributed processing of vast amount of● data on clusters. Hbase: which is a column oriented distributed NoSQL database used for random read/write access.● Pig: analyzing data of Hadoop computation pig is a high level data programming language● Hive: Is a data warehousing application which provides a SQL like access and relational model.● Sqoop: A project used for transferring/importing data between relational databases and Hadoop.● Oozie: An orchestration and a workflow management for dependent Hadoop jobs.●

### IV. BIG DATA HADOOP'S TRADITIONAL SECURITY

#### A.OVERVIEW

Originally Hadoop was developed without any security in mind, no security model, no authentication of users and services and no data privacy, so anybody could submit arbitrary code to be executed. Auditing and authorization controls (HDFS file permissions and ACLs) used during earlier distributions were easily evade because any user could impersonate other user. So various security controls measures that did subsist were not very effective. Later authorization and authentication was added, but had some weakness. All programmers users and

had the same level of access privileges to all the data in the cluster, any one could access any of the data in the cluster, and any user could read any data set [4]. MapReduce had no concept of authentication or authorization, user could lower the priorities of other Hadoop jobs to make his job complete faster or to be executed first – or worse, he could kill the other jobs. Supplementary security cannot keep up. The Hadoop supports some security features with current Kerberos implementation, firewalls, and HDFS permissions and ACLs.

## B. THREATS TO SECURITY

The related threats associated with processing data in Hadoop ecosystem are as follows: 1. Unauthorized client: An unauthorized client may write/read a data block of a file at a Data Node using the pipeline streaming Data-transfer protocol. And if gained access privileges and can submit a job to a queue or delete or change priority of the job. And can access intermediate data of Map job via its task trackers HTTP shuffle protocol. 2. Task: A task in execution may make use of host OS interfaces to access other tasks, or would access local data which include intermediate Map output or the local storage of the Data Node that runs on the same physical node. Similarly, A task or node may

masquerade as a Hadoop service component such as a Name Node, job tracker, Data Node, task tracker etc.3. Unauthorized user: An unauthorized user could execute arbitrary code or carry out further attacks by accessing an HDFS file via the RPC or via HTTP protocols. Similarly, he may sniff/ eavesdrop to data packets being sent by Data nodes to client and can submit a workflow to Oozie as another user. Data Nodes does not impose access control, he could bypass the various access control mechanism or restrictions and read arbitrary data blocks from Data Nodes or writing garbage data to Data Node to corrupt it.

## V. SECURITY ISSUES

Hadoop present security issues for data centre managers and security professionals. The security issues are as below [5, 6, 18]: 1. Fragmented Data: Big Data clusters contain data that allow multiple copies moving to-and-fro various nodes ensuring redundancy and resiliency. The data that is available for fragmentation and can be shared across multiple servers more complexity is added as a result of the fragmentation which poses a security issue due to the absence of a security model. 2. Distributed Computing: the data source is not fixed resources are processed where available, these lead to large levels of parallel computation. Complicated environments are created that are at high risks of attacks than their

counterparts of repositories that are centrally managed and monolithic. 3. Controlling Data Access: big data only provides access control at schema level. There is no finer granularity in addressing proposed users in terms of roles and access related scenarios. 4. Node-to-node communication: Hadoop don't implement secure communication; they use the RPC (Remote Procedure Call) over TCP/IP. 5. Client Interaction: Communication of client takes place with resource manager, data nodes. Clients that have been compromised tend to propagate malicious data or links to either service. 6. Virtually no security: big data stacks where designed with no security in mind. There is no security for common web threats too.

## VI. SOLUTION FOR BIG DATA SECURITY

### IN HADOOP

Analyzing the security issues associated with big data Hadoop. Here in this paper I have

mentioned the solutions that help in ensuring the security of data. [18]

#### A. AUTHENTICATION

Hadoop have Kerberos as a primary authentication. Initially SASL/GSSAPI was used for implementing Kerberos and mutually authenticates users, their applications, and other Hadoop services over the RPC connections [7]. Hadoop supports “Pluggable” Authentication for HTTP Web Consoles, meaning implementers of web applications and web consoles can implement their own authentication mechanism for HTTP connections. HDFS communications that is between the Name Node and Data

Nodes is over RPC connection and mutual Kerberos authentication is performed between them [15]. HBase supports SASL Kerberos secure client authentication via RPC, HTTP. Delegation token which is a two party authentication protocol used between user and the Name Node for authenticating users, it is very simple and more effective than three party protocol used by Kerberos [7, 15]. Oozie and HDFS, MapReduce supports delegation token

#### B. ACLs AND AUTHORIZATION

In Hadoop, the access controls is implemented using file-based permissions following the UNIX permissions model. In HDFS, Access control to files could be enforced by the Name Node through file permissions and ACLs of users and groups. MapReduce gives ACLs for job queues;

defining which users or groups can submit jobs to a queue or change queue properties. Hadoop gives fine-grained authorization using file permissions in HDFS and resource level access control using ACLs for MapReduce and coarser grained access control at a service level

#### C. ENCRYPTION

The data needs protection during the transfer to and from the Hadoop system. The simple authentication and security layer (SASL) authentication framework is used to encrypt the data in motion in Hadoop ecosystem. SASL security provide guarantee of the data being

exchanged between client and servers and ensures that, the data is not readable by “man-in middle”. [15]. Hadoop also supports encryption capability to various channels like RPC, HTTP, and Data Transfer Protocol for data in motion

#### D. AUDIT TRAILS

To meet the security compliance requirements, auditing the entire Hadoop ecosystem on a periodic basis and deploy or implement a system that does log monitoring is necessary. HDFS and MapReduce have basic audit support. Apache Hive metastore does

maintain audit information for Hive interactions [13, 15]. Apache Oozie, the workflow engine, provides audit trail for services, Hue Supports audit logs. For that Hadoop component which does not have built-in audit logging, audit logs monitoring tools can be used.

#### VII. ZETTASET ORCHESTRATOR

##### A. PERIMETER SECURITY FAILS

Vendors of data security think that traditional perimeter security solutions such as intrusion detection/prevention technologies and firewalls can properly address the Hadoop and distributed cluster security. But all security solutions that rely on perimeter security fail to

provide effective security to the Hadoop cluster. Firewalls that attempt to map IP to actual AD credentials are problematic in Hadoop environment. Firewalls only restrict access on basis of IP/ports and do not know anything about the Hadoop file system.

##### B. MOVE SECURITY CLOSER TO

##### THE DATA

Zettaset Orchestrator gives a security solution for big data embedded in the data cluster itself which moves security as close to the data as possible, and protection that perimeter security devices such as firewalls fails to deliver. Orchestrator address the security gaps that open-source solutions ignore, with big data management solution that is hardened to address policy, access control, compliance, and risk management in the Hadoop

Zettaset Orchestrator solution is specifically designed to meet the security requirements of the distributed architectures that predominate in big data and Hadoop environments. Orchestrator creates

cluster environment. Orchestrator takes into account RBAC that strengthens user authentication process. Orchestrator makes simple the integration of Hadoop clusters into an existing security policy framework, and support for AD, LDAP. For organizations with compliance reporting requirements, Orchestrator provides extensive logging, search, and auditing capabilities.

security wrapper around Hadoop distribution and distributed computing environment which make it enterprise-ready. With Orchestrator, organizations deploy Hadoop in data center environments

## VIII. CONCLUSION

Today the size of data is increasing rapidly, in this generation of big data where source of data is not fixed there is a need to secure data coming from various sources. As Hadoop is used to process such data, in this paper I have made a study of various security issues associated

with big data in Hadoop environment and the possible solutions implemented. I think that if we focus on the data that are stored and processed, so that personal privacy is not lost, security can effectively work.

## REFERENCES

- [1] Cloud Security Alliance “Top Ten big Data Security and Privacy Challenges”.
- [2] Tom White O'Reilly |Yahoo! Press “Hadoop The definitive guide”.
- [3] Owen O'Malley, Kan Zhang, Sanjay Radia, Ram Marti, and Christopher Harrell “Hadoop Security Design”.
- [4] Mike Ferguson “Enterprise Information Protection - The Impact of Big Data”.
- [5] Volumetric “Securing Big Data: Security Recommendations for Hadoop and NoSQL Environments”, October 12, 2012.
- [6] Zettaset “The Big Data Security Gap: Protecting the Hadoop Cluster”.
- [7] Devaraj Das, Owen O'Malley, Sanjay Radia, and Kan Zhang “Adding Security to Apache Hadoop”.
- [8] Seref SAGIROGLU and Duygu SINANC “Big Data: A Review Collaboration Technologies and Systems (CTS)”, 2013 International Conference, May 2013.
- [9] Horton works “Technical Preview for Apache Knox Gateway”.
- [10] Kevin T. Smith “Big Data Security: The Evolution of Hadoop's Security Model”.
- [11] M. Tim Jones “Hadoop Security and Sentry”.
- [12] Victor L. Voydock and Stephen T. Kent “Security mechanisms in high-level network protocols”. ACM Comput. Surv.1983.
- [13] Vinay Shukla s “Hadoop Security: Today and Tomorrow”.
- [14] MahadevSatyanarayanan “Integrating security in a large distributed system”. ACM Trans. Comput. Syst. 1989.

- [15] Sudheesh Narayana, Packt Publishing “Securing Hadoop- Implement robust end-to-end security for your Hadoop ecosystem”.
- [16] S. Singh and N. Singh, "Big Data Analytics", 2012 International Conference on Communication, Information & Computing Technology Mumbai India, IEEE, October 2011.
- [17] Jeffhurlblog.com “three-vs.-of-big-data-as-applied-conferences”, July 7, 2012.
- [18] Priya P. Sharma, Chandrakant P. Navdeti “Securing Big Data Hadoop: A Review of Security Issues, Threats and Solution” (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 5 (2), 2014.

# Mobile Agent based Security using Big Data Management

A. S. Gousia Banu, Research Scholar, Department of Computer Science

D. Saritha, Research Scholar of JNTU Kakinada, Department of Computer Science

**Abstract—** *This paper deals with information security and safety issues in public open spaces. Public open spaces include high streets, street markets, shopping centers, community gardens, parks, and playgrounds, each of which plays a vital role in the social, cultural and economic life of a community. Those outdoor public places are mashed up with various ICT (Information and communication technologies) tools, such as video surveillance, smart phone apps, and biometric big data (called Cyber Parks). Security and safety in public places may include video surveillance of movement and the securing of personalized information and location-based services. The article introduces technologies used in Cyber Parks to achieve information security in big data era.*

**Keywords—** *Big Data, Cyber Security, Information Security.*

## 1. INTRODUCTION

The data volume used in Internet technologies is rising rapidly. This huge amount is known as big data and is characterized by three aspects according to Madden:

- The data are numerous,
- The data cannot be categorized into regular relational databases,
- The data are generated, captured, and evolved very quickly.

Big data has generated significant interest in various fields, along with the manufacturing of healthcare machines, banking transactions, social media, and satellite imaging. Big data challenges have been described by Michael and Miller, such as rapid data growth, transfer speeds, the diversity of data, and security issues. Big data is still in its infancy stage and has not been analyzed in general. Hence, this survey comprehensively surveys and classifies its various attributes, i.e. volume, management, analysis, security, nature, definitions, and rapid growth rate. The development of new IT technologies has rapidly increased the volume of information, which cannot be evolved using existing technologies and methods.

In computational sciences, big data presents critical problems that require serious attention. In the IT industry as a whole, the rapid growth of big data has generated new challenges with respect to data management and analysis.

According to Khan et al., five common issues involve: volume, variety, velocity, value, and complexity. Madden note additional issues such as the fast growth of volume, variety, value, management, security, and efficiency. In some fields, data have grown rapidly. However, the type of data that increases most rapidly is unstructured data. This type is characterized by “user information” such as high-definition videos, movies, photos, scientific simulations, financial transactions, phone records, genomic datasets, seismic images, geospatial maps, e-mails, tweets, website data, call-center conversations, mobile phone calls, documents, sensor data, telemetry information, medical records and images, climatology and weather records, log files, and text.

## 2. RELATED WORK

According to Khan et al., unstructured information may account for more than 70% to 80% of all data in organizations. Currently, 84% of IT managers evolve unstructured data, and this percentage is expected to drop by 44% in the near future. Most unstructured data are not modeled, are random, and are difficult to analyze.

Big data technology main objective to minimize hardware and evolving costs and to verify the value of information before committing significant company resources. Properly managed big data are accessible, reliable, secure, and manageable. Hence, such applications can be applied in various complex scientific disciplines (either single or

interdisciplinary), along with atmospheric science, astronomy, medicine, biology, genomics, and biogeochemistry.

Big Data analytics are increasingly used to discover potentially interesting patterns in large data sets. In this chapter, we discuss the potential of combining Big Data methods with those of agent-based simulations to support architectural and urban designs, for agent-based models allow for the generation of novel datasets to survey hypothetical situations and thus designs. Specifically, we present two conceptual studies that investigate the utility of agent-based models in conjunction with Big Data analytics in the context of multi-level pedestrian areas and current office designs, respectively. The analyses of the case studies suggest that it will be worthwhile, both for urban designers and architects, to pursue a combined agent-based simulation Big Data analytics approach.

In distributed systems and in open systems such as the Internet, often mobile code has to run on unknown and potentially hostile hosts. Mobile code, such as a mobile agent is vulnerable when executing on remote hosts. The mobile agent may be subjected to various attacks such as tampering, inspection, and replay attack by a malicious host. Much research has been done to provide solutions for various security problems, such as authentication of mobile agent and hosts, integrity and confidentiality of the data carried by the mobile agent. Many of such recommended solutions in literature are not suitable for open systems whereby the mobile code arrives and executes on a host which is not known and trusted by the mobile agent owner. In this paper, we propose the adoption of the reference monitor by hosts in an open system for providing trust and security for mobile code execution. A secure protocol for the distribution of the reference monitor entity is described as well as a novel approach to assess the authenticity and integrity of the reference monitor running on the destination agent platform before any mobile agent migrates to that destination. This reference monitor entity on the remote host may provide several security services such as authentication, Integrity and confidentiality of the agent's code and/or data.

In cloud computing infrastructure, usually cloud user has to rely on cloud service provider for transfer of data into it. It is still a matter of great concern for a cloud user to trust security and reliability of cloud services. There is major need of bringing reliability, transparency and security in cloud model for client satisfaction. The cloud user data resides on virtual machines which are located on a shared environment which makes it vulnerable to

many attacks. In this paper we propose a trust model for cloud architecture which uses mobile agent as security agents to acquire useful information from the virtual machine which the user and service provider can utilize to keep track of privacy of their data and virtual machines. These agents monitor virtual machine integrity and authenticity. Security agents can dynamically move in the network, replicate itself according to requirement and perform the assigned tasks like accounting and monitoring of virtual machines.

### 3. TECHNICAL CONCEPTS

Khan et al. have suggested a new data life cycle that uses the technologies and terminologies of big data. This new approach to data management and handling required in science is reflected in the scientific data life cycle management (SDLM) model. With this model, existing practices are analyzed in different scientific communities. The generic life cycle of scientific data is composed of sequential stages, including experiment planning (for research projects), data collection and evolving, discussion, feedback, and archiving. The suggested data life cycle having the following stages: collection, filtering and classification, data analysis, storing, sharing and publishing, data retrieval and discovery. In evolving big data, users face several challenges.

Applications requires a huge storage capacity, rapidly search engines, sharing and analysis capabilities, and in some areas data visualization. These and others challenges need to avoid to maximize big data. Currently, various techniques and technologies are used, such as SAS, R, machine learning platforms and Mat lab to handle extensive data analysis. However, the schemes are limited in managing big data effectively and are still lacking.

According to Khan et al. , others challenges to big data analysis include data inconsistency and incompleteness, scalability, timeliness, and security. This paper introduces a new scheme for big data management based on agent oriented cyber security in public spaces.

### 4. BIG DATA GENERATION

In various area devices with enormous of sensors networks are used in different fields, such as, security and privacy, social network, transportation, medical care, industry, traffic, and public department. Devices are grown up quickly and collect the most important part of big data. An important source of big data.

*4.1. Security and Privacy Indoor:* Video surveillance system is the most important issue in homeland security field because of its ability to track and to detect a particular person. To avoid the lack of the conventional video surveillance system that is based on user perception this paper introduces a novel cognitive video surveillance system (CVS) that is based on mobile agents. CVS offers important attributes such as suspect objects detection, smart camera cooperation for person tracking.

According to many studies, an agent-based approach is appropriate for distributed systems, since mobile agents can transfer copies of themselves to other servers in the system. Various numbers of papers in the literature have been suggested and focused on computer vision problems in the context of multi-camera surveillance systems. The main problems highlighted in these papers are object detection and tracking and site-wide, multi-target, multi-camera tracking. The importance of accurate detection and tracking is obvious, since the extracted tracking information can be directly used for site activity/event detection. Furthermore, tracking data is needed as a first step toward controlling a set of security cameras to acquire high-quality images, and toward, for example, building biometric signatures of the tracked targets automatically.

The security camera is controlled to track and capture one target at a time, with the next target chosen as the nearest one to the current target. These heuristics-based algorithms provide a simple way of computing. Here the scenario is considered that the smart camera captures two similar objects (e.g. twins), then each object selects different path. The tracking evolve will become confused. Furthermore, the smart camera is limited to cover certain zone in public place (indoor). The suggested solutions to improve the conventional video surveillance system are extended in various ways. A part of the approaches was to use an active camera to track a person automatically, thus the security camera moves in a synchronized motion along with the projected movement of the targeted person.

These approaches are capable of locating and tracking small number of people. Another common approach was to position the camera at strategic surveillance locations. This is not possible in some situations due to the number of cameras that would be necessary for full coverage, and in such cases, this approach is not feasible due to limited resources. A third approach is to identify and track numerous targeted people at the same time involves image evolving and installation of video cameras at

any designated location. Such image evolving increases server load.

The limitation of user perception system in conventional video surveillance system increases the demand to develop cognitive surveillance application. Many of the suggested video surveillance system are expensive and lack the capability of cognitive monitoring system (such as no image analysis) and ability to send warning signal autonomous in real-time and before the incidents happen.

Furthermore, it is difficult and might take a long time for the user to locate the suspects in the video after the incidents did happen. The problem may get more completely in the larger scale surveillance system. The next generation video surveillance systems expected not only to solve the issues of detection and tracking but also to solve the issue of user body analysis. In the literature, it can be found many references in development.

In such area, the CVS main objective to offer meaningful characteristics like automatic, autonomy, real-time surveillance such as face recognition, suspects object, target detection, and tracking using cooperative smart cameras. Many face recognition systems have a video sequence as the input. Those systems may require being capable of not only detecting but also tracking faces. Face tracking is essentially a motion estimation problem. Face tracking can be performed using many different methods, e.g., head tracking, feature tracking, image-based tracking, model-based tracking. These are different ways to classify these algorithms.

*4.2. Model of CVS System:* In this section we introduce the system model of the video surveillance system. Video surveillance system has been used for monitoring, real-time image capturing, evolving, and surveillance information analyzing. The infrastructure of the system model is divided in three main layers: mobile agents that are used to track suspect objects, cognitive video surveillance management (CVS), and protocol for communication as shown in Fig. 1. Each end device, smart camera, covers a certain zone or cell. Smart camera used for collecting parameters of user face.

In the system model has been introduced two communication protocols. The first protocol is used for agent-to-agent communication protocol. The protocol is based on messages exchange as shown in Fig. 1. The goal is to update the agents. The second protocol is used for communication between CVS and mobile agent protocol. Mobile agents are placed

in smart camera stations and main objective to track the suspect object from smart camera station to others.

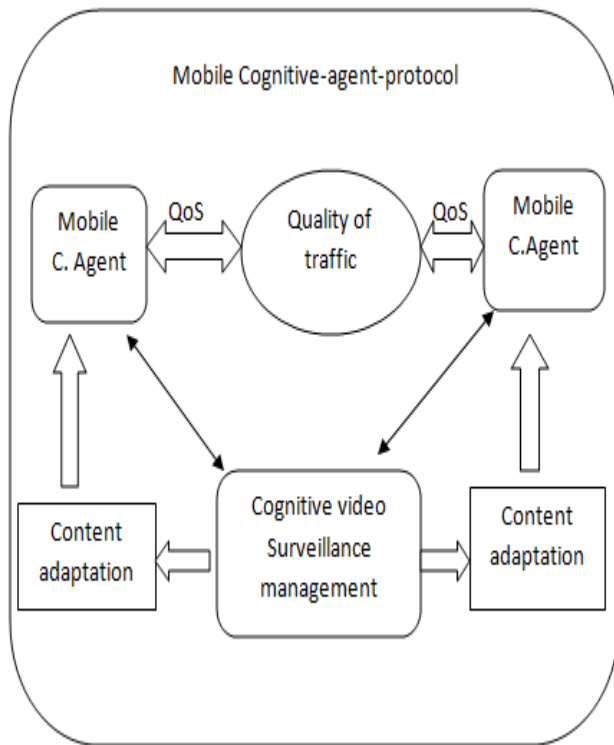


Fig:1 CVS system model concept

Mobile agent offers various characteristics, e.g. negotiation, making decision, roaming, and cloning. CVS provide the mobile agent with information. Based on received information mobile agents make decision when and where to move to next smart camera station. In order to track moving objects, two strategies are used. The first is based on messaging protocol (msg protocol) informing the mobile agent about the position of the suspect objects. The second strategy uses the protocol to help the mobile agent to roaming from point to others.

#### 4.3. Security and Privacy:

Outdoor Modern cities offer various kinds of public places, and are created for different targets, i.e. public places for students and others on academic campuses, for visitors to historical sites, and for families and tourists. Public open spaces that are supported by various kinds of modern information communication technologies are called Cyber Parks. Such places providing connectivity services to users on their personal computers, smart phones, tablets, and other mobile end-devices. Many users use

Internet technologies for storing private data. Furthermore, Internet technologies are used for communication in business, the military, medicine, education, and government and public services. Over the last decade, as well, crime in virtual life has increased.

Cyber attacks are performed through Internet networks that target individual machines, mobile devices, communications protocols, or smart phone application services. Cyber attacks are performed by spreading malware, by creating phishing websites, and by other means. To implement information security policies and safety in Cyber Parks, security models are needed that lay out guidelines for securing information and communication. Cyber Park security models are based on formal models of access rights to smart phone applications and web services.

In addition, an adaptive agent recognizes the applications that being used, and a mobile agent platform creates mobile agents to serve the Cyber Park visitors. By monitoring the behavior of users, detection systems ensure information privacy. A mobile agent main objective to fulfill user's preferences based on a dynamic environment. The mobile agent's structure is divided to three parts, as follows:

- *Source code* – the program resides of several classes to define the agent's behavior. In the source code, the backbone of the agent is created, which contains the basic rules. The agent then grows and develops itself according to the requirements of its environment;
- *State* – the agent's internal variables enable it to resume it is activities when it is found to be in one of the following states: offline (sleeping, in an evolution evolve), online (awake), busy, waiting (standby), or dead;
- *Attributes* – attributes reside of information describing the agent, its movement history, its resource requirements, and authentication keys. In order to mediate useful tasks, a communication model to establish communication between mobile end users and the Cyber Park service provider is used. The agents in the system should be able to understand each other, and they should use the same message transport protocol. Messages are a data oriented communication mechanism, generally used to transfer data between evolves. Communication is either asynchronous or synchronous.

#### 4.4. Concept of Secured Information

Authentication refers to evolve of obtaining a confirmation that a person who is requesting a service, is a valid user. It is accomplished via the presentation of an identity and credentials, such as passwords, tokens, digital certificates, and phone numbers. To increase information security, users need a password to log in. The system starts the identification evolve and creates a mobile agent for each user, as shown in Fig. 2. The mobile agent is responsible for communication security in the system.

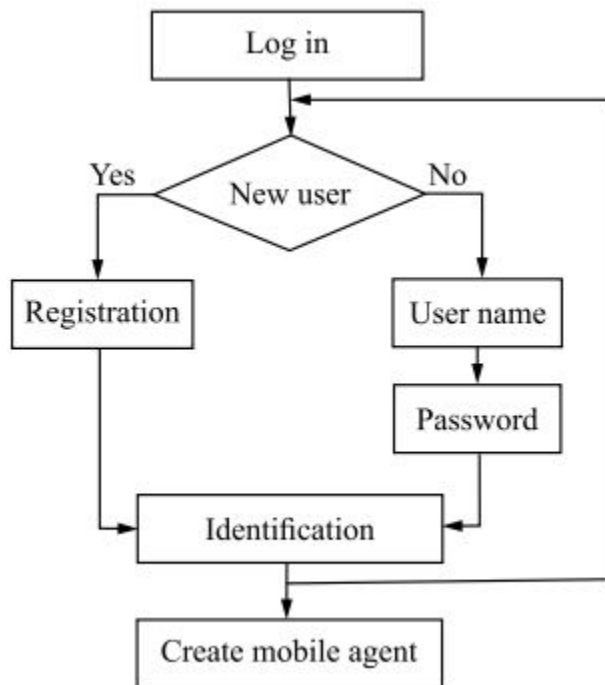


Fig 2:- Block diagram of authentication evolve

Messages are a data oriented communication mechanism. Request-response mechanism is used to transfer data between a user end device and a service provider:

- *Inform Message* – includes the mobile device ID and the kind of information requested,
- *Re-Inform Message* – includes information about Cyber Park resources,
- *Request Message* – includes the sender's name, a time stamp that indicates the time the request message was generated, the receiver's name and the requested resource,

- *Response Message* – includes the sender's name, a time stamp, and the requested resource.

#### 5. SIMULATION/GRAPH ANALYSIS

As organizations adopt social media and platforms and the digital footprint of its customers increases, the amount of data that is available for organizations to analyze and use increases exponentially.

Over 69% respondents employ big data analytics to model for and identify information security threats.

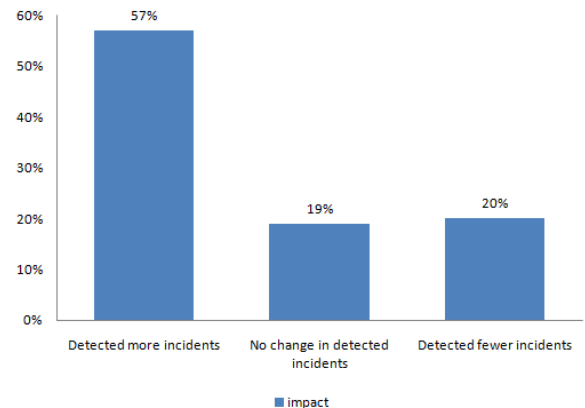


Fig:- Impact of big data analytics on information security

#### 6. CONCLUSION

The Internet and mobile technology are growing rapidly, and the data accumulated over twenty years have become big data. We have considered security big data indoor and outdoor, which is generated by devices. The privacy approach main objective to protect private position data. The mobile agent works to hide the identity of the user and his or her activity in Cyber Park (outdoor) services while the location for the user is visible. This prevents a cyber attacker from detecting the user's location.

#### REFERENCES

- [1] D. Che, M. Safran, and Z. Peng, "From big data to big data mining: challenges, issues, and opportunities", in Database Systems for Advanced Applications, Berlin: Springer, 2013, pp. 1–15.
- [2] S. Madden, "From databases to big data", IEEE Internet Comput., vol. 16, no. 3, pp. 4–6, 2012 (doi: 10.1109/MIC.2012.50).

- [3] K. Michael and K. W. Miller, "Big data: new opportunities and new challenges", *Computer*, vol. 46, no. 6, pp. 22–24, 2013 (doi: 10.1109/MC.2013.196).
- [4] J. Raiyn, "Using cognitive radio scheme for big data traffic management in cellular systems", *Int. J. of Inform. Technol. & Manag.*, vol. 14, no. 2-3, 2015.
- [5] J. Raiyn, "Toward developing real-time online course based interactive technology tools", *Adv. in Internet of Things*, vol. 4, no. 3, pp. 13–19, 2014 (doi: 10.4236/ait.2014.43003).
- [6] J. Raiyn, "A Survey of cyber attack detection strategies", *Int. J. of Secur. & Its Appl.*, vol. 8, no. 1, pp. 247–256, 2014.
- [7] C. A. Steed et al., "Big data visual analytics for exploratory earth system simulation analysis", *Computers & Geosciences*, vol. 61, pp. 71–82, 2013 (doi: 10.1016/j.cageo.2013.07.025)
- [8] N. Khan, I. Yaqoob, I. A. T. Hashem, Z. Inayat, W. K. M. Ali, M. Shiraz, and A. Gani, "Big data: Survey, technologies, opportunities, and challenges", *Scientif. World J.*, pp. 1–18, 2014 (doi: 10.1155/2014/712826).
- [9] J. Raiyn, "Information security and safety in Cyberpark", *Global J. of Adv. Engin.*, vol. 2, no. 8, pp. 73–78, 2015.
- [10] J. Raiyn, "Modern information and communication technology and their application in Cyberpark", *J. of Multidiscip. Sci. & Technol.*, vol. 2, no. 8, pp. 2178–2183, 2015.
- [11] M. Wooldridge and N. R. Jennings, "Intelligent agents: Theory and practice", *The Knowl. Engin. Rev.*, vol. 10, no. 2, pp. 115–152, 1995.
- [12] S. Russel and P. Norvig, *Artificial Intelligence: A Modern Approach*, Englewood Cliffs, NJ: Prentice Hall, 1995.
- [13] M. Luck, V. Marik, O. Stepankova and R. Trappl, Eds., *MultiAgent Systems and Applications. 9th ECCAI Advanced Course ACAI 2001 and Agent Link's 3rd European Agent Systems Summer School, EASSS 2001*, Prague, Czech Republic, July 2-13,

2001. *Selected Tutorial Papers*, LNAI, vol. 2086. Springer, 2001.

- [14] T. Springer, T. Ziegert, and A. Schill, "Mobile agents as an enabling technology for mobile computing applications", *Kuenstliche Intelligenz*, vol. 14, no. 4, pp. 55–61, 2000.

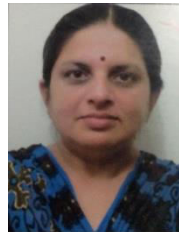
[15] Z. Lin and K. Carley, "Proactive or reactive: An analysis of the effect of agent style on organization decision-making performance", *Intell. System in Account. Finance and Manag.*, vol. 2, no. 4, pp. 271–287, 1993.

#### First A. Author



A. S. Gousia Banu  
Research Scholar  
Department of Computer  
Science  
[gbanuzia@gmail.com](mailto:gbanuzia@gmail.com)

#### Second B. Author



D. Saritha  
Research Scholar of JNTU  
Kakinada  
Department of Computer  
Science  
[Ziauddinb17@gmail.com](mailto:Ziauddinb17@gmail.com)

# GOOGLE PROJECT LOON

**V.Naveena,**  
vemula.naveena1995@gmail.com  
Malla Reddy College of Engineering

**Dr.T.Sunil,** *Prof. & Dean*

Malla Reddy College Engineering,  
sunil.tekale2010@gmail.com

## Abstract

This paper describes an overview of a Balloon-powered Internet for everyone. Currently we are using the internet service through Internet Service Providers to connect globally. Loon purpose is to provide wireless network to remote areas through of a set of high altitude balloon equipped with advanced sophisticated wireless transceivers to connect people globally. This technology could allow developing countries to avoid the using of expensive underground infrastructure. This project loon helpful to connect many areas and will also help to share new ideas and techniques for the development of countries.

## INTRODUCTION

### What is Project Loon?

In the evolution of the Internet nowadays, some population of the world enjoys the benefits of the Internet. According to Google™, two-thirds of people on the earth, reliable Internet connection is still out of reach. To solve this global problem, Google™ developed an innovative project called the “LOON” , to provide broadband for free in rural and remote areas, as well as to improve communication during and after natural disasters or a humanitarian crisis. During

because information in itself is really lifesaving. Here the key concept is a set of a crisis, connectivity is really significant high-altitude balloons ascends to the stratosphere and creates an aerial wireless network (see Fig. 1). The technology designed in the project could allow countries to avoid using expensive underground infrastructure



**Figure 1.** Balloon-based network [1].

### History

The idea may sound a bit crazy, initially everyone thought it as a prank by Google™. The Project Loon unofficial development began in 2011 and officially announced as a Google™ project in 2013. A pilot experiment was happened in New Zealand’s South Island where about 30 balloons were launched (see Fig. 2). [1] After this initial trial, Google™ plans on sending up 300 balloons around the world at the 40th parallel south that could provide coverage to New Zealand, Australia, Chile, and Argentina. Google™ hopes to eventually have thousands of balloon fly in the stratosphere.



**Figure 2.** Balloon launching in New Zealand [1].

### How Loon moves?

Project Loon balloons positioned in the stratosphere winds at an altitude of about 20 km, twice as high as airplane flights and the weather changes (see Fig. 3). In the stratosphere, there are many layers of wind, and each layer of wind varies in direction and speed. Why the stratosphere means? It is situated on the edge of space, Loon balloons are also unique in that they are steerable and entirely solar powered. The balloons and equipment can be reused, and each loon has an approximately 2 years of life time.

In loon design there are

between 10 km and 60 km in altitude having steady winds below 20 mph. This spherical layer is great for solar panels because there are no clouds to block the sun. Loon balloons are directed by rising or descending into a layer of wind blowing in the desired direction of travel by using wind data from NOAA. [9] By moving with the wind, the balloons can be arranged to form one large communication network. Each balloon is equipped with a GPS for tracking its location. Project Loon has complex algorithms to determine where its balloons need to go, then moves each one into a layer of wind blowing in the right direction.



Loon moves [1].

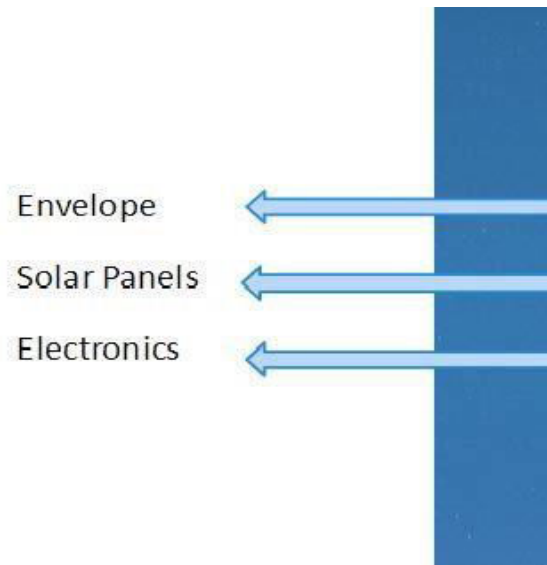
### Loon design

**Figure 3.** How

three main components (see Fig. 4):

1. Envelope
2. Solar Panels

### 3. Electronics



**Figure 4.** Loon Design [1].

#### Envelope:

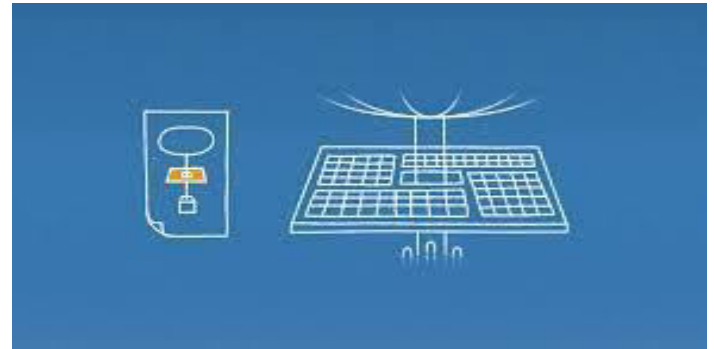
The inflatable part of the balloon is called envelope. Each super-pressure balloon is made of polyethylene plastic material and filled with helium. When fully inflated, the balloon height is 12 m and its width is 15 m. The envelope is designed to resistant exposure to UV rays and is capable to function at dramatic temperature swings as low as -80oC. A well-made polyethylene plastic balloon envelope is critical for allowing a balloon to last around 100 days in the stratosphere. A parachute is attached to the top of the envelope, which is used for bringing down the balloon safely.



**Figure 5:** Envelop

#### Solar panels:

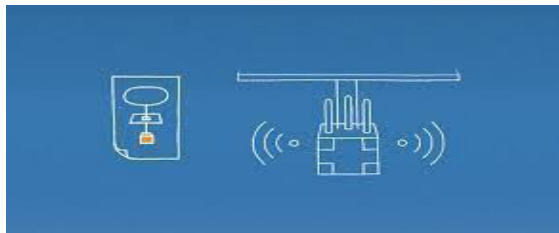
Each balloon's solar panel provides power to its own electronics. The solar array is made of flexible plastic laminate supported by a light-weight aluminum frame. It uses high efficiency monocrystalline solar cells. The solar panels are mounted at a steep angle to effectively capture sunlight. The panels produce approximately 100 Watts of power in full sun (that power is sufficient to keep Loon's electronics running 24 hrs a day), and the additional power is stored in a rechargeable battery.



**Figure:** Solar Panels

#### Equipment:

A small electronics box (payload) hangs underneath the inflated envelope. This box contains circuit boards, Linux-based computer, radio antenna, GPS, sensors, and batteries. They have specific functions [6]: circuit boards to control the system, radio antenna for communication, GPS for tracking location, sensors to monitor and record weather conditions, and lithium ion batteries to store solar power. [1]



**Figure:**Equipment

### How Loon connects?

Each balloon has a radio antenna that provides constant connectivity to the ground and connects each balloon to other balloon. There is a special ground antenna that is installed on the home or workplace to access the internet from balloon. Google™ claims that each balloon can provide signal connectivity to a ground area about 40 km in diameter and able to deliver 3G comparable speeds (up to 10 Mbps). [1] These antennas use ISM bands of spectrum 2.4 GHz & 5.8 GHz. ISM radio bands (portions of the radio spectrum) reserved internationally for industrial, scientific, and medical purposes other than telecommunications. [4][7]



**Figure 5.**

Transmitting signals [5].

Google™ balloons are connected in the mesh topology to ensure reliability. The IEEE802.11s standard defines how wireless devices form the mesh network. Loon's protocol stack is not disclosed yet. [6]

There are two types of communications (see Fig. 5):

1. Balloon-to-balloon communication
2. Balloon-to-ground communication.

### Subscriber-to-ISP:

First, the specialized internet antenna (see Fig. 6) on the ground sends signals to a balloon. Then signal hops forward from the balloon to neighboring balloons. Finally, signals from the balloon reach a ground station which is connected to a local internet provider, or pre-existing internet infrastructure which provides service via the network of balloons.

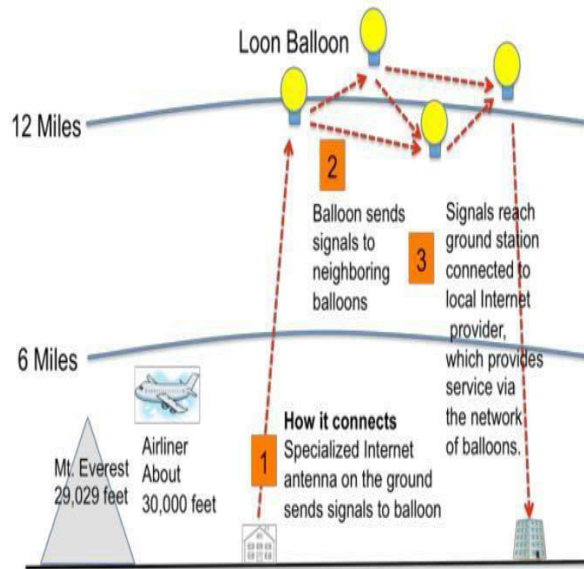
**Figure 6.**



**Figure 6.** Ground antenna [6].

### ISP-to-Subscriber:

The Internet Service Provider or pre-existing internet infrastructure sends response back to the Balloon network; then data travels through the balloon network. Finally, the closest balloon to the subscriber receives data and sends it back to the subscriber (see Fig. 7). [2][5]



**Figure 7.**

How it connects in Project Loon [5].

## Maintenance

If a balloon fails or needs maintenance, Google™ staff brings the balloon down. A trigger mechanism on the top of the balloon would deflate it by releasing gas from the envelope, and it releases a parachute that brings the balloon down to the Earth in a controlled descent. GPS equipment tracks where the balloon is landing. Google™ needs the dedicated staff across the globe for balloon maintenance. [1]

## Challenges

Google™ wants to build a network with no borders. Its biggest obstacle is not technology. Some countries unwilling to give permission. In addition to permissions, Google™ should negotiate with countries to purchase or borrow specific radio frequencies. There might be spying and security threat over data. [8]

*Pros:*

The most obvious avails of the project is that Google™ will provide the Internet for free. This may increase the Internet usage throughout the world. Ground antennas are easy to use and install. No extra underground infrastructure is required; the equipment is relatively cheap.

*Cons:*

This project is labor intensive and provides the limited internet speed. Balloons can work 100 days only. The main problem is that the hardware failures cannot be reached at the intended location. If a Loon balloon fails, it can either remain up in the air floating or it might go down in unwanted areas like sea. These scenarios are a huge concern to the stability as well as the safety of people. [1]

## Competing Ideas

Facebook™'s Drones is the competing idea for Google™ Loon. As compare to the balloons, drones provide more coverage area per drone, more internet speed, and can stay up in the air for years (~5 yrs.). But it requires expensive equipment and could do a lot of damage if it fails and fell out of the sky, and security and privacy have other concerns accompanied the use of drones. [10]

## Conclusion

Internet connectivity and communication become one of the basic needs in modern human daily life. An innovative and scalable idea like the Google™ Project Loon would aid and benefit remote areas of the world as well as population to reap the benefits of

modern communications. It would also provide backbone communications during and after natural disasters when ground infrastructure is scarce or destroyed.

## Abbreviations

NOAA - National Oceanic and Atmospheric Administration

GPS - Global Positioning System

UV rays - Ultraviolet rays

IEEE - Electronics and Electrical Engineering

ISP - Internet Service Provider

## References

- [1] Loon for all. [Online]. Retrieved from <http://www.google.com/loon/>
- [2] Project Loon. [Online]. Retrieved from <http://ipnsig.org/wp-content/uploads/2014/02/Project-Loon.pdf>
- [3] Jose Saldana. (2014, March 4). Origins of Project Loon. [Online]. Retrieved from <http://www.ietf.org/mail-archive/web/gaia/current/msg00068.html>
- [4] ISM BAND. [Online]. Retrieved from [http://www.princeton.edu/~achaney/tmve/wiki100k/docs/ISM\\_band.html](http://www.princeton.edu/~achaney/tmve/wiki100k/docs/ISM_band.html)
- [5] Ankur Sharma. (2013, October 28). Google Project Loon - What is it? [Online]. Retrieved from <http://www.lifengadget.com/lifengadget/googles-project-loon/>
- [6] Doowon Kim. (2013, December 20). A Survey of Balloon Networking Applications and Technologies. [Online Paper]. Retrieved from

<http://www.cse.wustl.edu/~jain/cse570-13/ftp/balloonn.pdf>

- [7] Project Loon. (2013, June 14). [Online]. Retrieved from

<https://www.youtube.com/user/ProjectLoon>

- [8] Kevin Fitchard. (2013, June 21). Project Loon: Google's biggest obstacle isn't technology. It's politics. [Online]. Retrieved from

<http://gigaom.com/2013/06/21/project-loon-googles-biggest-obstacle-isnt-technology-its-politics/>

- [9] Steven Levy. (2013, June 14). How Google Will Use High-Flying Balloons to Deliver Internet to the Hinterlands. [Online]. Retrieved from

[http://www.wired.com/2013/06/google\\_internet\\_balloons/all/google.com/loon#slideid-175682](http://www.wired.com/2013/06/google_internet_balloons/all/google.com/loon#slideid-175682)

- [10] Ben Popper. (2014, March 7). Google's balloons versus Facebook's drones. [Online]. Retrieved from

<http://www.theverge.com/2014/3/7/5473692/facebook-drone-titan-aerospace-project-loon>

## RESOURCE PLANNER

**P.Priskilla,**  
lovelyprisk514@gmail.com  
Malla Reddy College of Engineering

**Dr.T.Sunil,** Prof. & Dean  
Malla Reddy Engineering College for Women,  
sunil.tekale2010@gmail.com

### **Abstract**

*When a software development company wants to achieve its goals on time and efficiently use its staff on the projects, it is necessary for the company to have hands on information related to number of employees working on various projects along with their skill set and the number of employees still needed to complete the projects on time. Resource Planner is a convenient tool to handle various projects in a software company efficiently.*

*When a software development company wants to achieve its goals on time and efficiently use its staff on the projects, it is necessary for the company to have hands on information related to number of employees working on various projects along with their skill set and the number of employees still needed to complete the projects on time. Resource Planner is a convenient tool to handle various projects in a software company efficiently.*

*Resource Planner is an online tool to manage projects currently running with the company as well as future projects. This tool tracks the employees working for the existing projects and details of new projects like no. of employees required, location, etc. This tool is very useful in estimating revenue, etc which helps higher management to know the status of the various*

*projects and work force. With this tool HR can estimate the requirement of employees for the new projects and hence can recruit exact number of employees.*

*This application maintains the centralized database so that any changes done at a location reflects immediately. This is an online tool so more than one user can login into system and use the tool simultaneously.*

*The administrator of this software will be able to create new users and remove any user. He allots passwords and changes them. He can view the details of all employees in the company. He can also view the management reports where the information is presented project wise and location wise.*

### **Existing System**

*Current system is a manual one in which users are maintaining books etc to store the information like project details, requirement, availability and allocations of employees for the existing project as well as for the new projects. It is very difficult to maintain historical data. Also regular investments need to purchase stationary every year.*

*The following are the disadvantages of current system*

- *It is difficult to maintain important information in books*
- *More manual hours need to generate required reports*

- *It is tedious to manage historical data which needs much space to keep all the previous years books etc*

*Daily transactions are to be entering into different books immediately to avoid conflicts which are very difficult*

## **Proposed System**

*Proposed system is a software application which avoids more manual hours that need to spend in record keeping and generating reports. This application keeps the data in a centralized way which is available to all the users simultaneously. It is very easy to manage historical data in database. No specific training is required for the employees to use this application. They can easily use the tool that decreases manual hours spending for normal things and hence increases the performance. As the data is centralized it is very easy to maintain the currently running projects with the company as well as future projects.*

*The following are the advantages of proposed system*

- ✓ *Easy to manage all the daily transactions*
- ✓ *Can generate required reports easily*

1. **EMPLOYEE MODULE:** *This module deals with major and crucial part which tracks the details of employees currently working with the company. It allows the HR Manager only to add a new employee record into the database and it allows HR User only to easily remove an employee*

- ✓ *Easy to manage historical data in a secure manner*
- ✓ *Centralized database helps in avoiding conflicts*

*Easy to use GUI that does not requires specific training.*

## **Scope of the Project**

*The proposed system's user interface is developed in a browser specific environment to have web based architecture. The web documents are designed using HTML standards and JSP power the dynamic of the page design. The communication architecture is designed by concentrated on the standards of Servlets and JSP. The database connectivity is established using the Java Database connectivity (JDBC). Any changes made to the requirements in the future will have to go through formal change approval process.*

### **Number of Modules**

*The system after careful analysis has been identified to present itself with the following modules:*

*from the database. It allows all types of users to view the list of users current existing in our company. It facilitates us to convert the employee report into excel format just by clicking download to excel button*

2. **EMPLOYEE MODULE:** *This*

*module deals with major and crucial part which tracks the details of employees currently working with the company. It allows the HR Manager only to add a new employee record into the database and it allows HR User only to easily remove an employee from the database. It allows all types of users to view the list of users current existing in our company. It facilitates us to convert the employee report into excel format just by clicking download to excel button*

3. **PROJECTS MODULE:** *This module deals with major and crucial part which maintains the details of projects currently with the company & future projects. It allows the project manager to add new projects details to the database. It provides a user-friendly interface to add new projects. It allows PM to view and remove the details related to a project very easily. It provides an option to convert projects report into excel format*

1. **REQUIREMENT & ALLOCATION MODULE:** *This module deals with major and crucial part which provides Info about project-wise requirements which includes onsite and offshore that was entered by the project managers of different projects. It allows any type of user to view these project requirements. It helps the HR People to view project-wise requirements and start recruiting the people. It also helps in allocating the*

*people to a project after recruitment sothat HR people can idea about the gap between requirement and allocation at any point of time very easily by generating HRD GAP Summary report. It provides all these reports to be converted and stored permanently in excel sheets.*

2. **PROJECTS MODULE:** *This module deals with major and crucial part which maintains the details of projects currently with the company & future projects. It allows the project manager to add new projects details to the database. It provides*

*a user-friendly interface to add new projects. It allows PM to view and remove the details related to a project very easily. It provides an option to convert projects report into excel format*

### 3. **REQUIREMENT & ALLOCATION**

**MODULE:** This module deals with major and crucial part which provides Info about project-wise requirements which includes onsite and offshore that was entered by the project managers of different projects. It allows any type of user to view these project requirements. It helps the HR People to view project-wise requirements and start recruiting the people. It also helps in allocating the people to a project after recruitment so that HR people can idea about the gap between requirement and allocation at any point of time very easily by generating HRD GAP Summary report. It provides all these reports to be converted and stored permanently in excel sheets.

**ADMIN & REPORTING MODULE:** It Provides interfaces to manage this tool like add/remove users, change privileges of users etc. This module used to provide different reports required by the higher management for better analysis. It generates dynamic reports like Role-Location

which displays role-wise employees report in different locations, Project-Location report which displays project-wise employees reports in different locations, Role-Skill reports which displays skill-wise role based employees list in different locations, Project-Skill reports which displays project-wise skills report of different employees and Skill-Location report which displays skill-wise employees report in different locations etc.

#### **Features to be implemented**

- ☐ Session management
- ☐ Connection pooling
- ☐ Normalized database
- ☐ Prevention of duplication login
- ☐ Design patterns

#### **Technologies to be used**

- ☐ Web Presentation: HTML, CSS
- Client – side Scripting: JavaScript
- ☐ Programming Language: Java
- ☐ Web based Technologies: JNDI, Servlets, JSP
- ☐ Database Connectivity API: JDBC
- ☐ Build Tool: ANT
- ☐ Debug Tool: Log 4J
- ☐ CASE tool: Rational Rose, Visual Paradigm, Enterprise Architect

- *Backend Database: Oracle/SQL Server/MY SQL/MS Access*
- *Operating System: Windows XP/2000/2003, LINUX, Solaris*
- *J2EE Web/Application Server: Tomcat/Weblogic/Websphere/JBoss/Glass Fish*
- *IDEs: Eclipse with My Eclipse plug-ins/Net Beans/RAD*

*References:*

- [1] Devarakonda Krishna  
“A Safety Summons on  
Cloud Data”IETE-  
Chandigarh,PP:340-  
343,17-18 29 July.2017
- [2] BIRRELL, A. D., AND  
NELSON, B. J.  
Implementing remote  
procedure calls. ACM  
Trans. Comput. Syst. 2,  
1 (Feb. 1984), 39-59.

# Data Storage Security in a Hosted Environment

Anitha Bejugama<sup>1</sup>, Shravani Reddy<sup>2</sup>

Department of Computer Science and Engineering  
Mallareddy College of Engineering  
Kompally, Hyderabad, Telangana, India.  
[Anithabejugama@gmail.com](mailto:Anithabejugama@gmail.com) , [shravaniannam@gmail.com](mailto:shravaniannam@gmail.com)

## Abstract

The introduction of hosted environments to computing has brought a lot of change in the industry. With the advent of Infrastructure-as-a-Service, the concept of renting data stores is catching up. Security for the data in a hosted environment is an alarming issue. The paper focuses on cloud data storage security. We propose a security model to maintain the data store. This model ensures the safety of data in the hosted environment. By utilizing the homomorphic token with random masking the model achieves the abstraction of data stored on the Hosted Environment.

**Keywords:** Data Store, Security, Hosted Environment, Infrastructure-as-a-Service.

## 1. Introduction

Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction [1]. Cloud computing not only reduces the cost of service delivery but also increases the speed and agility with which services are deployed. Cloud computing incorporates virtualization, on-demand deployment, Internet delivery of services, and open source software [2]. The characteristics of cloud computing include On-demand self-service, Broad network access, Resource pooling, Rapid elasticity and Measured service. Even though cloud offers a lot of advantages there are some issues that cloud has to deal with. Security in cloud computing is a major concern in the industry. Many companies are still waiting

for an amicable solution for this major problem to be solved as they don't want to take any risks with their critical & sensitive data. The unique issues associated with cloud computing security have not been resolved yet. Access to your information from anywhere at any time is the specialty of hosted environments. You don't need to be in the same physical location as the hardware that stores your data. The Cloud provider houses the hardware and software necessary to run your applications.

Cloud computing raises a lot of security threats. Traditional cryptographic primitives can not be directly adopted for data security in a hosted environment because the user loses control over the data. Verification of correct data storage in the cloud must be conducted without explicit knowledge of the whole data. The problem of verifying correctness of data storage in the cloud is a challenge because various kinds of data for each user is stored in the cloud and also there is a demand for long term continuous assurance of this data. Data storage in a Hosted environment is not just a third party data

warehouse. The data stored in the Hosted environment may be frequently updated by the users. The various operations include insertion, deletion, modification, appending, reordering, etc. Ensuring storage correctness under dynamic data update is a significant task. Traditional integrity insurance techniques don't work here therefore new solutions are required. The deployment of data Hosted environment is powered by data centers running in a simultaneous, cooperated and distributed manner. In order to reduce the data integrity threats, data is redundantly stored in multiple physical locations. Therefore, distributed protocols for storage correctness, assurance can be employed for cloud data storage system. But the field is still evolving. Recent research is revolving around the importance of ensuring the remote data integrity .The techniques discussed in [3]-[7], can be useful to ensure the storage correctness without having users possessing data, can not address all the security threats in cloud data storage, since they are all focusing on single server scenario and most of them do not consider dynamic data operations. None of the proposed distributed schemes can be applied to dynamic data scenario. As a result, their applicability in cloud data storage can be drastically limited. In this paper, to achieve the abstraction of data stored on a Hosted Environment, we propose a security model that utilizes the homomorphic token with random masking technique. The rest of the paper is organized as follows: Section 2 throws light on basic architecture of a data store in a hosted environment. Section 3 gives insight of Cloud Data Storage using

Homomorphic Authenticator. Section 4 deals with the results & analysis part. Section 5.6 talk about the future works & Conclusions respectively.

## 2. BASIC ARCHITECTURE

### 2.1. System Model

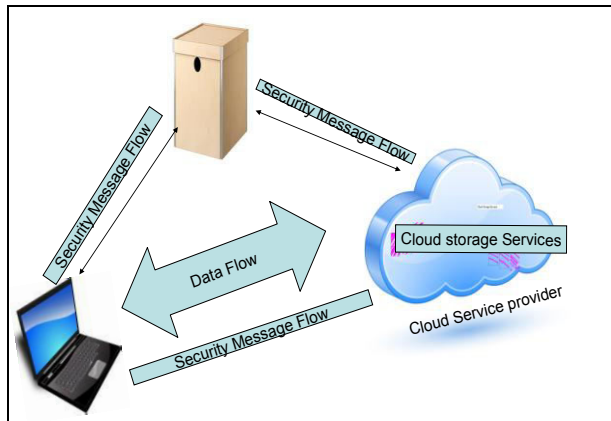
Figure 1.illustrates the representative network architecture for cloud data storage. Three different network entities exist:

**Clients:** clients comprise of users who have data to be stored in the cloud. Clients rely on the cloud for data computation. They include individual consumers and organizations.

**Cloud Service Provider (CSP):** A CSP has resources and expertise in managing and building distributed cloud storage servers. A CSP owns and operates live Cloud Computing systems.

**Third Party Auditor (TPA):** An optional TPA is employed to assess and expose risk of cloud storage services on behalf of the clients. A client stores his data through a CSP into a set of cloud servers, which are running in a simultaneous, cooperated and distributed manner. Faults or server crashes can be tolerated by employing Data redundancy with technique of erasure-correcting code Thereafter, for application purposes, the user interacts with the cloud servers via CSP to access or retrieve his data. The basic operations that a client performs include block update, delete, insert and append. As Client no longer possesses their data locally, it is of critical importance to assure users that their data are being correctly stored and maintained. The clients must be equipped with security means to make continuous correctness assurance of their stored data even if the client doesn't posses a local copy. A TPA on behalf of the client can monitor the data, if the client doesn't posses the resources, time.

Fig. 1: Hosted Environment Data Storage Architecture



### 3. Cloud Data Storage using Homomorphic Authenticator

In cloud data storage system, the data is stored in the cloud and the client no longer has access to the data. The availability & correctness of the data being stored in the Hosted Environment must be assured. Effective detection of any unconstitutional data variation and corruption is the key issue. Hence we propose this model.

The proposed model contains four modules:-

1. **Key generation module:** - In this module we use a key generation algorithm. User runs this module and sets up the scheme.
2. **Signature generation module:** - In this module we use a signature generation algorithm. Verification metadata is generated by the client through this module. Metadata may consist of MAC, signatures or other information used for auditing.
3. **Proof Generation module:** - In this module we use a Proof generation algorithm. Cloud server generates a proof of data storage correctness by running this module.
4. **Proof Verification Module:** - In this module we use a Proof Verification

algorithm. TPA audits the proof from the cloud server by running this module.

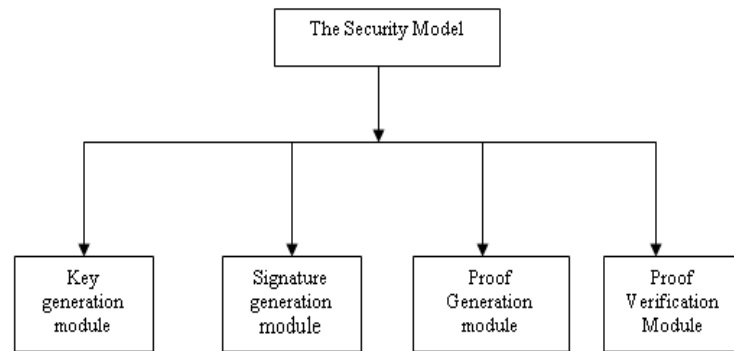


Fig2. The Security Framework

The following are the steps used in the framework:-

1. Public key & secret key generation.
2. File blocks code generation.
3. Blocks Migration in to cloud.
4. Challenge message generation.
5. Cloud Service Provider authentication.
6. Verification.

The first three steps are termed as the setup phase & the last three are termed as the audit phase.

1. **Public key & secret key generation:** - The user generates public and secret parameters. Key generation algorithm is used.
2. **File blocks code generation:-** A code is generated by the user for each file block using homomorphic authenticator. We also use a random mask achieved by a Pseudo Random Function (PRF). By looking only at the aggregated authenticator, a linear combination of data blocks can be checked.

$$\mu' = \sum_{i \in I} \nu_i m_i$$

Where  $\nu_i$  are random number,  $m_i$  are file blocks.

If TPA finds several linear combinations of the similar blocks, it might be able to infer the file blocks. A random mask given by the Pseudo Random Function (PRF) is used.

$$\mu = r + \gamma \mu'$$

Where  $r$  is the mask.

Fig 3. Depicts the generation of file blocks using homomorphic authenticator.

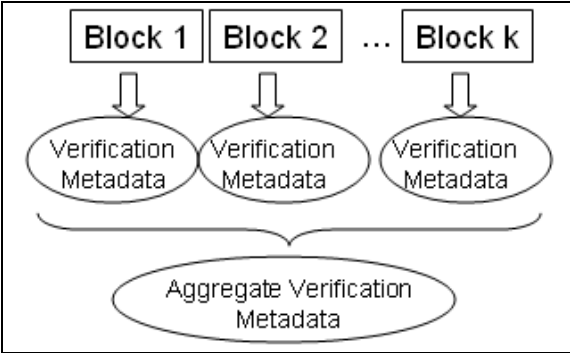


Fig.3 Homomorphic authenticator verifying file blocks

The PRF function masks the data. Verification Metadata is not affected by PRF. The Blocks without PRF mask & with PRF mask are verified. If both of them are equal then only they are authenticated by the aggregate authenticator. Figure 4 depicts this.

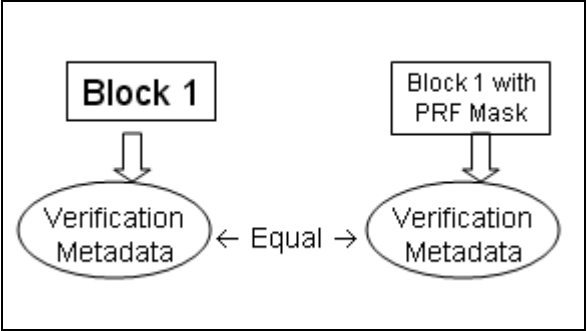


Fig.4 Random Mask by PRF

### 3. Blocks Migration in to cloud: - The codes, file blocks are migrated to the cloud.

Figure 5 depicts the setup phase.

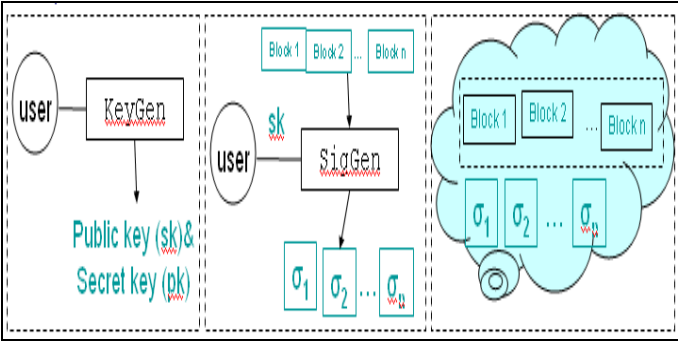


Fig.5 Setup phase

4. Challenge Message Generation: - The third party auditor sends a challenge message to the Cloud Service Provider. The position of the blocks is present in the challenge message. These positions will be checked in this phase.
5. Cloud Service Provider authentication: - A proof generation algorithm is run by the Cloud server. It generates a proof of data storage correctness. The CSP picks the file blocks generated in the challenge, applies the Proof Generation algorithm and generates the proof.

The selected blocks are linearly combined by the CSP. These blocks are applied a mask. Separate PRF key is used for each audit. The CSP sends aggregate authenticator & masked combination of the blocks to TPA for further processing.

6. Verification: - In this step the proof is verified by the Verification algorithm. The Third Party Authenticator generates an aggregate authenticator. It verifies aggregate authenticator & masked combination of blocks received from the Cloud Service Provider by comparing it with the obtained Aggregate authenticator. Figure 6 depicts it.

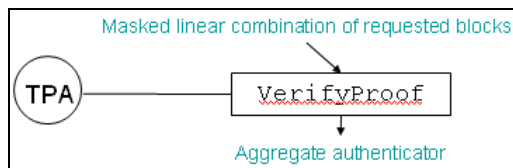


Fig 6 TPA Verification

```

Algorithm keyGen (p,s)
{
    //user generates the public, secret keys.
    Input public parameters p;
    Input private parameters s;
    Generate Public key p(k);
    Generate Secret key s(k);
    End;
}
    
```

```

Algorithm SigGen (m1,m2,m3...mn)
{
    //client generates aggregated authenticator for the
    file blocks.
    Use homomorphic authenticator.
    Use random mask.
    Generate  $\sigma_1, \sigma_2, \sigma_3, \dots, \sigma_n$  for m1,m2,m3...mn
    and aggregate them.
    Send it to CSP with respective file blocks.
}
    
```

```

Algorithm Genproof (TPAchallenge blocks)
{
    //CSP generates aggregated authenticator
    Receive challenge blocks from TPA.
    Combine linearly the blocks.
    Apply mask on the received blocks
    Generate aggregate authenticator.
    Send the masked combination of blocks + aggregate
    authenticator to TPA for verification
    End
}
    
```

```

Algorithm Verifyproof(aggregate authenticator,
masked combination of blocks)
{
    //TPA verifies the received aggregate authenticator
    and generated authenticator.
    Receive CSP computed aggregate authenticator.
    If Received aggregated authenticator == generated
    aggregated authenticator
    Return File block secure Else
    
```

File block tampered.

## 4. Analysis & Results:

| Factors                   | Our Model |        | Public Auditing Model |        |
|---------------------------|-----------|--------|-----------------------|--------|
|                           |           |        |                       |        |
| Selected Blocks           | 380       | 300    | 380                   | 300    |
| Server Computing Time(ms) | 339.52    | 270.20 | 407.66                | 265.87 |
| TPA Computing Time(ms)    | 419.47    | 476.81 | 504.25                | 472.55 |
| Cost per Byte             | 132       | 40     | 132                   | 40     |

Table 1. Performance analysis of the proposed model.

The Table 1 compares our model with the existing Public Auditing Model. We need 300 or 380 blocks to detect that with a probability larger than 95% or 99%, respectively if the server is missing 1% of the data. The data transmitted from CSP to TPA is independent of the data size and the Linear combination with mask.

## 5. Future Work

We can extend our model into a multi-user setting, where the TPA can perform multiple auditing tasks in a batch manner for better efficiency. Even the dynamics of the data on the cloud can be modified so as to adapt to any type of application.

## REFERENCES

- [1] Peter Mell, Timothy and Grance “The NIST Definition of Cloud Computing” <http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf>
- [2] *White Paper Presents General Concepts, Architectural Models & Considerations* June 29, 2009, Volume 137, Issue 1
- [3] A. Juels and J. Burton S. Kaliski, “PORs: Proofs of Retrievability for Large Files,” *Proc. of CCS '07*, pp. 584– 597, 2007.
- [4] H. Shacham and B. Waters, “Compact Proofs of Retrievability,” *Proc. of Asiacrypt '08*, Dec. 2008.
- [5] K. D. Bowers, A. Juels, and A. Oprea, “Proofs of Retrievability: Theory and Implementation,” *Cryptology ePrint Archive, Report 2008/175*, 2008, <http://eprint.iacr.org/>.
- [6] G. Ateniese, R. Burns, R. Curtmola, J. Herring, L. Kissner, Z. Peterson, and D. Song, “Provable Data Possession at Untrusted Stores,” *Proc. Of CCS '07*, pp. 598–609, 2007.
- [7] G. Ateniese, R. D. Pietro, L. V. Mancini, and G. Tsudik, “Scalable and Efficient Provable Data Possession,” *Proc. of SecureComm '08*, pp. 1– 10, 2008.
- [8] T. S. J. Schwarz and E. L. Miller, “Store, Forget, and Check: Using Algebraic Signatures to Check Remotely Administered Storage,” *Proc. of ICDCS '06*, pp. 12–12, 2006.
- [9] M. Armbrust et al., “Above the Clouds: A Berkeley View of Cloud Computing.” 2009.
- [10] A. Ruiz-Alvarez and M. Humphrey, “An Automated Approach to Cloud Storage Service Selection,” in *2nd Workshop on Scientific Cloud Computing (Science Cloud 2011)*, 2011.
- [11] M. Berkelaar, K. Eikland, and P. Notebaert, “Ipsolve: Open source (mixed-integer) linear programming system,” 2011. [Online]. Available: <http://lpsolve.sourceforge.net/>.

## 6. CONCLUSION

In this paper, we investigated the problem of security of data storage in a hosted environment. We have proposed an effective model that ensures the correctness of clients’ data in hosted environment. By utilizing the Homomorphic token with Random masking, our model achieves the abstraction of the data stored on the cloud. Our model eliminates the burden of client from the tedious and possibly expensive auditing task. It alleviates the clients’ fear of their outsourced data leakage. Through detailed security and performance analysis, we show that our model is highly efficient and robust. It prevents any malicious data modification attack, and even server colluding attacks.

- [12] W. W. Chu, "Optimal File Allocation in a Multiple Computer System," *IEEE Transactions on Computers*, vol. 18, no. 10, pp. 885- 889, Oct. 1969.
- [13] R. G. Casey, "Allocation of copies of a file in an information network," *In Proceedings of the AFIPS Joint Computer Conferences*, 1972, pp. 617-625.
- [14] E. Grapa and G. G. Belford, "Some theorems to aid in solving the file allocation problem," *Communications of the ACM*, vol. 20, no. 11, p. 878, 1977.
- [15] K. Lam and C. T. Yu, "An approximation algorithm for a file allocation problem in a hierarchical distributed system," in *In Proceedings of the 1980 ACM SIGMOD International Conference on Management of Data*, 1980, pp. 125 - 132.
- [16] S. Mahmoud and J. S. Riordon, "Optimal allocation of resources in distributed information networks," *ACM Transactions on Database Systems (TODS)*, vol. 1, no. 1, p. 66, 1976.
- [17] L. W. Dowdy and D. V. Foster, "Comparative Models of the File Assignment Problem," *ACM Computing Surveys (CSUR)*, vol. 14, no. 2, p. 287, 1982.

## *Agricultural update via SMS*

*Supriya Chinna,*

[Supriyachinna75@gmail.com](mailto:Supriyachinna75@gmail.com)

Malla Reddy College of Engineering

*Dr.Raja Sekar*

St.Peter College of Engineering,

[rajasekaratr@gmail.com](mailto:rajasekaratr@gmail.com)

### **Abstract**

Agriculture has always been India's most important economic sector. India is one of the fastest growing economies of the world and is currently the focus of a great deal of international attention. In the mid-1990s, it provides approximately one-third of the GDP (gross domestic product) and employs roughly two-thirds of the population. It is the seventh largest country in the world in terms of its geographical size. Agriculture still provides the bulk of wage goods required by the nonagricultural sector as well as numerous raw materials for industry.[7]

The indirect share of agricultural products in total exports, such as cotton textiles and jute goods, is taken into account, the percentage is much higher. With current population growth by 2025 India may even have caught up with China according to the UN. In this paper we focus on agriculture and especially on agriculture trade. India has a large and diverse agriculture and is one of the world's leading producers. [7]

It is also a major consumer, with an expanding population to feed. For this reason and agricultural trade policy, its

presence on the world market has been modest. Given the size of Indian agriculture, changes in its balance sheets for key commodities have a potentially large impact on world markets [8].

Agriculture plays an important, though declining role in the economy. Its share in overall GDP fell from 30% in the early nineties, to below 17.5% in 2006.[8]

Agriculture will continue to play a central role as Asia pursues the complementary goals of poverty reduction, sustainable food security, environmental conservation, and increasing trade competitiveness. According to the surveys new technologies including crop biotechnology, will be essential to meet these challenges. The prospects for their utilization are particularly promising. [8]

Plant biotechnology will facilitate the farming of crops with multiple durable resistances to pests and diseases, particularly in the absence of pesticides. This is expected to be very much useful in the countries like India. There is a lot of work going on this field. Some examples like Golden Rice, BT Brinjal, and BT Cotton etc. can be considered. Now a day's

various organizations, research organization  
Institutes, Universities & Government  
bodies are working on this. [7]

**Keywords: GDP, CTS, IETF**

## REQUIREMENTS ANALYSIS

The requirement gathering process is intensified and focused specifically on application. To understand the nature of the application to be built, the software engineer must understand the information for the software, as well as the required function, performance, and interfacing. Requirements have been documented and reviewed from the user's point of view. [2]

## Tools and Technology

### Android 2.5 Platform

This document enumerates the requirements that must be met in order for mobile phones to be compatible with Android 2.5. The use of "must", "must not", "required", "shall", "shall not", "should", "should not", "recommended", "may" and "optional" is per the IETF standard defined in RFC.

As used in this document, a "device implementer" or "implementer" is a person or

organization developing a hardware/software solution running Android 2.5. A "device implementation" or "implementation" is the hardware / software solution so developed.[1]



To be considered compatible with Android 2.5, device implementations:

- MUST meet the requirements presented in this Compatibility Definition, including any documents incorporated via reference.[2]
- MUST pass the most recent version of the Android Compatibility Test Suite (CTS) available at the time of the device implementation's software is completed. The CTS tests many, but not all, of the components

outlined in this document.[2]

Where this definition or the CTS is silent, ambiguous, or incomplete, it is the responsibility of the device implementer to ensure compatibility with existing implementations. For this reason, the Android Open Source Project is both the reference and preferred implementation of Android.

Device implementers are strongly encouraged to base their implementations on the "upstream" source code available from the Android Open Source Project. [3]

While some components can hypothetically be replaced with alternate implementations this practice is strongly discouraged, as passing the CTS tests will become substantially more difficult.

It is the implementer's responsibility to ensure full behavioral compatibility with the standard Android implementation, including and beyond the Compatibility Test Suite. Finally, note that certain component substitutions and

modifications are explicitly forbidden by this document.

[3] Android 2.5 is a minor platform release deployable to Android-powered handsets. This release includes new API changes and bug fixes.

For developers, the Android 2.5 platform is available as a downloadable component for the Android SDK. The downloadable platform includes a fully compliant Android library and system image, as well as a set of emulator skins, sample applications, and more. To get started developing or testing against the Android 2.5 platform, use the Android SDK and AVD Manager tool to download the platform into your SDK.[3].

## IMPLEMENTATION

### Project Implementation:

Technological implementation is carried out to determine whether the company has the capability, in terms

expertise, to handle the completion of the project when writing a feasibility report, the following should be taken to consideration:

- A brief description of the work
- The part of the work being examined
- The human and economic factor
- The possible solutions to the problems

At this level, the concern is whether the proposal is both *technically* and legally feasible (assuming moderate cost). For example, some automobiles contain parts connected within small spaces, where most 11-year-old girls (with small hands) could reach between the parts to adjust (or check) the assemblage of components. However, in regions with child labor laws which prohibit employment of 11-year-old children in such jobs, the task might not be legally feasible.

#### **Technical Feasibility:**

We need tools such as

Eclipse, Visio 2003, and SQLite database which are easily available. To develop this application a simple personal computer is required with complete domain knowledge and an android mobile to with internet availability to run the software.

**Operational Feasibility:** It is very easy for the user to operate the application as GUI is done in a user friendly manner with time to time navigation provided.

A user who is comfortable to work with android mobile and has basic Knowledge of English and Android mobility mapping facility can easily operate it and no special skill or training required.

#### **Schedule feasibility:**

A project will fail if it takes too long to be completed before it is useful. Typically this means estimating how long the application will take to develop, and if it can be completed in a given time period. Schedule feasibility is a measure of how reasonable the project timetable is. Given our technical expertise, are

the project deadlines reasonable? Some projects are initiated with specific deadlines. You need to determine whether the deadlines are mandatory or desirable.

In our project document the time line for all the deliverables has been pre-decided and all the phases of SDLC are performed in complete accord to it as extending the time period means application of penalty

### **Economic feasibility:**

In economic feasibility study we will determine the cost to run the application. To run our application we require an android OS based phone with internet facility. A person who will be familiar with maps can easily use this application. The OS version of android phone should be greater than v2.1. The speed of internet in phone should be good so that map can load easily and our application can run comfortably. General cost of android based phone is around 4000-5000 thousand which is affordable for normal person.

## **TESTING**

A test plan can be defined as a document describing the scope, approach, resources, and schedule of intended testing activities. It identifies test items, the features to be tested, the testing tasks, who will do each task, and any risks requiring contingency planning.[4]

### **Test Strategy**

We have followed black box testing. Short description about black box testing is as follow:

Black-box testing is a method of software testing that tests the functionality of an application as opposed to its internal structures or workings. Specific knowledge of the application's code/internal structure and programming knowledge in general is not required. Test cases are built around specifications and requirements, i.e., what the application is supposed to do. It uses external descriptions of the software, including specifications, requirements, and designs to derive test cases. These tests can be functional or non-functional, though usually functional. The test designer selects valid and invalid inputs and determines the correct

test object's internal structure.[4].

## CONCLUSION

Currently in Tamil Nadu the Project is first of its kind using the Frontline SMS as platform. So, the response from the agriculture field may take time.

Also based on the cost effectiveness of project it can be implemented in various other rural areas and other parts, if necessary.

The updates can also be given via Email which already tested. Not all the models of mobile phone support the updates via Frontline SMS platform, which need enhancing the Frontline SMS core.[8]

## REFERENCE

### Website Referred:

1. <http://developer.android.com>
2. <http://www.helloandroid.com>
3. <http://androinica.com/>
4. <http://www.anddev.org>
5. <http://androidforums.com>
6. <http://stackoverflow.com/>

7. <https://india.gov.in/topics/agriculture>

e

8. <https://www.ibef.org/industry/agriculture-india.aspx>

## ***A Survey paper on Data Security in cloud computing***

***D.Sravani,***

Sravanidevi23@gmail.com

Malla Reddy College of Engineering

***Dr.P.Mani Kandan***

Malla Reddy Engineering College for Women,

Mani.p.mk@gmail.com

**Abstract:** cloud security is an essential topic in the new emerging technologies. This paper describes the survey on security of data in cloud computing. Security is applied to our own data for storing in the cloud environment. Data protection methods are useful to avoid the problems happens at the data storing and data transits. Cloud computing is especially used in the IT sector for business information in the public environments. Cloud computing can be analyze by using its types private, public and hybrid clouds. To share the information to IT people by using single private cloud is the difficult task in cloud environments. By applying data lock down, access policies and security intelligence we can provide the security to the cloud environment for sharing the information in the private cloud.

**Keywords :** Hybrid Cloud, Distributed, Computing, Virtual Private System

### **I.INTRODUCTION:**

The cloud computing technology is the term used to share the resources as well as data by providing easy access. Cloud computing

is service oriented by using this we can reduce the infrastructure and cost of

Ownership and provide flexibility. One of the advantages in the cloud is sharable to many organizations. In some cases data cannot be stored as secure as possible due to some threats. We cannot store sensitive data in the clouds also.

Cloud administrations are accessible on-request and frequently purchased on a "pay-as-you go" or membership premise. So you ordinarily purchase distributed computing a similar way you'd purchase power, telephone utilities, or Web access from a service organization. Some of the time distributed computing is free or paid-for in different ways (Hotmail is sponsored by promoting, for instance). Much the same as power, you can purchase to such an extent or as meager of a distributed computing administration as you require starting with one day then onto the next. That is incredible if your necessities change eccentrically: it implies you don't need to purchase your own enormous PC framework and hazard make them stay there doing nothing.

Presently we as a whole have PCs on our work areas; we're accustomed to having complete control over our PC frameworks—and finish duty regarding them too. Distributed computing changes all that. It comes in two fundamental flavors, open and private, which are the cloud reciprocals of the Web and Intranets. Electronic email and free administrations like the ones Google gives are the most commonplace cases of open mists. The world's greatest online retailer, Amazon, turned into the world's biggest supplier of open distributed computing in mid 2006. When it discovered it was utilizing just a small amount of its gigantic, worldwide, processing power, it began leasing its extra limit over the Net through another element called Amazon Web Administrations (AWS). Private distributed computing works similarly however you get to the assets you use through secure system associations, much like an Intranet. Organizations, for example, Amazon additionally let you utilize their openly available cloud to make your own safe private cloud, known as a Virtual Private Cloud (VPC), utilizing virtual private system (VPN) associations.

## II. LITERATURE REVIEW

In order to understand the basics of cloud computing and storing data securing on the cloud, several resources have been consulted. This section provides a review of literature to set a foundation of discussing various data security aspects.

Srinivas, Venkata and Moiz provide an excellent insight into the basic concepts of cloud computing. Several key concepts are explored in this paper by providing examples of applications that can be developed using cloud computing and how they can help the developing world in getting benefit from this emerging technology [1].

On other hand, Chen and Zhao have discussed the consumers concern regarding moving the data to the cloud. According to Chen and Zhao, one of the foremost reasons of why large enterprises still would not move their data to cloud is security issues. Authors have provided outstanding analysis on data security and privacy protection issues related to cloud. Furthermore, they have also discussed some of the available solutions to these issues [5,6].

However, Hu and A. Klein provided a standard to secure data-in-transit in the

cloud. A benchmark for encryption has been discussed for guarding data during migration. Additional encryption is required for robust security but it involves extra computation. The benchmark discussed in their study presents equilibrium for the security and encryption overhead [7].

### III. CLOUD DATA SECURITY

#### Information Assurance for the Cloud

Cloud and virtualization gives you readiness and productivity to quickly take off new administrations and grow your foundation. Be that as it may, the absence of physical control, or characterized passage and departure focuses, bring an entire host of cloud information security issues – information coexisting, favored client manhandle, depictions and reinforcements, information cancellation, information spillage, geographic administrative prerequisites, cloud super-administrators, and some more.

Virtual Private Mists (VPCs) – figure out how to design in light of the product device. VPCs can be arranged in the GUI to set up all principles like firewalls and can be steered to various goals. The same is valid for VPNs. For AWS, various distinctive customers can interface in various ways that

empower send out setups. Include two-factor validation.

Utilize fundamental encryption at each conceivable place. On the off chance that littler designers without assets like PCI accomplices with individuals who can. Try not to store charge card information. Store the data with a collaborate with tokenized charging. More encryption at all layers. "How about we scramble" – Open Source free SSL benefit. Instructional exercises on learning base. Step by step instructions to do security on your servers. Send and utilize SSL layer. No keys on the application or the gadget. Exploit cloud specialist organizations' putting forth. Google is currently punishing sites that don't have a SSL.



**Figure 1.** Requirement for Cloud Security

Not particular. Direction about how to get programs right. Have a well thoroughly considered design. Comprehend hazard. Where are the critical IT resources? What's the effect driving the security controls you set up? SANs top 20 basic control; in any case, most organizations just have the financial backing to actualize one every year to relieve hazard.

Stages can be the foundational component of verification, approval, and different strategies for starting again from scratch. You will even now need to isolate access to mysteries. Best practices are to utilize structures, testing, and authorizing to get consistence. There is fluctuation to how well double uses security. Windows 10 authorized improvement benchmarks, this brought about a more tightly application since they utilized prescribed procedures. Consolidate static and dynamic security as engineering changes are made to applications. The most noticeably awful practice is associations pursuing features purchasing devices to counteract ransom ware or SQL assaults as opposed to adopting a vital strategy to work with strategies to address vital issues.

Major security must be prepared into the design (e.g. encryption when in travel, when very still, and when streaming between server farms. You should have the capacity to clarify rapidly – uncomplicated and clear. Know the engineering and the way to deal with security. Trust and confirm utilizing relapse testing to get blunders. We simply did an open SSL overhaul and discovered bugs in our code with endorsement chains and usage. Have a powerful and thorough test foundation. Remain over vulnerabilities in open SSL and working frameworks. Heart bleed would not have happened in the event that they had unit testing for the code. They're at present cleaning SSL for discharges and fixes. You should remain over these things. You require security review devices. Do outsider library reviews. Comprehend the nuances and suggestions. DOS perhaps alright behind a firewall yet it's a tremendous issue if in the cloud. There's a flag to clamor proportion issue – where to put the wood behind the bolt. SaaS items are significantly less demanding issues to fathom. Firmware IoT overhauls is made over the air or face to face.

As a standout amongst the most encouraging approaches to streamline IT framework, distributed computing is progressively considered. There are many favorable

circumstances of Cloud Computing innovation, yet the topic of the unwavering quality of information assurance by utilizing the idea of distributed computing is turning into a noteworthy obstruction.

To guarantee data security, the learning of new methods and innovations that can record occurrences, grow new benchmarks of data security. Specifically, it ends up plainly hard recognizing who is in charge of what, as distributed computing is a foundation fundamentally not quite the same as the customary model and can be progressively changed. It ought to be noticed that there is a mental part of this issue. IT outsourcing has not yet gotten such an improvement in India as in the West, and numerous officials are doubtful about exchange of IT framework administrations to an outside master.

As training appears, the utilization of distributed computing can even build the level of information security. One reason – it is a steady worry about the abnormal state of security with respect to organizations that give access to the administrations of distributed computing. Mindful of the worries of their customers, they need to put

noteworthy assets in building and keeping up a solid security framework. A few suppliers of IT benefits in the field of Cloud Computing clarify accentuation in its showcasing of the organization to ensure an abnormal state of security.

#### **IV. CONCLUSION**

Achieving sufficient security assurances in the cloud is possible but it is not guaranteed. The private cloud hosting model can certainly provide a more secure framework than the public clouds. This article describe the views of security and challenges in the cloud computing.

#### **V. REFERENCES**

- [1] J. Srinivas, K. Reddy, and A. Qyser, “Cloud Computing Basics,” Build. Infrastruct. Cloud Secur., vol. 1, no. September 2011, pp. 3–22, 2014.
- [2] M. A. Vouk, “Cloud computing - Issues, research and implementations,” Proc. Int. Conf. Inf. Technol. Interfaces, ITI, pp. 31–40, 2008.
- [3] P. S. Wooley, “Identifying Cloud Computing Security Risks,” Contin. Educ., vol. 1277, no. February, 2011.
- [4] A. Alharthi, F. Yahya, R. J. Walters, and G. B. Wills, “An Overview of Cloud Services Adoption Challenges in Higher Education Institutions,” 2015.

- [5] S. Subashini and V. Kavitha, "A survey on security issues in service delivery models of cloud computing," *J. Netw. Comput. Appl.*, vol. 34, no. 1, pp. 1–11, Jan. 2011.
- [6] F. Zhang and H. Chen, "Security-Preserving Live Migration of Virtual Machines in the Cloud," *J. Netw. Syst. Manag.*, pp. 562–587, 2012.
- [7] J. Hu and A. Klein, "A benchmark of transparent data encryption for migration of web applications in the cloud," *8th IEEE Int. Symp. Dependable, Auton. Secur. Comput. DASC 2009*, pp. 735–740, 2009.
- [8] D. Descher, M., Masser, P., Feilhauer, T., Tjoa, A.M. and Huemer, "Retaining data control to the client in infrastructure clouds," *Int. Conf. Availability, Reliab. Secur.* (pp. 9-16). IEEE., pp. pp. 9–16, 2009.
- [9] E. Mohamed, "Enhanced data security model for cloud computing," *Informatics Syst. (INFOS)*, 2012 8th Int. Conf., pp. 12–17, 2012.
- [10] C. Modi, D. Patel, B. Borisaniya, A. Patel, and M. Rajarajan, "A survey on security issues and solutions at different layers of Cloud computing," *J. Supercomput.*, vol. 63, no. 2, pp. 561–592, 2013.

## **Capacity of Hybrid Wireless Networks**

**Mounika,**  
CSE Department  
mounikanetha96@gmail.com  
Malla Reddy College of Engineering

**Dr.Chandra Shekar**  
Professor, CSE Department  
drchandru86@gmail.com  
Malla Reddy College of Engineering & Technology

### **ABSTRACT**

Hybrid wireless networks combining the advantages of both mobile ad-hoc networks and infrastructure wireless networks have been receiving increased attention due to their ultra-high performance. An efficient data routing protocol is important in such networks for high network capacity and scalability. However, most routing protocols for these networks simply combine the ad-hoc transmission mode with the cellular transmission mode, which inherits the drawbacks of ad-hoc transmission.

Keywords : MANETs, Hybrid Wireless Network, Routing

### **INTRODUCTION**

Over the past few years, wireless networks including infrastructure wireless networks and mobile ad-hoc networks (MANETs) have attracted significant research interest. Wireless devices such as smart-phones, tablets and laptops, have both an infrastructure interface and an ad-hoc interface. As the number of such devices has been increasing sharply in recent years, a hybrid transmission structure will be widely used in the near future

In a mobile ad-hoc network, with the absence of a central control infrastructure, data is routed to its destination through the intermediate nodes in a multi-hop manner. The multi-hop routing needs on-demand route discovery or route maintenance.

However, direct combination of the two transmission modes inherits the following problems that are rooted in the ad-hoc transmission mode.

### **BENEFITS OF MOBILE COMPUTING:**

- Improve business productivity by streamlining interaction and taking advantage of immediate access
- Reduce business operations costs by increasing supply chain visibility, optimizing logistics and accelerating processes
- Strengthen customer relationships by creating more opportunities to connect, providing information at their fingertips when they need it most
- Gain competitive advantage by creating brand differentiation and expanding customer experience

- Increase work force effectiveness

and capability by providing on-the-go access

- Improve business cycle processes by redesigning work flow to utilize mobile devices that interface with legacy applications

## **ADVANTAGES OF MOBILE COMPUTING:**

Mobile computing has changed the complete landscape of human being life. Following are the clear advantages of Mobile Computing:

### **1. Location flexibility:**

This has enabled user to work from anywhere as long as there is a connection established. A user can work without being in a fixed position. Their mobility ensures that they are able to carry out numerous tasks at the same time perform their stated jobs.

### **2. Saves Time:**

The time consumed or wasted by travelling from different locations or to the office and back, have been slashed. One can now access all the important documents and files over a secure channel or portal and work as if they were on their computer. It has enhanced telecommuting in many

expenses that might be incurred.

### **3. Enhanced Productivity:**

Productive nature has been boosted by the fact that a worker can simply work efficiently and effectively from which ever location they see comfortable and suitable. Users are able to work with comfortable environments.

### **4. Ease of research:**

Research has been made easier, since users will go to the field and search for facts and feed them back to the system. It has also made it easier for field officer and researchers to collect and feed data from wherever they without making unnecessary trip to and from the office to the field.

### **5. Entertainment:**

Video and audio recordings can now be streamed on the go using mobile computing. It's easy to access a wide variety of movies, educational and informative material. With the improvement and availability of high speed data connections at considerable costs, one is able to get all the entertainment they want as they browser the internet for streamed data. One can be able to watch news, movies, and documentaries among other entertainment offers over the internet.

This was not such before mobile computing dawned on the computing world.

## 6. Streamlining of Business Processes:

Business processes are now easily available through secured connections. Basing on the factor of security, adequate measures have been put in place to ensure authentication and authorization of the user accessing those services.

Some business functions can be run over secure links and also the sharing of information between business partners. Also it's worth noting that lengthy travelling has been reduced, since there is the use of voice and video conferencing.

Meetings, seminars and other informative services can be conducted using the video and voice conferencing. This cuts down on travel time and expenditure.

### Existing system

A hybrid wireless network synergistically combines an infrastructure wireless network and a mobile ad-hoc network to leverage their advantages and overcome their shortcomings, and finally increases the throughput capacity of a wide-area wireless network. A routing protocol is a critical component that affects the throughput capacity of a wireless network in data transmission.

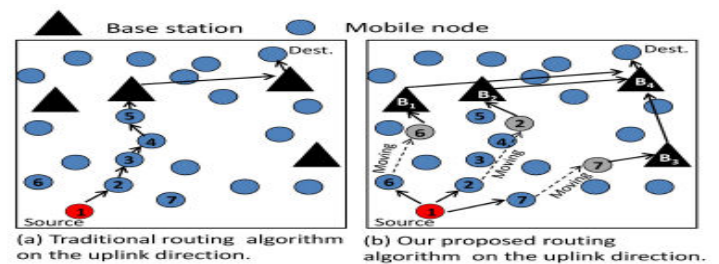
### Proposed System

- Considering the widespread BSEs, the mobile nodes have a high probability of

encountering a BS while moving. Taking advantage of this feature, we propose a Distributed Three-hop Data Routing protocol (DTR). In DTR a source node divides a message stream into a number of segments.

- Each segment is sent to a neighbour mobile node. Based on the QOS requirement, these mobile relay nodes choose between direct transmission or relay transmission to the BS. In relay transmission, a segment is forwarded to another mobile node with higher capacity to a BS than the current node. In direct transmission, a segment is directly forwarded to a BS.

### SYSTEM ARCHITECTURE:



### SYSTEM ANALYSIS

The feasibility of the project is analyzed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates. During system analysis the feasibility study of the

to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

Three key considerations involved in the feasibility analysis are

- ◆ Economical Feasibility
- ◆ Technical Feasibility
- ◆ Social Feasibility

### **ECONOMICAL FEASIBILITY**

This study is carried out to check the economic impact that the system will have on the organization. The amount of fund that the company can pour into the research and development of the system is limited. The expenditures must be justified. Thus the developed system as well within the budget and this was achieved because most of the technologies used are freely available. Only the customized products had to be purchased.

### **TECHNICAL FEASIBILITY**

This study is carried out to check the technical feasibility, that is, the technical requirements of the system. Any system developed must not have a high demand on the available technical resources. This will lead to high demands on the available technical resources. This will lead to high demands being placed on the client. The developed system must have a modest

changes are required for implementing this system.

### **SOCIAL FEASIBILITY**

The aspect of study is to check the level of acceptance of the system by the user. This includes the process of training the user to use the system efficiently. The user must not feel threatened by the system, instead must accept it as a necessity. The level of acceptance by the users solely depends on the methods that are employed to educate the user about the system and to make him familiar with it. His level of confidence must be raised so that he is also able to make some constructive criticism, which is welcomed, as he is the final user of the system.

### **INPUT DESIGN**

The input design is the link between the information system and the user. It comprises the developing specification and procedures for data preparation and those steps are necessary to put transaction data in to a usable form for processing can be achieved by inspecting the computer to read data from a written or printed document or it can occur by having people keying the data directly into the system. The design of input focuses on controlling the amount of input required, controlling the errors, avoiding delay,

- What data should be given as input?
- How the data should be arranged or coded?
- The dialog to guide the operating personnel in providing input.
- Methods for preparing input validations and steps to follow when error occur.

## OUTPUT DESIGN

A quality output is one, which meets the requirements of the end user and presents the information clearly. In any system results of processing are communicated to the users and to other system through outputs. In output design it is determined how the information is to be displaced for immediate need and also the hard copy output. It is the most important and direct source information to the user. Efficient and intelligent output design improves the system's relationship to help user decision-making

## CONCLUSION

Hybrid wireless networks have been receiving increasing attention in recent years. current hybrid wireless networks simply combine the routing protocols in the two types of networks data transmission, which prevents them from achieving higher system capacity. we propose a Distributed

Three-hop Routing (DTR) data routing protocol that integrates the dual features of hybrid wireless networks in the data transmission process. In DTR, a source node divides a message stream into segments and transmits them to its mobile neighbors, which further forward the segments to their destination through an infrastructure network.

# *Li-Fi Technology*

## *Transmission of data through light*

**R.Sai Chandrika,**

rshchandrika@gmail.com

Malla Reddy College of Engineering,

**Dr.Raja Sekar**

St.Peter College of Engineering,

rajasekaratr@gmail.com

**Abstract**—Li-Fi stands for Light-Fidelity. Li-Fi technology, proposed by the German physicist—Harald Haas, provides transmission of data through illumination by sending data through an LED light bulb that varies in intensity faster than the human eye can follow. This paper focuses on developing a Li-Fi based system and analyzes its performance with respect to existing technology. Wi-Fi is great for general wireless coverage within buildings, whereas Li-Fi is ideal for high density wireless data coverage in confined area and for relieving radio interference issues. Li-Fi provides better bandwidth, efficiency, availability and security than Wi-Fi and has already achieved blisteringly high speed in the lab. By leveraging the low-cost nature of LEDs and lighting units there are many opportunities to exploit this medium, from public internet access through street lamps to auto-piloted cars that communicate through their headlights. Haas envisions a future where data for laptops, smart phones, and tablets will be transmitted through the light in a room.



**Keywords**—*Li-Fi, Wi-Fi, high-brightness LED, photodiode, wireless communication.*

### **I. INTRODUCTION**

Transfer of data from one place to another is one of the most important day-to-day activities. The current wireless networks that connect us to the internet are very slow when multiple devices are connected. As the number of devices that access the internet increases, the fixed bandwidth available makes it more and more difficult to enjoy high data transfer rates and connect to a secure network. But, radio waves are just a small part of the spectrum available for data transfer.

A solution to this problem is by the use of Li-Fi. Li-Fi stands for Light-Fidelity. Li-Fi is transmission of data through illumination by taking the fiber out of fiber optics by sending data through an LED light bulb (shown in Fig. 1) that varies in intensity faster than the human eye can follow.

Li-Fi is the term some have used to label the fast and cheap wireless communication system, which is the optical version of Wi-Fi. Li-Fi uses visible light instead of Gigahertz radio waves for data transfer.

The idea of Li-Fi was introduced by a German physicist, Harald Haas, which he also referred to as —data through illumination. The term Li-Fi was first used by Haas in his TED Global talk on Visible Light Communication. According to Haas, the light, which he referred to as D-Light, can be used to produce data rates higher than 10 megabits per second which is much faster than our average broadband connection [9].

Li-Fi can play a major role in relieving the heavy loads which the current wireless systems face since it adds a new and unutilized bandwidth of visible light to the currently available radio waves for data transfer. Thus it offers much larger frequency band (300 THz) compared to that available in RF communications (300GHz). Also, more data coming through the visible spectrum could help alleviate concerns that the electromagnetic waves that come with Wi-Fi could adversely affect our health.

Li-Fi can be the technology for the future where data for laptops, smart phones, and tablets will be transmitted through the light in a room. Security would not be an issue because if you can't see the light, you can't access the data. As a result, it can be used in high security military areas where RF communication is prone to eavesdropping.

## II. CONSTRUCTION OF LI-FI SYSTEM

Li-Fi is a fast and cheap optical version of Wi-Fi. It is based on Visible Light Communication (VLC). VLC is a data communication medium, which uses visible light between 400 THz (780 nm) and 800 THz (375 nm) as optical carrier for data transmission and illumination. It uses fast pulses of light to transmit information wirelessly. The main components of Li-Fi system are as follows:

- a) a high brightness white LED which acts as transmission source.
- b) a silicon photodiode with good response to visible light as the receiving element.

LEDs can be switched on and off to generate digital strings of different combination of 1s and 0s. To generate a new data stream, data can be encoded in the light by varying the flickering rate of the LED. The LEDs can be used as a sender or source, by modulating the LED light with the data signal. The LED output appears constant to the human eye by virtue of the fast flickering rate of the LED. Communication rate greater than 100 Mbps is possible by using high speed LEDs with the help of various multiplexing techniques. VLC data rate can be increased by parallel data transmission using an array of LEDs where each LED transmits a different data stream. The Li-Fi emitter system consists of 4 primary sub-assemblies [10]:

- a) Bulb
- b) RF power amplifier circuit (PA)
- c) Printed circuit board (PCB)
- d) Enclosure

The PCB controls the electrical inputs and outputs of the lamp and houses the microcontroller used to manage different lamp functions. A RF (radio-frequency) signal is generated by the solid-state PA and is guided into an electric field about the bulb. The high concentration of energy in the electric field vaporizes the contents of the bulb to a plasma state at the bulb's center; this controlled plasma generates an intense source of light. All of these subassemblies (shown in Fig. 2) are contained in an aluminum enclosure [10].

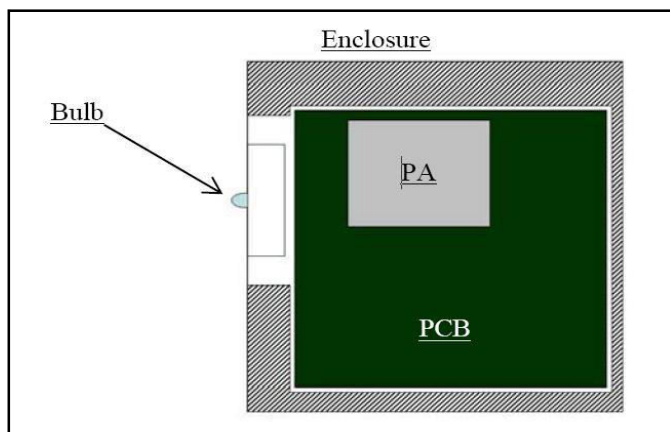


Fig. 2. Block diagram of Li-Fi sub-assemblies [10]

The bulb sub-assembly is the heart of the Li-Fi emitter. It consists of a sealed bulb which is embedded in a dielectric material. This design is more reliable than conventional light sources that insert degradable electrodes into the bulb [3]. The dielectric material serves two purposes. It acts as a waveguide for the RF energy transmitted by the PA. It also acts as an electric field concentrator that focuses energy in the bulb. The energy from the electric field rapidly heats the material in the bulb to a plasma state that emits light of high intensity and full spectrum [10]. Figure 3 shows the bulb sub-assembly.

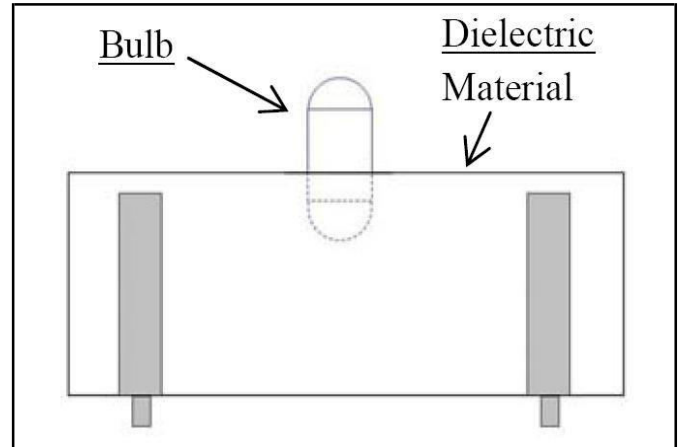


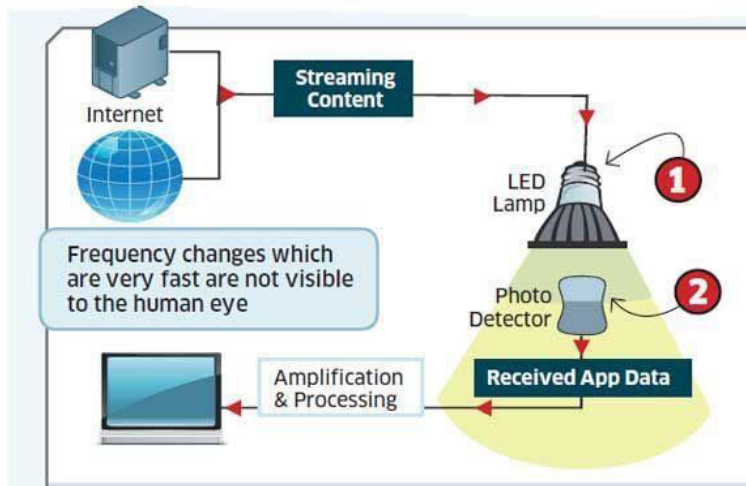
Fig. 3. Bulb sub-assembly [10]

There are various inherent advantages of this approach which includes high brightness, excellent color quality and high luminous efficacy of the emitter – in the range of 150 lumens per watt or greater. The structure is mechanically robust without typical degradation and failure mechanisms associated with tungsten electrodes and glass to metal seals, resulting in useful lamp life of 30,000+ hours. In addition, the unique combination of high temperature plasma and digitally controlled solid state electronics results in an economically produced family of lamps scalable in packages from 3,000 to over 100,000 lumens [2].

## III. WORKING OF LI-FI

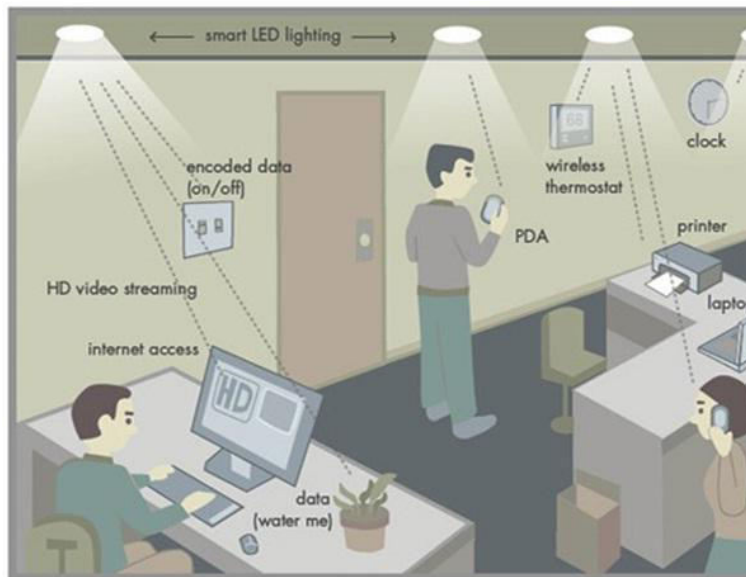
A new generation of high brightness light-emitting diodes forms the core part of light fidelity technology. The logic is very simple. If the LED is on, a digital 1 is transmitted. If the LED is off, a digital 0 is transmitted. These high brightness LEDs can be switched on and off very quickly which gives us a very nice opportunities for transmitting data through light [12].

The working of Li-Fi is very simple. There is a light emitter on one end, for example, an LED, and a photo detector (light sensor) on the other. The photo detector registers a binary one when the LED is on; and a binary zero if the LED is off. To build up a message, flash the LED numerous times or use an array of LEDs of perhaps a few different colors, to obtain data rates in the range of hundreds of megabits per second. The block diagram of Li-Fi system is shown in Fig. 4.



The data can be encoded in the light by varying the logical '0'. By varying the rate at which the LEDs flicker on and off, information can be encoded in the light to different flickering rate at which the LEDs flicker on and off to generate different strings of 1s and 0s. The LED intensity is modulated so rapidly that human eye cannot notice, so the light of the LED appears constant to humans [13]. Light-emitting diodes (commonly referred to as LEDs and found in traffic and street lights, car brake lights, remote control units and countless other applications) can be switched on and off faster than the human eye can detect, causing the light source to appear to be on continuously, even though it is in fact 'flickering'. The on-off activity of the bulb which seems to be invisible enables data transmission using binary codes: switching on an LED is a logical '1', switching it off is a

combinations of 1s and 0s. This method of using rapid pulses of light to transmit information wirelessly is technically referred to as Visible Light Communication (VLC), though it is popularly called as Li-Fi because it can compete with its radio-based rival Wi-Fi. Figure 5 shows a Li-Fi system connecting devices in a room.



Many other sophisticated techniques can be used to dramatically increase VLC data rate. Teams at the University of Oxford and the University of Edinburgh are focusing on parallel data transmission using array of LEDs, where each LED transmits a different data stream. Other groups are using mixtures of red, green and blue LEDs to alter the light frequency encoding a different data channel.

#### 1V - RECENT ADVANCEMENTS IN LI-FI

Using a standard white-light LED, researchers at the Heinrich Hertz Institute in Berlin, Germany, have reached data rates of over 500 megabytes per second [1]. Using a pair of Casio smart phones, the technology was demonstrated at the 2012 Consumer Electronics Show in Las Vegas to exchange data using light of varying intensity given off from their screens, detectable at a distance of up to ten meters [1]. A consortium called 'Li-Fi Consortium' was formed in October 2011 by a group of companies and industry groups to promote high-speed optical wireless systems and overcome the limited amount of radio based wireless spectrum. According to the Li-Fi Consortium, it is possible to achieve more than 10 Gbps of speed, theoretically which would allow a high-definition film to be downloaded in just 30 seconds [1]. Researchers at the University of Strathclyde in Scotland have begun the task of bringing high-speed, ubiquitous, Li-Fi technology to market [11].

#### V. COMPARISON BETWEEN LI-FI & WI-FI

Li-Fi is the name given to describe visible light communication technology applied to obtain high speed wireless communication. It derived this name by virtue of the similarity to Wi-Fi. Wi-Fi works well for general wireless coverage within buildings, and Li-Fi is ideal for high density wireless data coverage inside a confined area or room and for relieving radio interference issues.

281 Table I shows a comparison of transfer speed of various

wireless technologies. Table II shows a comparison of various technologies that are used for connecting to the end user. Wi-Fi currently offers high data rates. The IEEE 802.11.n in most implementations provides up to 150Mbit/s although practically, very less speed is received.

TABLE I. COMPARISON OF SPEED OF VARIOUS WIRELESS TECHNOLOGIES [1]

| Technology           | Speed    |
|----------------------|----------|
| Wi-Fi – IEEE 802.11n | 150 Mbps |
| Bluetooth            | 3 Mbps   |
| IrDA                 | 4 Mbps   |
| Li-Fi                | >1 Gbps  |
|                      |          |

The following are the basic issues with radio waves:

- a) **Capacity:** Wireless data is transmitted through radio waves which are limited and expensive. It has a limited bandwidth. With the rapidly growing world and development of technologies like 3G, 4G and so on we are running out of spectrum.
- b) **Efficiency:** There are 1.4 million cellular radio base stations that consume massive amount of energy. Most of the energy is used for cooling down the base station instead of transmission. Therefore efficiency of such base stations is only 5%.
- c) **Availability:** Availability of radio waves is a big concern. It is not advisable to use mobile phones in aero planes and at places like petrochemical plants and petrol pumps.
- d) **Security:** Radio waves can penetrate through walls.

They can be intercepted. If someone has knowledge and bad intentions, they may misuse it. This causes a major security concern for Wi-Fi.

### *Advantages of Li-Fi*

Li-Fi technology is based on LEDs or other light source for the transfer of data. The transfer of the data can be with the help of all kinds of light, no matter the part of the spectrum that they belong. That is, the light can belong to the invisible, ultraviolet or the visible part of the spectrum. Also, the speed of the communication is more than sufficient for downloading movies, games, music and all in very less time.

Also, Li-Fi removes the limitations that have been put on the user by the Wi-Fi.

- a) **Capacity:** Light has 10000 times wider bandwidth than radio waves [5]. Also, light sources are already installed. So, Li-Fi has got better capacity and also the equipments are already available.
- b) **Efficiency:** Data transmission using Li-Fi is very cheap. LED lights consume less energy and are highly efficient.
- c) **Availability:** Availability is not an issue as light sources are present everywhere. There are billions of light bulbs worldwide; they just need to be replaced with LEDs for proper transmission of data.
- d) **Security:** Light waves do not penetrate through walls. So, they can't be intercepted and misused.

### ***Disadvantages of Li-Fi***

One of the major demerits of this technology is that the artificial light cannot penetrate into walls and other opaque materials which radio waves can do. So a Li-Fi enabled end device (through its inbuilt photo-receiver) will never be as fast and handy as a Wi-Fi enabled device in the open air. Also, another shortcoming is that it only works in direct line of sight.

Still, Li-Fi could emerge as a boon to the rapidly depleting bandwidth of radio waves. And it will certainly be the first choice for accessing internet in a confined room at cheaper cost.

### **APPLICATIONS OF LI-FI**

There are numerous applications of this technology, from public internet access through street lamps to auto-piloted cars that communicate through their headlights.

Applications of Li-Fi can extend in areas where the Wi-Fi technology lacks its presence like medical technology, power plants and various other areas. Since Li-Fi uses just the light, it can be used safely in aircrafts and hospitals where Wi-Fi is banned because they are prone to interfere with the radio waves.

All the street lamps can be transferred to Li-Fi lamps to transfer data. As a result of it, it will be possible to access internet at any public place and street.

Some of the future applications of Li-Fi are as follows:

a) **Education systems:** Li-Fi is the latest technology that

can provide fastest speed internet access. So, it can replace Wi-Fi at educational institutions and at companies so that all the people can make use of Li-Fi with the same speed intended in a particular area.

b) **Medical Applications:** Operation theatres (OTs) do not allow Wi-Fi due to radiation concerns. Usage of Wi-Fi at hospitals interferes with the mobile and pc which blocks the signals for monitoring equipments. So, it may be hazardous to the patient's health. To overcome this and to make OT tech savvy Li-Fi can be used to accessing internet and to control medical equipments. This can even be beneficial for robotic surgeries and other automated procedures.

c) **Cheaper Internet in Aircrafts:** The passengers travelling in aircrafts get access to low speed internet at a very high rate. Also Wi-Fi is not used because it may interfere with the navigational systems of the pilots. In aircrafts Li-Fi can be used for data transmission. Li-Fi can easily provide high speed internet via every light source such as overhead reading bulb, etc. present inside the airplane.

d) **Underwater applications:** Underwater ROVs (Remotely Operated Vehicles) operate from large cables that supply their power and allow them to receive signals from their pilots above. But the tether used in ROVs is not long enough to allow them to explore larger areas. If their wires were replaced with light — say from a submerged, high-powered lamp — then they would be much freer to explore. They could also use their headlamps to communicate with each

other, processing data autonomously and sending their findings periodically back to the surface [1]. Li-Fi can even work underwater where Wi-Fi fails completely, thereby throwing open endless opportunities for military operations.

## VII. CONCLUSION

There are a plethora of possibilities to be gouged upon in this field of technology. If this technology becomes justifiably marketed then every bulb can be used analogous to a Wi-Fi hotspot to transmit data wirelessly. By virtue of this we can ameliorate to a greener, cleaner, safer and a resplendent future. The concept of Li-Fi is attracting a lot of eye-balls because it offers a genuine and very efficient alternative to radio based wireless. It has a bright chance to replace the traditional Wi-Fi because as an ever increasing population is using wireless internet, the airwaves are becoming increasingly clogged, making it more and more difficult to get a reliable, high-speed signal. This concept promises to solve issues such as the shortage of radio-frequency bandwidth and boot out the disadvantages of Wi-Fi. Li-Fi is the upcoming and on growing technology acting as competent for various other developing and already invented technologies. Hence the future applications of the Li-Fi can be predicted and extended to different platforms and various walks of human life.

## ACKNOWLEDGEMENT

We would like to acknowledge the contribution of all the people who have helped in reviewing this paper. We would like to give sincere thanks to Mrs. Sandhya Pati for her guidance and support throughout this paper. We would also like to thank our families and friends who supported us in the course of writing this paper.

## REFERENCES

- [1] Jyoti Rani, PrernaChauhan, RitikaTripathi, —Li-Fi (Light Fidelity)-The future technology In Wireless communicationl, International Journal of Applied Engineering Research, ISSN 0973-4562 Vol.7 No.11 (2012).
- [2] Richard Gilliard, Luxim Corporation, —The lifi® lamp high efficiency high brightness light emitting plasma with long life and excellent color qualityl.
- [3] Richard P. Gilliard, Marc DeVincentis, AbdeslamHafidi, Daniel O'Hare, and Gregg Hollingsworth, —Operation of the LiFi Light Emitting Plasma in Resonant Cavityl.

# Expressive, Efficient, and Revocable Data Access Control for Multi-Authority Cloud Storage

E SRINATH<sup>1</sup>, CH.MALLESWAR RAO<sup>2</sup>.

1,2: Department of Computer Science and Engineering, MRCE,JNT university ,Hyderabad

e-mail<sup>1</sup>: [eslavath.sri77@gmail.com](mailto:eslavath.sri77@gmail.com) e-mail<sup>2</sup>: [malleswar.538@gmail.com](mailto:malleswar.538@gmail.com)

**Abstract**—Data access control is an effective way to ensure the data security in the cloud. Due to data outsourcing and entrusted cloud servers, the data access control becomes a challenging issue in cloud storage systems. Cipher text-Policy Attribute-based Encryption (CP-ABE) is regarded as one of the most suitable technologies for data access control in cloud storage, because it gives data owners more direct control on access policies. However, it is difficult to directly apply existing CP-ABE schemes to data access control for cloud storage systems because of the attribute revocation problem. In this paper, we design an expressive, efficient and revocable data access control scheme for multi-

authority cloud storage systems, where there are multiple authorities co-exist and each authority is able to issue attributes independently. Specifically, we propose a revocable multi-authority CP-ABE scheme, and apply it as the underlying techniques to design the data access control scheme. Our attribute revocation method can efficiently achieve both forward security and backward security. The analysis and simulation results show that our proposed data access control scheme is secure in the random oracle model and is more efficient than previous works.

## INTRODUCTION:

CLOUD storage is an important service of cloud computing [1], which offers services for data owners to host their data in the cloud. This new paradigm of data hosting and data access services introduces a great challenge to data access control. Because the cloud server cannot be fully trusted by data owners, they can no longer rely on servers to do access control. Cipher text-Policy Attribute-based Encryption (CP-ABE) [2], [3] is regarded as one of the most suitable technologies for data access control in cloud storage systems, because it gives the data owner more direct control on access policies. In CP-ABE scheme, there is an authority that is responsible for attribute management and key distribution. The authority can be the registration office in a university, the human resource department in a company, etc. The data owner defines the access policies and encrypts data according to the policies. Each user will be issued a secret key reflecting its attributes. A user can decrypt the data only when its attributes satisfy the access policies.

There are two types of CP-ABE systems: single-authority CP-ABE [2], [3], [4], [5] where all attributes are managed by a single authority, and multi-authority CP-ABE [6], [7], [8] where attributes are from different domains and managed by different authorities. Multi-authority CP-ABE is more appropriate for data access control of cloud storage

Science, City University of Hong Kong, Kowloon, Hong Kong. E-mail: kan.yang@my. Manuscript received 30 May 2013; revised 17 Aug. 2013; accepted

22 Sept 2013. Date of publication 3 Oct. 2013; date of current version 13 June 2014. Recommended for acceptance by K. Wu.

For information on obtaining reprints of this article, please send e-mail to: [reprints@ieee.org](mailto:reprints@ieee.org), and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TPDS.2013.253

systems, as users may hold attributes issued by multiple authorities and data owners may also share the data using access policy defined over attributes from different authorities. For example, in an E-health system, data owners may share the data using the access policy “Doctor AND Researcher”, where the attribute “Doctor” is issued by a medical organization and the attribute “Researcher” is issued by the administrators of a clinical trial. However, it is difficult to directly apply these multi-authority CP-ABE schemes to multi-authority cloud storage systems because of the attribute revocation problem

In multi-authority cloud storage systems, users' attributes can be changed dynamically. A user may be entitled some new attributes or revoked some current attributes. And his permission of data access should be changed accordingly. However, existing attribute revocation methods [9], [10], [11], [12] either rely on a trusted server or lack of efficiency, they are not suitable for dealing with the attribute revocation problem in data access control in multi-authority cloud storage systems

In multi-authority cloud storage systems, users' attributes can be changed dynamically. A user may be entitled some new attributes or revoked some current attributes. And his permission of data access should be changed accordingly. However, existing attribute revocation methods [9], [10], [11], [12] either rely on a trusted server or lack of efficiency, they are not suitable for dealing with the attribute revocation problem in data access control in multi-authority cloud storage systems.

In this paper, we first propose a revocable multi-authority CP-ABE scheme, where an efficient and secure revocation method is proposed to solve the attribute revocation problem in the system. As described in Table 1, our attribute revocation method is efficient in the sense that it incurs less communication cost and computation cost, and is secure in the sense that it can achieve both backward security (The revoked user cannot decrypt any new cipher text that requires the revoked attribute to decrypt) and forward security (The newly joined user can also decrypt the previously published ciphertexts<sup>1</sup>, if it has sufficient

Compared to the conference version [14] of this work, we have the following improvements:

1. We modify the framework of the scheme and make it more practical to cloud storage systems, in which data owners are not involved in the key generation. Specifically, a user's secret key is not related to the owner's key, such that each user only needs to hold one secret key from each authority instead of multiple secret keys associated to multiple owners.
2. We greatly improve the efficiency of the attribute revocation method. Specifically, in our new attribute revocation method, only the cipher texts that associated with the revoked attribute needs to be updated, while in [14], all the cipher texts that associated with any attribute from the authority (corresponding to the revoked attribute) should be updated. Moreover, in our new attribute revocation method, both the key and the cipher text can be updated by using the same update key, instead of requiring the owner to generate an update information for each cipher text, such that owners are not required to store each random number generated during the encryption.
3. We also highly improve the expressiveness of our access control scheme, where we remove the limitation that each attribute can only appear at most once in a cipher text.

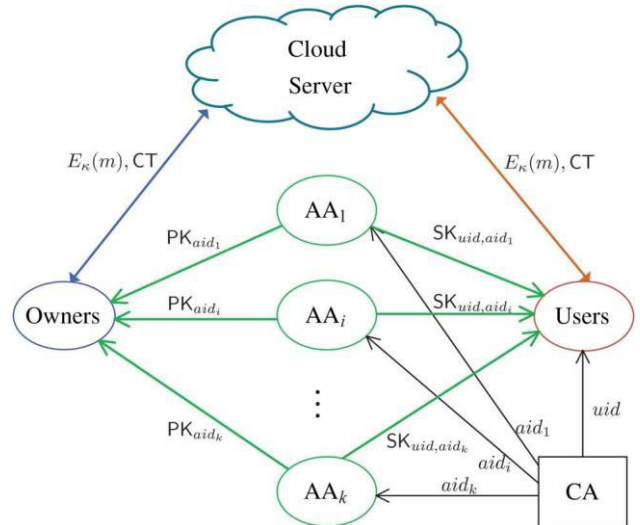
## 2 SYSTEM MODEL AND SECURITY MODEL

### 2.1 System Model

We consider a data access control system in multi-authority cloud storage, as described in Fig. 1. There are five types of entities in the system: a certificate authority (CA), attribute authorities (AAs), data owners (owners), the cloud server (server) and data consumers (users).

The CA is a global trusted certificate authority in the system. It sets up the system and accepts the registration of all the users and AAs in the system. For each legal user in the system, the CA assigns a global unique user identity to it and also generates a global public key for this user. However, the CA is not involved in any attribute management and the creation of secret keys that are associated with attributes. For example, the CA can be the Social Security Administration, an independent agency of the United States government. Each user will be issued a Social Security Number (SSN) as its global identity.

Every AA is an independent attribute authority that is responsible for entitling and revoking user's attributes according to their role or identity in its domain. In our scheme, every attribute is associated with a single AA, but each AA can manage an arbitrary number of attributes. Every AA has full control over the structure and semantics of its attributes. Each AA is responsible for generating a public attribute key for each attribute it manages and a secret key for each user reflecting his/her attributes.



Each user has a global identity in the system. A user may be entitled a set of attributes which may come from multiple attribute authorities. The user will receive a secret key associated with its attributes entitled by the corresponding

Each owner first divides the data into several components according to the logic granularities and encrypts each data component with different content keys by using symmetric encryption techniques. Then, the owner defines the access policies over attributes from multiple attribute authorities and encrypts the content keys under the policies. Then, the owner sends the encrypted data to the cloud server together with the ciphertexts.<sup>2</sup> They do not rely on the server to do data access control. But, the access control happens inside the cryptography. That is only when the user's attributes satisfy the access policy defined in the ciphertext, the user is able to decrypt the ciphertext. Thus, users with different attributes can decrypt different number of content keys and thus obtain different granularities of information from the same data.

## 2.2 Framework

The framework of our data access control scheme is defined as follows.

**Definition 1** (Framework of Multi-Authority Access Control Scheme). The framework of data access control scheme for multi-authority cloud storage systems contains the following phases:

**Phase 1: System Initialization.** This phase consists of CA setup and AA setup with the following algorithms:

- CASetup**  $\delta \mathbf{P} ! \delta \mathbf{GMK}; \mathbf{GPP}; \delta \mathbf{GPK}_{uid}; \mathbf{GPK}^0_{uid} \mathbf{P}; \delta \mathbf{GSK}_{uid}; \mathbf{GSK}^0_{uid} \mathbf{P}; \mathbf{Certificate}_{uid} \mathbf{P}$ . The CA setup algorithm is run by the CA. It takes no input other than the implicit security parameter  $\mathbf{P}$ . It generates the global master key  $\mathbf{GMK}$  of the system and the global public parameters  $\mathbf{GPP}$ . For each user  $uid$ , it generates the user's global public keys  $\delta \mathbf{GPK}_{uid}; \mathbf{GPK}^0_{uid} \mathbf{P}$ , the user's global secret keys  $\delta \mathbf{GSK}_{uid}; \mathbf{GSK}^0_{uid} \mathbf{P}$  and a certificate  $\mathbf{Certificate}_{uid} \mathbf{P}$  of the user.
- AASetup**  $\delta \mathbf{U}_{aid} \mathbf{P} ! \delta \mathbf{SK}_{aid}^0; \mathbf{PK}_{aid}^0; \mathbf{fVK}_{x_{aid}^0}; \mathbf{PK}_{x_{aid}^0} \mathbf{g}_{x_{aid}^0} \mathbf{2U}_{aid} \mathbf{P}$ . The attribute authority setup algorithm is run by each attribute authority. It takes the attribute universe  $\mathbf{U}_{aid}$  managed by the  $\mathbf{AA}_{aid}$  as input. It outputs a secret and public key pair  $\delta \mathbf{SK}_{aid}^0; \mathbf{PK}_{aid}^0$  of the  $\mathbf{AA}_{aid}$  and a set of version keys and public attribute keys  $\mathbf{fVK}_{x_{aid}^0}; \mathbf{PK}_{x_{aid}^0} \mathbf{g}_{x_{aid}^0} \mathbf{2U}_{aid}$  for all the attributes managed by the  $\mathbf{AA}_{aid}$ .

**Phase 2: Secret Key Generation by AAs.**

- SKeyGen**  $\delta \mathbf{GPP}; \mathbf{GPK}_{uid}; \mathbf{GPK}^0_{uid}; \mathbf{GSK}_{uid}; \mathbf{SK}_{aid}^0; \mathbf{fVK}_{x_{aid}^0}; \mathbf{PK}_{x_{aid}^0} \mathbf{g}_{x_{aid}^0} \mathbf{2Suid}; \mathbf{P} ! \mathbf{SK}_{uid;aid}$ . The secret key generation algorithm is run by each AA. It takes as inputs the global public parameters  $\mathbf{GPP}$ , the global public keys  $\delta \mathbf{GPK}_{uid}; \mathbf{GPK}^0_{uid} \mathbf{P}$  and one global secret key  $\mathbf{GSK}_{uid}$  of the user  $uid$ , the secret key  $\mathbf{SK}_{aid}^0$

- In this paper, we simply use the ciphertext to denote the encrypted content keys with CP-ABE. of the  $\mathbf{AA}_{aid}$ , a set of attributes  $\mathbf{S}_{uid;aid}$  that describes the user  $uid$  from the  $\mathbf{AA}_{aid}$  and its corresponding version keys  $\mathbf{fVK}_{x_{aid}^0} \mathbf{g}$  and public attribute keys  $\mathbf{fPK}_{x_{aid}^0} \mathbf{g}$ . It outputs a secret key  $\mathbf{SK}_{uid;aid}$  for the user  $uid$  which is used for decryption.

**Phase 3: Data Encryption by Owners.** Owners first encrypt the data  $m$  with content keys by using symmetric encryption

methods, then they encrypt the content keys by running the following encryption algorithm:

- Encrypt**  $\delta \mathbf{GPP}; \mathbf{fPK}_{aidk} \mathbf{g}_{aidk2IA}; \mathbf{P}; \mathbf{A} \mathbf{P} ! \mathbf{CT}$ . The encryption algorithm is run by the data owner to encrypt the content keys. It takes as inputs the global public parameters  $\mathbf{GPP}$ , a set of public keys  $\mathbf{fPK}_{aidk} \mathbf{g}_{aidk2IA}$  for all the AAs in the encryption set  $\mathbf{I}_A^3$ , the content key  $\mathbf{P}$  and an access policy  $\mathbf{A}$ .<sup>4</sup> The algorithm encrypts  $\mathbf{P}$  according to the access policy and outputs a ciphertext  $\mathbf{CT}$ . We will assume that the ciphertext implicitly contains the access policy  $\mathbf{A}$ .

**Phase 4: Data Decryption by Users.** Users first run the decryption algorithm to get the content keys, and use them to further decrypt the data.

- Decrypt**  $\delta \mathbf{CT}; \mathbf{GPK}_{uid}; \mathbf{GSK}^0_{uid}; \mathbf{fSK}_{uid;aidk} \mathbf{g}_{aidk2IA} \mathbf{P} ! \mathbf{P}$ . The decryption algorithm is run by users to decrypt the ciphertext. It takes as inputs the ciphertext  $\mathbf{CT}$  which contains an access policy  $\mathbf{A}$ , a global public key  $\mathbf{GPK}_{uid}$  and a global secret key  $\mathbf{GSK}^0_{uid}$  of the user  $uid$ , and a set of secret keys  $\mathbf{fSK}_{uid;aidk} \mathbf{g}_{aidk2IA}$  from all the involved AAs. If the attributes of the user  $uid$  satisfy the access policy  $\mathbf{A}$ , the algorithm will decrypt the ciphertext and return the content key  $\mathbf{P}$ .

**Phase 5: Attribute Revocation.** This phase contains three steps: Update Key Generation by AAs, Secret Key Update by Non-revoked Users<sup>5</sup> and Ciphertext Update by Server.

- UKeyGen**  $\delta \mathbf{SK}_{aid}^0; \mathbf{x}_{aid}^0; \mathbf{VK}_{x_{aid}^0} \mathbf{P} ! \delta \mathbf{VK}_{x_{aid}^0} \mathbf{P}; \mathbf{UK}_{s;x_{aid}^0} \mathbf{P}$ . The update key generation algorithm is run by the corresponding  $\mathbf{AA}_{aid}^0$  that manages the revoked attribute  $\mathbf{x}_{aid}^0$ . It takes as inputs the secret key  $\mathbf{SK}_{aid}^0$  of  $\mathbf{AA}_{aid}^0$ , the revoked attribute  $\mathbf{x}_{aid}^0$  and its current version key  $\mathbf{VK}_{x_{aid}^0}$ . It outputs a new version key  $\mathbf{VK}_{x_{aid}^0}$  and the update key  $\mathbf{UK}_{s;x_{aid}^0}$  (for secret key update) and the update key  $\mathbf{UK}_{c;x_{aid}^0}$  (for ciphertext update).<sup>f</sup>
- SKUpdate**  $\delta \mathbf{SK}_{uid;aid}^0; \mathbf{UK}_{s;x_{aid}^0} \mathbf{P} ! \mathbf{SK}_{uid;aid}^0$ . The secret key update algorithm is run by each non-revoked user  $uid$ . It takes as inputs the current secret key of the non-revoked user  $\mathbf{SK}_{uid;aid}^0$  and the update key  $\mathbf{UK}_{s;x_{aid}^0}$ . It

3. Note that not all the AAs are involved in the encryption. We use encryption set  $\mathbf{I}_A$  to denote the set of those AAs involved in the encryption.

4. The access policy is a LSSS structure  $\delta \mathbf{M}; \mathbf{P}$ , which is defined in the supplemental file available online.

5. We denote those users who possess the revoked attributes  $\mathbf{x}_{aid}^0$  but have not been revoked as the non-revoked users.

## 2.3 Security Model

In multi-authority cloud storage systems, we make the following assumptions:

- The CA is fully trusted in the system. It will not collude with any user, but it should be prevented from decrypting any ciphertexts by itself.
- Each AA is trusted but can be corrupted by the adversary.
- The server is curious but honest. It is curious about the content of the encrypted data or the received message, but will execute correctly the task assigned by each attribute authority.

We now describe the security model for our revocable multi-authority CP-ABE systems by the following game between a challenger and an adversary. Similar to the identity-based encryption schemes [15], the security model allows the adversary to query for any secret keys and update keys that cannot be used to decrypt the challenge ciphertext.

queries are made adaptively. Let  $S_A$  denote the set of all the attribute authorities. The security game is defined as follows.

Challenger generates the public keys by running the attribute authority setup algorithm and generates the secret keys by

running the secret key generation algorithm. For uncorrupted attribute authorities in  $S_A \setminus S_{A0}$ , the challenger only sends the public keys to the adversary. For corrupted authorities in  $S_{A0}$ , the challenger sends both the public keys and secret keys to the adversary. The adversary can also get the global public parameters.

### 3 OUR DATA ACCESS CONTROL SCHEME

In this section, we first give an overview of the challenges and techniques. Then, we propose the detailed construction of our access control scheme which consists of five phases: System Initialization, Key Generation, Data Encryption, Data Decryption and Attribute Revocation.

#### 3.1 Overview

To design the data access control scheme for multi-authority cloud storage systems, the main challenging issue is to construct the underlying Revocable Multi-authority CP-ABE protocol. In [6], Chase proposed a multi-authority CP-ABE protocol, however, it cannot be directly applied as the underlying techniques because of two main reasons: 1) **Security Issue**: Chase's multi-authority CP-ABE protocol allows the central authority to decrypt all the cipher texts, since it holds the master key of the system; 2) **Revocation Issue**: Chase's protocol does not support attribute revocation.

We propose a new revocable multi-authority CP-ABE protocol based on the single-authority CP-ABE proposed by Lewko and Waters in [16]. That is we extend it to multi-authority scenario and make it revocable. We apply the techniques in Chase's multi-authority CP-ABE protocol [6] to tie together the secret keys generated by different authorities for the same user and prevent the collusion attack. Specifically, we separate the functionality of the authority into a global certificate authority (CA) and multiple attribute authorities (AAs). The CA sets up the system and accepts the registration of users and AAs in the system.<sup>7</sup> It assigns a global user identity  $uid$  to each user and a global authority identity  $aid$  to each attribute authority in the system. Because the  $uid$  is globally unique in the system, secret

key issued by CA for the same  $uid$  cannot be used together for decryption. Also, because each AA is associated with an  $aid$ , every attribute is distinguish-able even though some AAs may issue the same attribute.

To deal with the security issue in [6], instead of using the system unique public key (generated by the unique master key) to encrypt data, our scheme requires all attribute authorities to generate their own public keys and uses them to encrypt data together with the global public parameters. This prevents the certificate authority in our scheme from decrypting the cipher texts.

To solve the attribute revocation problem, we assign a version number for each attribute. When an attribute revocation happens, only those components associated with the revoked attribute in secret keys and cipher texts need to be updated. When an attribute of a user is revoked from its corresponding AA, the AA generates a new version key for this revoked attribute and generates an update key. With the update key, all the users, except the revoked user, who hold the revoked attributes can update its secret key (Backward Security). By using the update key, the components associated with the revoked attribute in the ciphertext can also be updated to the current version. To improve the efficiency, we delegate the workload of

#### 3.2 AA Setup

Let  $S_{aid}$  denote the set of all attributes managed by each attribute authority  $AA_{aid}$ . It chooses three random numbers  $r_{aid}$ ,  $s_{aid}$ ,  $z_p$  as the authority secret key

$$SK_{aid} = (r_{aid}, s_{aid}, z_p);$$

$r_{aid}$  is used for data encryption,  $s_{aid}$  is used to distinguish attributes from different AAs and  $z_p$  is used

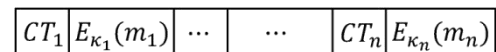


Fig. 2. Format of data on cloud server.

for attribute revocation. It also generates the public key  $PK_{aid}$  as

$$PK_{aid} = (e^{r_{aid}}, g^{s_{aid}}, g^{z_p});$$

For each attribute  $x_{aid} \in S_{aid}$ , the  $AA_{aid}$  generates a public attribute key as

$$PK_{1,x_{aid}} = (H(x_{aid})^{r_{aid}}, H(x_{aid})^{s_{aid}});$$

by implicitly choosing an attribute version key  $v_{x_{aid}}$ . All the public attribute keys  $PK_{1,x_{aid}}$  are published on the public bulletin board of the  $AA_{aid}$ , together with the public key  $PK_{aid}$  of the  $AA_{aid}$ .

#### 3.3 Secret Key Generation

Each user  $uid$  is required to authenticate itself to the  $AA_{aid}$  before it can be entitled some attributes from the  $AA_{aid}$ . The user submits its certificate  $Cert_{uid}$  to the  $AA_{aid}$ . The  $AA_{aid}$  then authenticates the user by using the verification key issued by the CA.

If it is a legal user, the  $AA_{aid}$  entitles a set of attributes



published ciphertexts are encrypted under attributes with old version. The ciphertext update algorithm in our protocol can update previously published cipher-texts into the latest attribute version, such that newly joined users can still decrypt previously published ciphertexts, if their attributes can satisfy access policies associated with ciphertexts. This guarantees the forward security.

However, when some AAs is corrupted by the adversary, the collusion resistance becomes more complicated. Specifically, the adversary may launch **Attribute Forge Attack**, defined as follows. Suppose a user  $uid_0$  possesses an attribute “ $x_{aid_0}$ ” from  $AA_{aid_0}$ , while the adversary does not hold the attribute “ $x_{aid_0}$ ” from  $AA_{aid_0}$ . The adversary attempts to forge (“clone”) the attribute “ $x_{aid_0}$ ” from the user  $uid_0$ ’s secret key by colluding with some other AAs.

In our scheme, the item  $g^{u_{uid} t_{uid;aid} - aid}$  in the secret key construction helps to resist this attack. When the adversary corrupts any AAs, he/she can get all the global secret key  $GSK_{uid}$ . Actually, the Forward Security and Backward Security are two basic requirements of attribute revocation. Now we prove that our scheme can achieve this two requirements as follows.

**Backward Security:** During the secret key update phase, the corresponding AA generates an update key for each non-revoked user. Because the update key is associated with the user’s global identity  $uid$ , the revoked user cannot use update keys of other non-revoked users to update its own secret key, even if it can compromise some non-revoked users. Moreover, suppose the revoked user can corrupt some other AAs (not the AA corresponding to the revoked attributes), the item  $H_{x_{aid}}^{v_{x_{aid} - aid} - aid}$  in the secret key can prevent users from updating their secret keys with update keys of other users, since  $_{aid}$  is only known by the  $AA_{aid}$  and kept secret to all the users. This guarantees the back-ward security.

**Forward Security:** After each attribute revocation operation, the version of the revoked attribute will be updated. When new users join the system, their secret keys are associated with attributes with the latest version. However, previously published ciphertexts are encrypted under attributes with old version. The ciphertext update algorithm in our protocol can update previously published ciphertexts into the latest attribute version, such that newly joined users can still decrypt previously published ciphertexts, if their attributes can satisfy access policies associated with ciphertexts. This guarantees the forward security. g

**Theorem .** Our access control scheme can resist the collusion attack, even when some AAs are corrupted by the adversary.

**Proof.** Users may collude and combine their attributes to decrypt the ciphertext, although they are not able to decrypt the ciphertext alone. Due to the random number  $t$  and the  $aid$  in the secret key, each component associated with the attribute in the secret key is distinguishable from

$GSK_{uid}$  for all the users in the system (because each AA has full knowledge on one of the user’s global secret keys  $GSK_{uid}$ ). Suppose all the  $K_{x_{aid};uid}$  in the secret key is constructed without this item. The adversary can successfully forge the attribute “ $x_{aid_0}$ ” as

In our scheme, the item  $g^{u_{uid} t_{uid;aid} - aid}$  in the secret key construction helps to resist this attack. When the adversary corrupts any AAs, he/she can get all the global secret key  $GSK_{uid}$  for all the users in the system (because each AA has full knowledge on one of the user’s global secret keys  $GSK_{uid}$ ). Suppose all the  $K_{x_{aid};uid}$  in the secret key is constructed without this item. The adversary can successfully forge the attribute “ $x_{aid_0}$ ” as

$$K_{x_{aid_0};uid_0} = \frac{1}{4} \delta K_{x_{aid_0};uid_0} \cdot GSK = GSK_{uid_0} ;$$

By adding the item  $g^{u_{uid} t_{uid;aid} - aid}$ , such attribute forge attack will be eliminated. g

**Privacy-Preserving Guarantee:** Although the CA holds the global master key GMK, it does not have any secret key issued from the AA. Without the knowledge of  $g^{-aid}$ , the CA cannot decrypt any ciphertexts in the system. Our scheme can also prevent the server from getting the content of the cloud data by using the proxy-encryption method.

## 5 PERFORMANCE ANALYSIS

In this section, we analyze the performance of our scheme by comparing with the Ruj’s DACC scheme [13] and our previous scheme in the conference version [14], in terms of storage overhead, communication cost and computation efficiency.

We conduct the comparison under the same security level. Let  $|p|$  be the element size in the  $G$ ;  $G_T$ ;  $Z_p$ . Suppose there are  $n_A$  authorities in the system and each attribute authority  $AA_{aid}$  manages  $n_{aid}$  attributes. Let  $n_U$  and  $n_O$  be the total number of users and owners in the system respectively. For a user  $uid$ , let  $n_{uid;aid_k} = \frac{1}{4} |S_{uid;aid_k}|$  denote the number of attributes that the user  $uid$  obtained from  $AA_{aid_k}$ . Let ‘ $v$ ’ be the total number of attributes in the ciphertext.

### 5.1 Storage Overhead

The storage overhead is one of the most significant issues of the access control scheme in cloud storage systems. Let  $n_{system}$  and  $n_P$  denote the total number attributes in the system and  $n_P$  denote the total number of attributes the user holds from all the  $s$  in the system. We compare the storage overhead on each entity in the system, as shown in Table 2.

#### 5.1 Storage Overhead

The storage overhead is one of the most significant issues of the access control scheme in cloud storage systems. Let  $n_{system}$  and  $n_P$  denote the total number attributes in the system and  $n_P$  denote the total number of attributes the user holds from all the  $s$  in the system. We compare the storage overhead on each entity in the system, as shown in Table 2.

1) **Storage Overhead on Each AA:** Each AA needs to store the information of all the attributes in its domain. Besides, in [14], each AA also needs to store the secret keys from all the owners, where the storage overhead on each AA is also linear to the total number of owners  $n_O$  in the system. In our scheme, besides the storage of attributes, each AA also needs to store a public key and a secret key for each user in the system. Thus, the storage overhead on each AA in our scheme is also linear to the number of users  $n_U$  in the system.

2) **Storage Overhead on Each Owner:** The public parameters contribute the main storage overhead on the owner. Besides the public parameters, in [13], owners are required to re-encrypt the ciphertexts and in [14] owners are required to generate the update information during the revocation, where the owner should also hold the encryption secret for every ciphertext in the system. This incurs a heavy storage overhead on the owner, especially when the number of ciphertext is large in cloud storage systems.

3) **Storage Overhead on Each User:** The storage overhead on each user in our scheme comes from the secret keys issued by all the AAs. However, in [13], the storage overhead on each user consists of both the secret keys issued by all the AAs and the ciphertext components that associated with the revoked attribute  $x$ , because when the ciphertext is re-encrypted, some of its components related to the revoked attributes should be sent to each non-revoked user who holds the revoked attributes. In [14], the user needs to hold multiple secret keys for multiple data owners, which means that the storage overhead on each user is also linear to the number of owners  $n_O$  in the system.

4) **Storage Overhead on Server:** The ciphertexts contribute the main storage overhead on the server (here we do not consider the encrypted data which are encrypted by the symmetric content keys).

## 5.2 Communication Cost

The communication cost of the normal access control is almost the same. Here, we only compare the communication cost of attribute revocation, as shown in Table 3. The communication cost of attribute revocation in [13] is linear to the number  $n$  in [14], the communication overhead is linear to the total number of attributes  $n_{c,aid}$  belongs to the AA<sub>aid</sub> in all the ciphertexts. It is not difficult to find that our scheme incurs much less communication cost during the attribute revocation.

## 5.3 Computation Efficiency

We implement our scheme and DACC scheme [13] on a Linux system with an Intel Core 2 Duo CPU at 3.16GHz and 4.00 GB RAM. The code uses the Pairing-Based Cryptography (PBC) library version 0.5.12 to implement the access control schemes. We use a symmetric elliptic curve  $\text{secp256k1}$ , where the base field size is 512-bit and the embedding degree is 2. The  $\text{secp256k1}$  curve has a 160-bit group order, which means  $p$  is a 160-bit length prime. All the simulation results are the mean of 20 trials.

We compare the computation efficiency of both encryption and decryption in two criteria: the number of authorities and the number of attributes per authority. Fig. 3a describes the comparison of encryption time versus the number of authorities, where the involved number of attributes per authority is set to be 10. Fig. 3c gives the encryption time comparison versus the number of attributes per authority, where the involved number of authority is set to be 10. It is easy to find that our scheme incurs less encryption time than DACC scheme in [13].

Fig. 3b shows the comparison of decryption time versus the number of authorities, where the number of attributes the user holds from each authority is set to be 10. Suppose the user has the same number of attributes from each authority, Fig. 3d describes the decryption time comparison versus the number of attributes the user holds from each authority. In Fig. 3d, the number of authority for the user is fixed to be 10. It is not difficult to see that our scheme incurs less decryption on the user than DACC scheme in [13].

Fig. 3e describes the time of ciphertext update/re-encryption versus the number of revoked attributes, and our scheme is more efficient than [13]. The ciphertext update/re-encryption contributes the main computation overhead of the attribute revocation. In our conference version [14], when an attribute is revoked from its corresponding authority AA<sub>aid</sub>, all the ciphertexts which are associated with any attributes from AA<sub>aid</sub> should be updated. In this paper, however, the attribute revocation method only requires the update of ciphertexts which are associated with the revoked attribute.

## 6 RELATED WORK

Ciphertext-Policy Attribute-Based Encryption (CP-ABE) [2]-[3] is a promising technique that is designed for access control of

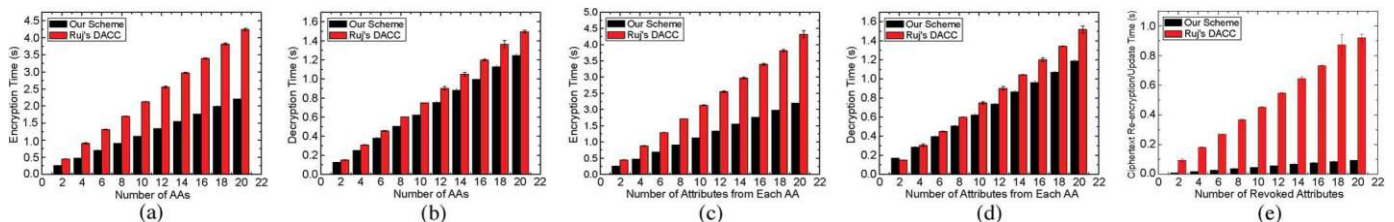


Fig. 3. Comparison of Computation Time. (a) Encryption. (b) Decryption. (c) Encryption. (d) Decryption. (e) Re-encryption.

encrypted data. There are two types of CP-ABE systems: single-authority CP-ABE [2], [3], [4], [5] where all attributes are managed by a single authority, and multi-authority CP-ABE [6], [7], [8] where attributes are from different domains and managed by different authorities. Multi-authority CP-ABE is more appropriate for the access control of cloud storage systems, as users may hold attributes issued by multiple authorities and the data owners may share the data using access policy defined over attributes from different authorities. However, due to the attribute revocation problem, these multi-authority CP-ABE schemes cannot be directly applied to data access control for such multi-authority cloud storage systems.

To achieve revocation on attribute level, some re-encryption-based attribute revocation schemes [9], [11] are proposed by relying on a trusted server. We know that the cloud server cannot be fully trusted by data owners, thus traditional attribute revocation methods are no longer suitable for cloud storage systems.

Ruj, Nayak and Ivan proposed a DACC scheme [13], where an attribute revocation method is presented for the Lewko and Waters' decentralized ABE scheme [8]. Their attribute revocation method does not require a fully trusted server. But, it incurs a heavy communication cost since it requires the data owner to transmit a new ciphertext component to every non-revoked user.

## 7 CONCLUSION

In this paper, we proposed a revocable multi-authority CP-ABE scheme that can support efficient attribute revocation. Then, we constructed an effective data access control scheme for multi-authority cloud storage systems. We also proved that our scheme was provable secure in the random oracle model. The revocable multi-authority CP-ABE is a promising technique, which can be applied in any remote storage systems and online social networks etc.

## REFERENCES

- [1] P. Mell and T. Grance, "The NIST Definition of Cloud Computing," National Institute of Standards and Technology, Gaithersburg, MD, USA, Tech. Rep., 2009.
- [2] J. Bethencourt, A. Sahai, and B. Waters, "Ciphertext-Policy Attribute-Based Encryption," in Proc. IEEE Symp. Security and privacy (S&P'07), 2007, pp. 321-334.
- [3] B. Waters, "Ciphertext-Policy Attribute-Based Encryption: An Expressive, Efficient, and Provably Secure Realization," in Proc. 4th Int'l Conf. Practice and Theory in Public Key Cryptography (PKC'11), 2011, pp. 53-70.
- [4] V. Goyal, A. Jain, O. Pandey, and A. Sahai, "Bounded Ciphertext Wide Web, Journal of Combinatorial Optimization, etc. He is the General Policy Attribute Based Encryption," in Proc. 35th Int'l Colloquium Chair of ACM MobiHoc 2008, TPC Co-Chair of IEEE MASS 2009, Area-on Automata, Languages, and Programming (ICALP'08), 2008, pp. 579-591.
- [5] A.B. Lewko, T. Okamoto, A. Sahai, K. Takashima, and B. Waters, "Fully Secure Functional Encryption: Attribute-Based Encryption and (Hierarchical) Inner Product Encryption," in Proc. Advances in Cryptology-EUROCRYPT'10, 2010, pp. 62-91.
- [6] M. Chase, "Multi-Authority Attribute Based Encryption," in Proc. 4th Theory of Cryptography Conf. Theory of Cryptography (TCC'07), 2007, pp. 515-534.
- [7] M. Chase and S.S.M. Chow, "Improving Privacy and Security in Multi-Authority Attribute-Based Encryption," in Proc. 16th ACM Conf. Computer and Comm. Security (CCS'09), 2009, pp. 121-130.
- [8] A.B. Lewko and B. Waters, "Decentralizing Attribute-Based Encryption," in Proc. Advances in Cryptology-EUROCRYPT'11, 2011, pp. 568-588.
- [9] S. Yu, C. Wang, K. Ren, and W. Lou, "Attribute Based Data Sharing with Attribute Revocation," in Proc. 5th ACM Symp. Information, Computer and Comm. Security (ASIACCS'10), 2010, pp. 261-270.
- [10] M. Li, S. Yu, Y. Zheng, K. Ren, and W. Lou, "Scalable and Secure Sharing of Personal Health Records in Cloud Computing Using Attribute-Based Encryption," IEEE Trans. Parallel Distributed Systems, vol. 24, no. 1, pp. 131-143, Jan. 2013.
- [11] J. Hur and D.K. Noh, "Attribute-Based Access Control with Efficient Revocation in Data Outsourcing Systems," IEEE Trans. Parallel Distributed Systems, vol. 22, no. 7, pp. 1214-1221, July 2011.
- [12] S. Jahid, P. Mittal, and N. Borisov, "Easier: Encryption-Based Access Control in Social Networks with Efficient Revocation," in Proc. 6th ACM Symp. Information, Computer and Comm. Security (ASIACCS'11), 2011, pp. 411-415.
- [13] S. Ruj, A. Nayak, and I. Stojmenovic, "DACC: Distributed Access Control in Clouds," in Proc. 10th IEEE Int'l Conf. TrustCom, 2011, pp. 91-98.
- [14] K. Yang and X. Jia, "Attribute-Based Access Control for Multi-Authority Systems in Cloud Storage," in Proc. 32th IEEE Int'l Conf. Distributed Computing Systems (ICDCS'12), 2012, pp. 1-10.
- [15] D. Boneh and M.K. Franklin, "Identity-Based Encryption from the Weil Pairing," in Proc. 21st Ann. Int'l Cryptology Conf.: Advances in Cryptology - CRYPTO'01, 2001, pp. 213-229.
- [16] A.B. Lewko and B. Waters, "New Proof Methods for Attribute-Based Encryption: Achieving Full Security through Selective Techniques," in Proc. 32st Ann. Int'l Cryptology Conf.: Advances in Cryptology - CRYPTO'12, 2012, pp. 180-198.

of ciphertexts which contain the revoked

# Oruta: Privacy-Preserving Public Auditing for Shared Data in the Cloud

***Byreddy Madhavi***

CSE Department

saimadhavireddy11@gmail.com

Malla Reddy College of Engineering

***Dr.V.Bhoopathy***

*Professor, CSE Department*

v.bhoopathy@gmail.com

Malla Reddy College of Engineering

## **Abstract:**

With cloud storage services, it is commonplace for data to be not only stored in the cloud, but also shared across multiple users. However, public auditing for such shared data— while preserving identity privacy — remains to be an open challenge. In this paper, we propose the first privacy-preserving mechanism that allows public auditing on shared data stored in the cloud. In particular, we exploit ring signatures to compute the verification information needed to audit the integrity of shared data. With our mechanism, the identity of the signer on each block in shared data is kept private from a third party auditor (TPA), who is still able to verify the integrity of shared data without retrieving the entire file. Our experimental results demonstrate the effectiveness and efficiency of our proposed mechanism when auditing shared data.

**Key Words:** Provable Data Possession, Third party Auditor, Hybrid Cloud

## **Existing System:**

The first provable data possession (PDP) mechanism [2] to perform public auditing is designed to check the correctness of data stored in an un trusted server, without retrieving the entire data. Moving a step forward, Wang *et al.* [3] (referred to as WWRL in this paper) is designed to construct a public auditing mechanism for cloud data, so that during public auditing, the content of private data belonging to a personal user is not disclosed to the third party auditor.

## **Disadvantage:**

Data is not in an encrypted format.

## **Proposed System:**

In this paper, we only consider how to audit the integrity of shared data in the cloud with *static groups*. It means the group is pre-defined before shared data is created in the cloud and the membership of users in the group is not changed during data sharing. The original user is responsible for deciding who is able to share her data before outsourcing data to the cloud. Another interesting problem is how to audit the integrity of shared data in the cloud with *dynamic groups* — a new user can be added

into the group and an existing group member can be revoked during data sharing — while still preserving identity privacy.

### **Advantage:**

Here we proposed the secured system and data owner can decide whether the user can access the system or not.

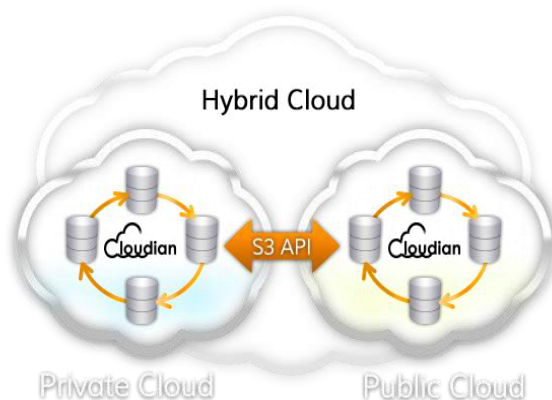
### **PROBLEM STATEMENT:**

In our model, privacy is accomplished by allowing the parties to upload their data in multi clouds and data is split into multiple parts so it gives more protection

### **Scope:**

We are going to raise the privacy level of the data owner and the confidentiality of the data in a better way through the multiple cloud environment.

### **Architecture:**



### **Modules :**

1. **Owner**
2. **Third Party Auditor**

### 3. **User**

### 4. **Data Sharing**

### **Modules Description**

#### **Owner Registration:**

In this module an owner has to upload its files in a cloud server, he/she should register first. Then only he/she can be able to do it. For that he needs to fill the details in the registration form. These details are maintained in a database.

#### **Owner Login:**

In this module, any of the above mentioned person have to login, they should login by giving their emailid and password .

#### **User Registration:**

In this module if a user wants to access the data which is stored in a cloud, he/she should register their details first. These details are maintained in a Database.

#### **User Login:**

If the user is an authorized user, he/she can download the file by using file id which has been stored by data owner when it was uploading.

#### **ThirdPartyAuditor Registration:**

In this module , if a third party auditor TPA(maintainer of clouds) wants to

do some cloud offer , they should register first. Here we are doing like, this system allows only three cloud service providers.

### ThirdPartyAuditor Login:

After third party auditor gets logged in, He/ She can see how many data owners have uploaded their files into the cloud. Here we are providing three tpa for maintaining three different clouds.

### Data Sharing:

we only consider how to audit the integrity of shared data in the cloud with *static groups*. It means the group is pre-defined before shared data is created in the cloud and the membership of users in the group is not changed during data sharing. The original user is responsible for deciding who is able to share her data before outsourcing data to the cloud. Another interesting problem is how to audit the integrity of shared data in the cloud with *dynamic groups* — a new user can be added into the group and an existing group member can be revoked during data sharing — while still preserving identity privacy.

### Proposed System:

To enable the TPA efficiently and securely verify shared data for a group of users, Oruta should be designed to achieve following properties: (1) **Public Auditing:** The third party auditor is able to verify the

integrity of shared data for a group of users without retrieving the entire data. (2) **Correctness:** The third party auditor is able to correctly detect whether there is any corrupted block in shared data. (3) **Enforceability:** Only a user in the group can generate valid verification information on shared data. (4) **Identity Privacy:** During auditing, the TPA cannot distinguish the identity of the signer on each block in shared data.

### System Configuration:-

#### H/W System Configuration:-

|       |                    |             |
|-------|--------------------|-------------|
| III   | Processor          | - Pentium – |
|       | Speed              | - 1.1 GHz   |
|       | RAM                | - 256 MB    |
| (min) | Hard Disk          | - 20 GB     |
|       | Floppy Drive       | - 1.44 MB   |
|       | Key Board          | - Standard  |
|       | Windows Keyboard   |             |
|       | Mouse              | - Two or    |
|       | Three Button Mouse |             |

**Monitor - SVGA**

❖ Scripts :  
JavaScript.

❖ Database :  
My sql

❖ Database Connectivity :  
JDBC.

**S/W System Configuration:-**

❖ Operating System  
:Windows95/98/2000/XP

❖ Front End :  
Swings & AWT

# ANDROID SECURITY

**A.Sindhu,**

Sindhurajender967@gmail.com  
Malla Reddy College of Engineering

**Dr.P.Mani Kandan**

Malla Reddy Engineering College for Women,  
Mani.p.mk@gmail.com

## ABSTRACT

In this application which is useful for the user when he is in some problem or needs any help. When the user opens this application, he can see a HELP button. Also he can store a message and 3 contact numbers. When the user is in some difficulty or needs any help,

### 1.1 Purpose of the Project

When the user is in some difficulty or needs any help, he needs to simply open the app and click on the “HELP” button. This application sends the message to those contact numbers which he has stored.

### 1.2 Scope of the Project

Android application which is useful for the user when she is in some problem or needs any help. When the user opens this application, he can see a HELP button. Also he can store a message and 3 contact numbers. When the user is in some difficulty or needs any help, he needs to simply open the app and click on the “HELP” button.

### 1.3 Features of the Project

To reduce user effort and solve problems inherent to the cellular phones small screen, several functions are provided on the cellular viewer.

- Supports multiple connections at the same time.
- Different work modes: "view only" and "full control".
- Different display modes: "windowed", "full screen", and "scaled".
- Runs as a service on the NT systems.
- Works through the firewalls and supports DHCP.
- Supports high screen resolutions and color depths.

## 2 SYSTEM ANALYSIS AND DESCRIPTION

### 2.1 Existing System:

In the existing system, the user has to write the message content and select the contacts and only then he can send the message but what if the user do not have that much time or unable to do it.

## **2.2 Proposed System:**

In this proposed system, the user writes the message content and also selects the contacts to which the message has to be sent and save it. So, when he is in some danger by just opening the app and pressing the HELP button, the message stored will be sent to those numbers he has added in this application. So that he can receive the help in correct time.

### **2.2.1 Advantages:**

- Through this web application we can save the time, we can efficiently get succeed in implementing all the requests within short span of time.

### **2.3.1 Modules and Functionalities**

#### **Modules Are.....**

Help Button

Adding Contacts

Messages

## **2.4 Feasibility Study:**

The next step in analysis is to verify the feasibility of the proposed system. "All projects are feasible given unlimited resources and infinite time". But in reality both resources and time are scarce. Project should confirm to time bounce and should be optimal in their consumption of resources. This place a constant is approval of any project.

Feasibility they are 3 types:

- Technical feasibility
- Operational feasibility
- Economic feasibility

### **2.4.1 Technical Feasibility:**

To determine whether the proposed system is technically feasible, we should take into consideration the technical issues involved behind the system.

This Application uses the web technologies, which is rampantly employed these days worldwide. The world without the web is incomprehensible today. That goes to proposed system is technically feasible.

### **2.4.2 Operational Feasibility:**

To determine the operational feasibility of the system we should take into consideration the awareness level of the users. This system is operational feasible since the users are familiar with the technologies and hence there is no need to gear up the personnel to use system. Also the system is very friendly and to use.

### **2.4.3 Economic Feasibility:**

To decide whether a project is economically feasible, we have to consider various factors as:

- Cost benefit analysis
- Long-term returns
- Maintenance costs

It requires average computing capabilities and access to internet, which are very basic requirements and can be afforded by any organization hence it

doesn't incur additional economic overheads, which renders the system economically feasible.

### 3. IMPLEMENTATION

### 3.1 JAVA SCRIPT

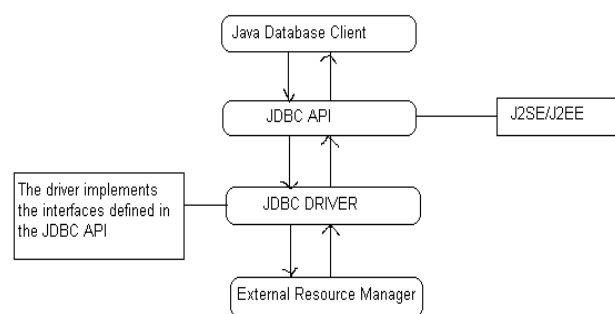
Java script originally supported by Netscape navigator is the most popular web scripting language today. Java script lets you embedded programs right in your web pages and run these programs using the web browser. You place these programs in a <SCRIPT> element, usually with in the <HEAD> element. If you want the script to write directly to the web page, place it in the <BODY> element. .

### 3.2 .JDBC DRIVERS:

The JDBC API only defines interfaces for objects used for performing various database-related tasks like opening and closing connections, executing SQL commands, and retrieving the results. We all write our programs to interfaces and not implementations. Either the resource manager vendor or a third party provides the implementation classes for the standard JDBC interfaces. These software implementations are called JDBC drivers. JDBC drivers transform the standard JDBC calls to the external resource manager-specific API calls. The diagram below depicts how a database client Written in java accesses an external resource manager using the JDBC API and

JDBC

driver:



**Fig: 3.2**

## jdbc drivers

**TYPE1:**

Type1 JDBC drivers implement the JDBC API on top of a lower level API like ODBC. These drivers are not generally portable because of the independency on native libraries. These drivers translate the JDBC calls to ODBC calls and ODBC sends the request to external data source using native library calls. The JDBC-ODBC driver that comes with the software distribution for J2SE is an example of a type1 driver.

**TYPE2:**

Type2 drivers are written in mixture of java and native code. Type2 drivers use vendors specific native APIs for accessing the data source. These drivers transform the JDBC calls to vendor specific calls using the vendor's native library.

These drivers are also not portable like type1 drivers because of the dependency on native code.

### **TYPE3:**

Type3 drivers use an intermediate middleware server for accessing the external data sources. The calls to the middleware server are database independent. However, the middleware server makes vendor specific native calls for accessing the data source. In this case, the driver is purely written in java.

### **TYPE4:**

Type4 drivers are written in pure java and implement the JDBC interfaces and translate the JDBC specific calls to vendor specific access calls. They implement the data transfer and network protocol for the target resource manager. Most of the leading database vendors provide type4 drivers for accessing their database servers.

### **3.2.1 DRIVER MANAGER AND DRIVER:**

The java.sql package defines an interface called Java.sql.Driver that makes to be implemented by all the JDBC drivers and a class called java.sql.DriverManager that acts as the interface to the database clients for performing tasks like connecting to external resource managers, and setting log streams. When a JDBC

client requests the Driver Manager to make a connection to an external resource manager, it delegates the task to an appropriate driver class implemented by the JDBC driver provided either by the resource manager vendor or a third party.

### **3.2.2 JAVA.SQL.DRIVERMANAGER:**

The primary task of the class driver manager is to manage the various JDBC drivers register. It also provides methods for:

- Getting connections to the databases.
- Managing JDBC logs

### **3.2.3 MANAGING DRIVERS:**

JDBC clients specify the JDBC URL, when they request a connection. The driver manager can find a driver that matches the request URL from the list of register drivers and delegate the connection request to that driver if it finds a match JDBC URLs normally take the following format:

**<Protocol>:<sub-protocol>:<resource>**

The protocol is always jdbc and the sub-protocol and resource depend on the type of resource manager. The URL for postgresSQL is in the format:

**Jdbc: postgres ://< host> :< port>/<database>**

Here host is the host address on which post master is running and database is the name

of the database to which the client wishes to connect.

### 3.3 JAVA SERVER PAGES (JSP)

#### 3.3.1 INTRODUCTION:

Java Server Pages (JSP) technology enables you to mix regular, static HTML with dynamically generated content. You simply write the regular HTML in the normal manner, using familiar Web-page-building tools. You then enclose the code for the dynamic parts in special tags, most of which start with `<%` and end with `%>`.

#### 3.3.2 THE NEED FOR JSP:

Servlets are indeed useful, and JSP by no means makes them obsolete. However,

- It is hard to write and maintain the HTML.
- You cannot use standard HTML tools.
- The HTML is inaccessible to non-Java developers.

#### 3.3.3 BENEFITS OF JSP:

JSP provides the following benefits over servlets alone:

- It is easier to write and maintain the HTML: In this no extra backslashes, no double quotes, and no lurking Java syntax.

- You can use standard Web-site development tools:

We use Macromedia Dreamweaver for most of the JSP pages. Even HTML tools that know nothing about JSP can be used because they simply ignore the JSP tags.

- You can divide up your development team:

The Java programmers can work on the dynamic code. The Web developers can concatenate on the representation layer. On large projects, this division is very important. Depending on the size of your team and the complexity of your project, you can enforce a weaker or stronger separation between the static HTML and the dynamic content.

#### 3.3.4 TYPES OF JSP SCRIPTING ELEMENTS:

JSP scripting elements allow you to insert Java code into the servlet that will be generated from the JSP page. There are three forms:

- **Expressions** of the form `<%=Java Expression %>`, which are evaluated and inserted into the servlet's output.
- **Scripts** of the form `<%Java code %>`, which are inserted

into the servlet's `_jspService` method (called by service).

- **Declarations** of the form `<%! Field/Method Declaration %>`, which are inserted into the body of the servlet class, outside any existing methods.

### 3.3.5 PREDEFINED VARIABLES:

To simplify expressions we can use a number of predefined variables (or “implicit objects”). The specialty of these variables is that, the system simply tells what names it will use for the local variables in `_jspService`. The most important ones of these are

## 4. TESTING

### 4.1 SOFTWARE TESTING

Software testing is a critical element of software quality assurance and represents the ultimate review of specification, design and code generation.

#### 4.1.1 TESTING OBJECTIVES

- To ensure that during operation the system will perform as per specification.
- To make sure that system meets the user requirements during operation
- To make sure that during the operation, incorrect input, processing and output will be detected

- To see that when correct inputs are fed to the system the outputs are correct
- To verify that the control incorporated in the same system as intended
- Testing is a process of executing a program with the intent of finding an error

A good test case is one that has a high probability of finding an as yet undiscovered error

The software developed has been tested successfully using the following testing strategies and any errors that are encountered are corrected and again the part of the program or the procedure or function is put to testing until all the errors are removed. A successful test is one that uncovers an as yet undiscovered error.

Note that the result of the system testing will prove that the system is working correctly. It will give confidence to system designer; users of the system prevent frustration during implementation process etc.

### 4.2 TEST CASE DESIGN:

#### 4.2.1 White box testing

White box testing is a testing case design method that uses the control structure of the procedure design to derive test cases. All independent paths in a

module are exercised at least once, all logical decisions are exercised at once, execute all loops at boundaries and within their operational bounds exercise internal data structure to ensure their validity. Here the customer is given three chances to enter a valid choice out of the given menu. After which the control exits the current menu.

**4.2.2 Black Box Testing**

Black Box Testing attempts to find errors in following areas or categories, incorrect or missing functions, interface error, errors in data structures, performance error and initialization and termination error. Here all the input data must match the data type to become a valid entry.

The following are the different tests at various levels:

**4.2.3 Unit Testing:**

Unit testing is essentially for the verification of the code produced during the coding phase and the goal is test the internal logic of the module/program. In the Generic code project, the unit testing is done during coding phase of data entry forms whether the functions are working properly or not. In this phase all the drivers are tested they are rightly connected or not.

**4.2.4 Integration Testing:**

All the tested modules are combined into sub systems, which are then tested. The goal is to see if the modules are properly integrated, and the emphasis being on the testing interfaces between the modules. In the generic code integration testing is done mainly on table creation module and insertion module.

**5. CONCLUSION**

In this project to use which is useful for the user when he is in some problem or needs any help. When the user opens this application, he can see a HELP button. Also he can store a message and 3 contact numbers. When the user is in some difficulty or needs any help button. So when the user opens this application, can see a HELP button. Click that button to send sms to register user.

**1.Advanced Java Programming**  
- Dietel and Dietel

**2.Mastering JAVA 2**  
- John Zukowski

**3.Java Server Programming**  
- Apress

**4.Software Engineering**  
- Roger S Pressman

**5.Análisis & Design of Information  
Systems – Senn**

|                             |                                     |
|-----------------------------|-------------------------------------|
| <b>IDE</b>                  | <b>-Eclipse with Adt<br/>Plugin</b> |
| <b>Operating<br/>System</b> | <b>-Windows</b>                     |
| <b>SDK</b>                  | <b>-Android Sdk 2.3</b>             |

Websites:

- 1.[www.eci.gov.in](http://www.eci.gov.in)
- 2.[www.google.com](http://www.google.com)
- 3.[www.apeci.com](http://www.apeci.com)

We will be using the information  
collected from websites for android

- 1.<http://developer.android.com/index.html>
- 2.<http://www.android-trainer.com/>
- 3.<http://stackoverflow.com/>
- 4.<http://www.google.co.in/>

## **Is Alcohol Affect Students Performance: Searching and Predicting using Data Mining Algorithms**

***P.Pravalika,***

CSE Department

Malla Reddy College of Engineering

Pravalikapinky10@gmail.com

***Dr.Raja Sekar***

Professor, CSE Department

St.Peter College of Engineering,

rajasekaratr@gmail.com

### **ABSTRACT**

*Chronic heavy drinking and alcoholism can have serious repercussions for the learning, and memory. Excessive drinking among college students is associated with a variety of negative consequences that include alcohol poisoning; fatal and nonfatal injuries; blackouts; violence, academic failure; including sexually transmitted diseases, rape and assault; unintended pregnancy; property damage; including HIV/AIDS; criminal consequences and vocational that could jeopardize future job prospects. The present research intends to the student achievement in Higher Education using Data Mining techniques. This real-world data was collected by using performance reports and questionnaires which was collected and analyzed by MCA department, VBS Purvanchal University, India. In this experimental dataset used data set about MCA student on their courses which holds 450 instances. Four Decisions Tree algorithms (BFTree, J48, RepTree and Simple Cart) are applied in this work. The results showed that BFTree algorithm mostly proper to classify and predict student's whose performance is excellent and who's poor during studying the subjects.*

***Keywords: Data Mining, Questionnaires, BFTree, Drinking and Alcoholism, Chronic.***

## INTRODUCTION

Data mining sometimes called data or knowledge discovery in databases is the extraction of hidden predictive information from large databases the process of analyzing data from different summarizing and perspectives it into useful information. Technically, data mining is the process of finding correlations or patterns among dozens of fields in large relational databases [1].

Educational data mining (also referred to as —EDMI) is an area of data mining defined as the scientific inquiry centered on the development of methods for making discoveries within the unique kinds of data that come from educational dataset, and using those methods to better understand students and the settings which they learn in. Prediction has two key uses within educational data mining. In this model it is important for prediction, giving information about to predict student educational outcomes (cf. Romero et al, 2008) without predicting mediating factors or intermediate first. In a second type of usage, prediction methods are used in order to predict what the output value would be in contexts where it is not desirable to directly obtain a label for that construct [2]. This paper presents a new model that enhances the Decision Tree accuracy in identifying student's performance. Four Decision Tree algorithms BFTree, J48, RepTree and Simple Cart are applied. The data set consists of two Comma Separated Values (CSV) files taken from UCI Machine Learning Repository for Students Alcohol Consumption of courses MCA students. The source data set files contained 450 instances with 31 attributes. WEKA 3.6.9 tool is used to implement Decision Trees. In our studies 10-fold cross validation method was used to measure the unbiased estimate of these prediction models with 66% of the tested data.

The organization of this paper is: section two viewed the related works and listed all the models of implementing the algorithms of data mining with education. Section three explained the concept of Educational Data Mining (EDM) briefly. Section four listed and explained the decision trees BFTree, J48, RepTree and Simple Cart which are implemented later in this model. Section five explained the machine learning

tool WEKA. Section six listed the decision trees model of steps and results of implementing. The final section concludes the extraction from the whole work.

## BACKGROUND

Pandey and Pal [3] conducted a study on new comer students will performer or not on the basis of student performance selecting 600 students from different colleges of Dr.R.M.L. Awadh University, Faizabad, India.

They applied Bayes Classification on Category, Language and background qualification, in conclusion they found that whether newcomers will perform or not.

Romero, Cristóbal, et al. [4] compared different data mining methods and techniques for classifying students based on the final marks obtained in their respective courses. They used real data from seven model courses with Cordoba University students. In their applied method they found that a classifier model is sufficient for educational use has to be both accurate and comprehensible for teachers for making decisions.

Galit [5] perform a case study that uses student's data to predict and analyze their learning behavior and to warn students at risk before their final exams.

Osmanbegović, Edin, and Mirza Suljić [6] from University of Tuzla perform data mining techniques and methods for comparing the prediction of students' success, applying the data collected from the surveys conducted during the summer semester by the Faculty of Economics, academic year 2010-2011, among first year students and investigated the result of students' achieved from high school and from the entrance exam, and effect on success.

Pandey and Pal [7] perform a study on the student performance by selecting 60 students from a degree college of Dr. R. M. L. Awadh University, Faizabad, India. They applied association rule for find the interest of student in opting class teaching language.

Cortez, Paulo, and Alice Maria Gonçalves Silva [8] used data mining techniques to predict the student's achievement of secondary school using real-world data. The two core classes (Mathematics and Portuguese) were modeled. Four Data Mining models (i.e. Neural Networks,

Decision Trees, Support Vector Machines and Random Forest) and three input selections (e.g. without previous grades and with) were tested. The results shows predictive accuracy can be achieved, although student achievement also depend other relevant.

Han and Kamber [9] describes data mining software that allow the users to analyze data from different dimensions, categorize it and summarize the relationships which are identified during the mining process.

Kumar, S. Anupama, and M. N. Vijayalakshmi [10] perform a study on student's internal assessment data to predict their performance in the final exam. They used C4.5 decision tree algorithm. The accuracy of the algorithm is compared with ID3 algorithm and found to be more efficient in form of the accurately predicting the time taken to derive the tree and outcome of the student.

Bhardwaj and Pal [11] perform a study on BCA (Bachelor of Computer Application) course of Dr. R. M. L. Awadh University, Faizabad, India. The student performance based by selecting 300 students from 5 different degree college. They used classification method on 17 attribute, and found that the factors like living location, medium of teaching, students\_ grade in senior secondary exam, students other habit, mother's qualification, student's family status were highly correlated with the student academic performance and family annual income.

Khan [12] conducted a performance study on 400 students comprising 200 boys and 200 girls selected from the senior secondary school of Aligarh Muslim University, Aligarh, India the aim of this research to analyze the rate of success of students in higher secondary in science group. Cluster sampling technique was used to analyze this study. He found that girls who had good income, education, occupation, wealth and place of residence, got greater achievement on the other hand the boys with low living status had relatively greater academic gain also.

### **EDUCATIONAL DATA MINING (EDM)**

Traditional data mining methods often differ from methods from the Educational data mining. Methods from the psychometrics literature are often integrated with methods from the machine. In recent years Educational data mining has

emerged as an independent research area, culminating in 2008 with the establishment of the annual International Conference on Educational Data Mining, and the Journal of Educational Data Mining [13].

There are popular methods within educational data mining. These methods are following categories: prediction, clustering, relationship mining, discovery with models, and distillation of data for human judgment. The prediction, clustering, relationship mining are largely acknowledged to be universal across types of data mining and the discovery with models, and distillation of data for human judgment categories achieve particular prominence within educational data mining.

#### **i. Prediction**

For intelligent decision making databases are rich with hidden information that can be used. Prediction is one of the forms of data analysis that can be used to predict future data trends or to describing important data classes from extract models. This type of analysis can help provide us with a better understanding of the large data. There are three types of prediction: density estimation, regression, and classification. In the classification method, the predicted variable is a binary variable or categorical variable. Some of the popular classification methods include support vector machines, decision trees, and logistic regression. In regression, the predicted variable is a continuous variable [14]. Some popular regression methods within educational data mining include neural networks, support vector machine regression and linear regression. In density estimation, the predicted variable is a probability density function. Density estimators can be based on Gaussian functions and a variety of kernel functions. For each type of prediction, the input variables can be either continuous or categorical; different prediction methods are more effective, depending on the type of input variables used.

#### **ii. Clustering**

The process of clustering is grouping the data into clusters or classes, so that objects within clusters have high similarity in comparison to one another but are very dissimilar to objects in other clusters. Clustering is particularly useful in cases where the most common categories within

Clusters can be created at schools or students or clustered together to search behavioral patterns or to investigate differences and similarities in between schools. In clustering algorithm each and every data point must contain exactly one cluster or it can postulate that it may contain no cluster or one cluster at that point [15].

### iii. Relationship mining

The relationship mining has a basic goal to find the relationships between variables, in a data set which contain a large number of variables. This take place to find out which variables may take place most strongly associated with a single variable of particular interest or to discover which relationships in between any two the variables are strongest.

There are four types of relationship mining: causal data mining, association rule mining, sequential pattern mining and correlation mining. In association rule mining, the goal is to find if-then rules of the form that another variable will generally have a specific value if some set of variable values is found. Relationships found through relationship mining must satisfy two criteria: interestingness and statistical significance. Statistical significance is generally assessed through F-tests or some standard statistical tests. It is necessary to control for obtaining relationships through chance. One method for doing this is the Bonferroni adjustment or, such as to use post-hoc statistical methods. An alternate method is using Monte Carlo methods to assess the overall probability of the pattern of results found. This method assesses how likely it is that the overall pattern of results arose due to chance.

In very large data sets, it is difficult to define the significant relationships because of the numerous of relationship may be found. For our interest measures attempt to determine which relationship are the well-supported by the data and most separated and distinctive, in some cases also try to prune overly similar findings [16]. There are a range of wide variety of interestingness measures, which including cosine, support, leverage, conviction, lift, coverage, correlation and confidence. Within educational data mining (Merceron & Yacef, 2008) lift and cosine may or may not be particularly relevant in some

### iv. Discovery with Models

This model is used as a component in another analysis, such as prediction or relationship mining. In discovery with a model, a model of a phenomenon is developed via prediction, clustering, or knowledge engineering.

For predicting a new variable predictor variables created a model's of predictions in the prediction case. The current knowledge component which learned within online learning have generally depended on assessments of the probability that the student analyses of complex constructs in gaming. These assessments of student generally expressed as a mapping between exercises within the learning software and knowledge components, knowledge have in turn depended on models of the knowledge components in a domain,. In the relationships between the created model's predictions and additional variables are studied the relationship mining case. This

can enable a researcher to study a wide variety of observable constructs and the relationship between a complex latent construct. For instance, how much a student would game the system used predictions of gaming the system whether state or trait factors were better predictors [17]. Generalization across contexts in this fashion relies upon appropriate validation that the model accurately generalizes.

### v. Distillation of Data for Human Judgment

Within educational data mining is the distillation of data for human judgment. The methods in this area of educational data mining are information visualization methods – however, within EDM the visualizations is most commonly used are often different than those most often used for other information visualization problems, owing to the embedded within that structure, and specific structure meaning, often present in educational data. Data is distilled for human judgment in educational data mining for two key purposes: classification and identification. When data is distilled for identification, which are nonetheless difficult to formally

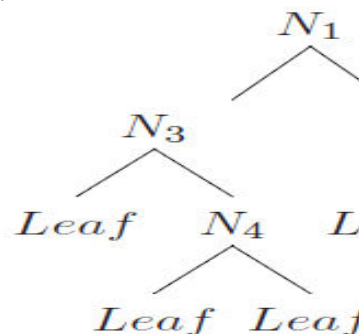
express data is displayed in ways that enable a human being to easily identify well-known patterns. For example, one of the classic educational data mining by which visualization is the learning curve, which displays the number of opportunities to practice a skill on the X axis, and performance on the Y axis. In this case, sub-sections of a data set are displayed in labeled by human coders, and visual or text format. These are the labels by which is used for the basis of the development of a predictor. This approach has been shown to speed up the development of prediction models of complex phenomena such as gaming the system by approximately 40 times, relative to prior methods for collecting the necessary data (Baker & de Carvalho, 2008).

## METHODOLOGY

In this paper we use the following data mining methods.

### i. BFTree

Best-first trees are constructed in a divide-and-conquer fashion similar to standard depth-first decision trees [18]. The basic idea behind a best-first tree is built is as follows. Firstly, select an attribute to put at the root node and make some branches for this attribute using some criteria. Then split training instances into subsets, one for each branch extending from the root node. Then, this step is repeated for a chosen branch, using only those instances that actually reach it. In each step we choose the best subset among all subsets that are available for expansions. This constructing process continues until all nodes are pure or a specific number of expansions is reached. Figure i.1 shows the difference in split order between a hypothetical



**Fig i.1: Decision tree: a hypothetical best first decision tree.**

binary best-first tree. The problem in

growing best-first decision trees is now how to determine which attribute to split on and how to split the data. Because the most important objective of decision trees is to seek accurate and small models, we try to find pure nodes as soon as possible. To measure purity, we can use its opposite, impurity. There are many criteria to measure node impurity. The goal is to aim an attribute to split on that can maximally reduce impurity. When expanding nodes in the best-first tree the information and the Gini gain are also used to determine node order. The best-first method always chooses the node for expansion whose corresponding best split provides the best information gain or Gini gain among all unexpanded nodes in the tree.

### ii. J48

The J48 Decision tree classifier follows simple algorithm. In order to classifying for new item, it first needs to create a decision tree which is based on the attribute values of the available training dataset. So, whenever it encounters as a set of items (training set) it identifies the attributes that discriminates the different instances most clearly way. For the gaining of the highest information the feature that is able to tell us most about the data instances so that we can classify them the best is known as highest information gain. Now, in between the possible values of this feature, if there is any other value for which there is no uncertainty of meaning, that is, for which the data instances falling within this category have the same value for the object variable, then we terminate that branch and assign to it the object value that we have obtained.

For the different cases, we then look forward in another attribute that gives the highest information gain. Continuing in this manner until us either we run out of attributes, or get a clear decision of what combination of attributes gives us a particular object value. In the event that we run out of attributes, or if we cannot get an unambiguous result from the available information, we assign this branch a object value that the majority of the items under this branch possess. Now we have the decision tree, we follow the order of attribute selection as we have obtained for the tree. By checking all the relative attributes and their values with those seen in the decision tree model, we can assign or make prediction for the target value of this

iii. RepTree

REPTree builds a regression or decision tree using runs it using reduced-error pruning and information gain/variance reduction [19]. For optimizing the speed, it only identifies numeric attributes and sorts values for once. It deals with missing values by breaking instances into pieces. We can set the minimum number of instances per leaf, minimum proportion of training set variance for a split (numeric classes only), maximum tree depth (useful when boosting trees), and number of folds for pruning.

iv. Simple Cart

CART (Classification and Regression trees) incorporated a decision tree inducer for discrete classes much like that of C4.5, which was developed independently, and a scheme for inducing regression trees. Many of the techniques described, such as the method of the surrogate device for dealing with missing values and handling nominal attributes, were included in CART. However, first described by Quinlan (1992) this model trees did not appear until much more recently. The average value of the predicted attribute for the training tuples that reach the leaf every regression tree leaf stores a prediction of continuous-valued. Since the terms numeric prediction and regression are used synonymously in statistics, the resulting trees will be a regression trees, although it did not use any equations of regression model. By contrast, each leaf holds a regression model in model trees,—a multivariate linear equation for the predicted attribute. Regression as well as model trees tend to be more accurate than linear regression when the data are not

**MACHINE LEARNING TOOL**

Here we used WEKA 3.6.9 tool for analyzing our dataset. WEKA is a data mining system developed by the University of Waikato in New Zealand that implements data mining algorithms. In a real-world data mining problems WEKA is a state-of-the-art facility for developing machine learning (ML) techniques and their applications. It is a collection of different machine learning algorithms for data mining tasks and predictions. The algorithms are applied directly to a dataset. WEKA implements algorithms for data association, preprocessing, classification, clustering, regression rules; it has another advantage as visualization tools. This package of WEKA used for developing new machine learning schemes. Under the GNU General Public License, WEKA is a open source software [20].

**ALCOHOL CONSUMING STUDENT'S DATA**

We have taken the attributes from primary data for secondary school student course level data from the Portuguese student which created by Paulo Cortez and Alice Silva, University of Minho, Portugal [21]. The data used in this study is obtained from VBS Purvanchal University students of MCA department on 31 variables from the session 2007 to 2017 having 450 instances data-sets related to the MCA course. The

attributes closed questions related to several demographic e.g. school related, family income, social/emotional  
 e.g. alcohol consumption and mother's education e.g. number of past class failures variables that affect the performance of the student [22]. All the related attributes are showing in following Table 1.

**TABLE 1**  
**Alcohol Consumption Student DATA SET**

| Attribute | Domain                                                                                                                                           |
|-----------|--------------------------------------------------------------------------------------------------------------------------------------------------|
| sex       | student's sex (binary: "F" - female or "M" - male)                                                                                               |
| age       | student's age (numeric: from 15 to 22)                                                                                                           |
| address   | student's home address type (binary: "U" - urban or "R" - rural)                                                                                 |
| famsize   | family size (binary: "LE3" - less or equal to 3 or "GT3" - greater than 3)                                                                       |
| Pstatus   | parent's cohabitation status (binary: "T" - living together or "A" - apart)                                                                      |
| Medu      | mother's education (numeric: 0 - none, 1 - primary education (4th grade), 2 – 5th to 9th grade, 3 – secondary education or 4 – higher education) |
| Fedu      | father's education (numeric: 0 - none, 1 - primary education (4th grade), 2 – 5th to 9th grade, 3 – secondary education or 4 – higher education) |

|             |                                                                                                                                  |
|-------------|----------------------------------------------------------------------------------------------------------------------------------|
| Mjob        | mother's job (nominal: "teacher", "health" care related, civil "services" (e.g. administrative or police), "at_home" or "other") |
| Fjob        | father's job (nominal: "teacher", "health" care related, civil "services" (e.g. administrative or police), "at_home" or "other") |
| reason      | reason to choose this Institution (nominal: close to "home", Institution "reputation", "course" preference or "other")           |
| guardian    | student's guardian (nominal: "mother", "father" or "other")                                                                      |
| travelttime | home to school travel time (numeric: 1 - <15 min., 2 - 15 to 30 min., 3 - 30 min. to 1 hour, or 4 - >1 hour)                     |
| studytime   | weekly study time (numeric: 1 - <2 hours, 2 - 2 to 5 hours, 3 - 5 to 10 hours, or 4 - >10 hours)                                 |
| failures    | number of past class failures (numeric: n if $1 \leq n < 3$ , else 4)                                                            |
| schoolsup   | extra educational support (binary: yes or no)                                                                                    |
| famsup      | family educational support (binary: yes or no)                                                                                   |
| paid        | extra paid classes within the course (binary: yes or no)                                                                         |
| activities  | extra-curricular activities (binary: yes or no)                                                                                  |
| higher      | wants to take doctoral education (binary: yes or no)                                                                             |
| internet    | Internet access at home (binary: yes or no)                                                                                      |
| romantic    | with a romantic relationship (binary: yes or no)                                                                                 |
| famrel      | quality of family relationships (numeric: from 1 - very bad to 5 - excellent)                                                    |
| freetime    | free time after school (numeric: from 1 - very low to 5 - very high)                                                             |
| goout       | going out with friends (numeric: from 1 - very low to 5 - very high)                                                             |
| Dalc        | workday alcohol consumption (numeric: from 1 - very low to 5 - very high)                                                        |
| Walc        | weekend alcohol consumption (numeric: from 1 - very low to 5 - very high)                                                        |
| health      | current health status (numeric: from 1 - very bad to 5 - very good)                                                              |
| absences    | number of institution absences (numeric: from 0 to 93)                                                                           |
| G1          | first year grade (numeric: from 1 to 3) {1. First $\geq 60\%$ 2. Second $\geq 45$ & $<60\%$ 3. Fail $< 45\%$ }                   |
| G2          | second year grade (numeric: from 1 to 3)                                                                                         |
| G3          | final year grade (numeric: from 1 to 3)                                                                                          |

The target attribute G3 has a strong correlation with attributes G2 and G1. This occurs because G3 is the final year grade (issued at the 3rd year), while G1 and G2 correspond to the 1st and 2nd year grades. It is more difficult to predict G3 without G2 and G1, but such prediction is much more useful.

#### DATA MINING MODEL

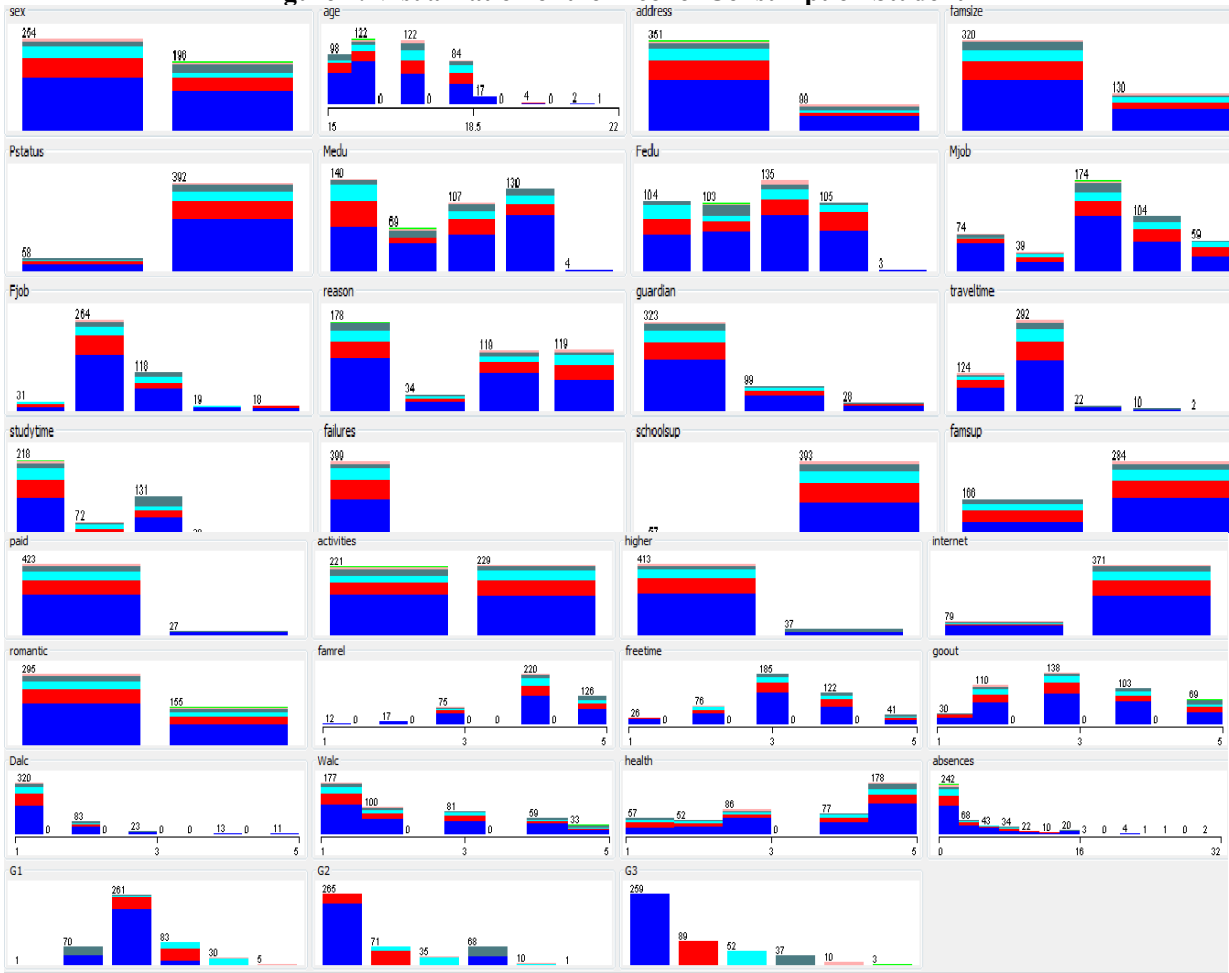
In the given survey BFTree, J48, RepTree and Simple Cart have been used to predict attributes such as school, sex, age, failures, study time and grades (G1, G2, G3) i.e. first year grade (numeric: from 1 to 3), second year grade (numeric: from 1 to 3), and third year grade (numeric: from 1 to 3), for chances of a student's getting perform

or not. We have used 10 folds cross validation to minimize any bias in the process and improve the efficiency of the process. The results show clearly that the proposed method performs well compared to other similar methods in the literature, taking into the fact that the attributes taken for analysis are not direct indicators of the performance of the students who is consuming alcohol during their study.

#### EXPERIMENTAL RESULT AND DISCUSSION

Here, we analyze Alcohol Consumption Student data set visually using different attributes and figure out the distribution of values. Figure 1 shows the distribution of values of Alcohol Consumption Student.

Figure 2: Visualization of the Alcohol Consumption Student



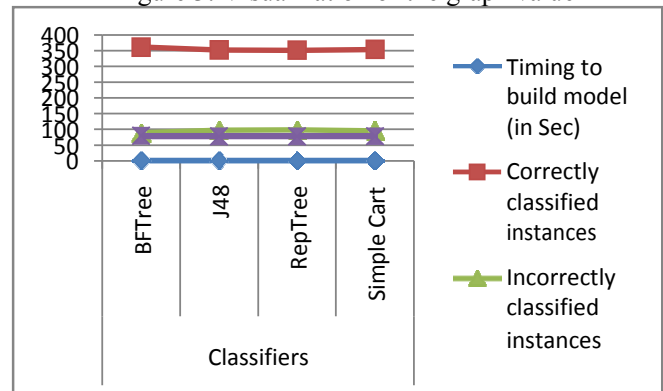
The above Figure 2 visualization of Alcohol Consumption Student data set shows that there are 450 students who studying MCA department and their final grade G3 shows 259 students are sufficient to their study, 89 students are good, 52 students are very good, 37 students are weak, 10 students are excellent and 03 students are poor.

Classifier model is a pruned decision tree in textual form that was produced on the full training data. Here we select Test mode:10-fold cross-validation. Table 2 shows the experimental result and figure 3 shows the corresponding graph values. We have carried out some weka experiments in order to evaluate the performance and usefulness of different classification algorithms for predicting student's performance.

Table 2: Evaluation of performance of students

| Evaluation Criteria              | Classifiers |         |         |             |
|----------------------------------|-------------|---------|---------|-------------|
|                                  | BFTree      | J48     | RepTree | Simple Cart |
| Timing to build model (in Sec)   | 0.83        | 0.09    | 0.05    | 0.66        |
| Correctly classified instances   | 361         | 352     | 351     | 354         |
| Incorrectly classified instances | 89          | 98      | 99      | 96          |
| Accuracy (%)                     | 80.2222     | 78.2222 | 78      | 78.6667     |

Figure 3: Visualization of the graph value



The above table 2 indicates that from the same source the results obtained from the training data are not optimistic compared with what might be obtained from the independent test set. In addition to the evaluation

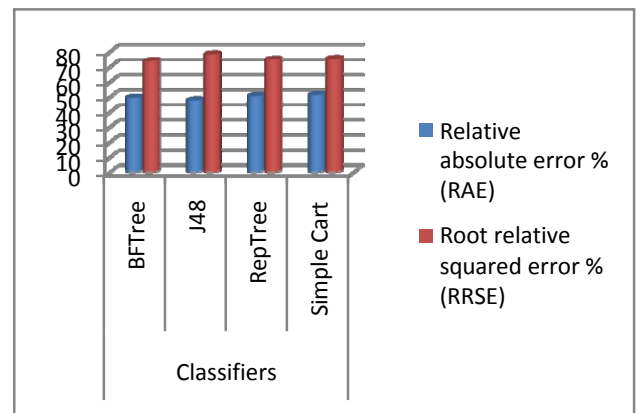
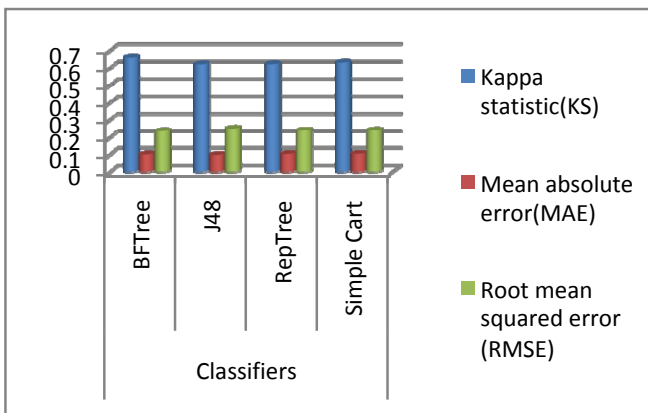
output measurements derived from the class probabilities assigned by the tree, classification error. In a similar way the mean absolute error calculated by using the absolute instead of squared difference. Because not all training instances are classified correctly that is the reason the errors are not 1 or 0. Here we can show that BFTree classifier has more accuracy than other classifiers that we used in our model. The percentage of correctly classified

instances is often called accuracy or sample accuracy of a model. Relative absolute error, kappa statistic, mean absolute error, root mean squared error and root relative squared error will be in numeric value only. From the following Table 3 shows the values in percentage value of relative absolute error and relative root squared error for evaluation and Figure 4 for the reference.

**Table 3: Training and Simulation Error**

| Evaluation Criteria                  | Classifiers |         |         |             |
|--------------------------------------|-------------|---------|---------|-------------|
|                                      | BFTree      | J48     | RepTree | Simple Cart |
| Kappa statistic(KS)                  | 0.6564      | 0.62    | 0.6205  | 0.63        |
| Mean absolute error(MAE)             | 0.1001      | 0.0967  | 0.1023  | 0.1038      |
| Root mean squared error (RMSE)       | 0.234       | 0.2482  | 0.237   | 0.2386      |
| Relative absolute error % (RAE)      | 49.0446     | 47.4012 | 50.1357 | 50.8614     |
| Root relative squared error % (RRSE) | 73.4482     | 77.8864 | 74.3877 | 74.8853     |

**Figure 4: Comparison between Parameters**





2. Amershi, S., Conati, C. (2006) Automatic Recognition of Learner Groups in Exploratory Learning Environments. Proceedings of ITS 2006, 8th International Conference on Intelligent Tutoring Systems.
3. U . K. Pandey, and S. Pal, .Data Mining: A prediction of performer or underperformer using classification., (IJCSIT) International Journal of Computer Science and Information Technology, Vol. 2(2), pp.686-690, ISSN:0975-9646, 2011.
4. Romero, Cristóbal, et al. "Data mining algorithms to classify students." Educational Data Mining 2008. 2008.
5. Galit.et.al, .Examining online learning processes based on log files analysis: a case study. Research, Reflection and Innovations in Integrating ICT in Education 2007.
6. Osmanbegović, Edin, and Mirza Suljić. "Data mining approach for predicting student performance." Economic Review 10.1 (2012).
7. U. K. Pandey, and S. Pal, .A Data mining view on class room teaching language., (IJCSI) International Journal of Computer Science Issue, Vol. 8, Issue 2, pp. 277-282, ISSN:1694-0814, 2011.
8. Cortez, Paulo, and Alice Maria Gonçalves Silva. "Using data mining to predict secondary school student performance." (2008).
9. J. Han and M. Kamber, .Data Mining: Concepts and Techniques., Morgan Kaufmann, 2000.
10. Kumar, S. Anupama, and M. N. Vijayalakshmi. "Efficiency of decision trees in predicting student's academic performance." First International Conference on Computer Science, Engineering and Applications, CS and IT. Vol. 2. 2011.
11. B.K. Bharadwaj and S. Pal. Data Mining: A prediction for performance improvement using classification., International Journal of Computer Science and Information Security (IJCSIS), Vol. 9, No. 4, pp. 136-140, 2011.
12. Z. N. Khan, .Scholastic achievement of higher secondary students in science stream., Journal of Social Sciences, Vol. 1, No. 2, pp. 84-87, 2005.
13. Baker, R.S.J.d., Corbett, A.T., Aleven, V. (2008) More Accurate Student Modeling Through Contextual Estimation of Slip and Guess Probabilities in Bayesian Knowledge Tracing. Proceedings of the 9th International Conference on Intelligent Tutoring Systems, 406-415.
14. Barnes, T., Bitzer, D., Vouk, M. (2005) Experimental Analysis of the Q-Matrix Method in Knowledge Discovery. Lecture Notes in Computer Science 3488: Foundations of Intelligent Systems, 603-611.
15. Beck, J.E., Mostow, J. (2008). How who should practice: Using learning decomposition to evaluate the efficacy of different types of practice for different types of students. Proceedings of the 9th International Conference on Intelligent Tutoring Systems, 353-362.
16. Cen, H., Koedinger, K., Junker, B. (2006) Learning Factors Analysis - A General Method for Cognitive Model Evaluation and Improvement. Proceedings of the 8th International Conference on Intelligent Tutoring Systems, 12-19.
17. HersHKovitz, A., Nachmias, R. (2008) Developing a Log-Based Motivation Measuring Tool. Proceedings of the First International Conference on Educational Data Mining, 226-233.
18. Frank, E. (2000). Pruning decision trees and lists. PhD thesis, The University of Waikato.
19. N. Landwehr, M. Hall, and E. Frank, —Logistic Model Trees, Machine Learning, pp. 161-205, 2005.
20. Weka Machine Learning Project, <http://www.cs.waikato.ac.nz/~ml/index.html>.
21. P. Cortez and A. Silva. Using data mining to predict secondary school student performance.[ in a. Brito and J. Teixeira eds. proceedings of 5th future business technology conference ( f ubutec2008) pp:5-12 porto portugal]. 2008.
22. U.K. Pandey and S. Pal, –Mining data to find adept teachers in dealing with students, International Journal of Intelligent Systems and Applications, vol. 4(3), 2012.
23. S. Pal and V. Chaurasia, –PERFORMANCE ANALYSIS OF STUDENTS CONSUMING ALCOHOL USING DATA MINING TECHNIQUES International Journal of Advance Research in Science & Engineering, vol.6(02), 2017, pp. 238-250.

# Design and Implementation Of Dynamic Trust Model for Individual Authorization

Thirumala Vasala, Danuka Nilima Priyadrshini

Department of Computer Science and Engineering  
Malla Reddy College of Engineering  
Kompally, Hyderabad, Telangana, India  
thirumala5b6@gmail.com, nilimadanuka31@gmail.com

**Abstract**—Development of authorization mechanisms for secure information access by a large community of users in an open environment is an important problem in the ever-growing Internet world. In this paper we propose a computational dynamic trust model for user authorization, rooted in findings from social science. Unlike most existing computational trust models, this model distinguishes trusting belief in integrity from that in competence in different contexts and accounts for subjectivity in the evaluation of a particular trustee by different trusters. Simulation studies were conducted to compare the performance of the proposed integrity belief model with other trust models from the literature for different user behavior patterns. Experiments show that the proposed model achieves higher performance than other models especially in predicting the behavior of unstable users.

**Index Terms**—Authorization, human factors, security, trust



## 1 INTRODUCTION

THE everyday increasing wealth of information available online has made secure information access mechanisms an indispensable part of information systems today. The mainstream research efforts for user authorization mechanisms in environments where a potential user's permission set is not predefined, mostly focus on role-based access control (RBAC), which divides the authorization process into the role-permission and user-role assignment. RBAC in modern systems uses digital identity as evidence about a user to grant access to resources the user is entitled to. However, holding evidence does not necessarily certify a user's good behavior. For example, when a credit card company is deciding whether to issue a credit card to an individual, it does not only require evidence such as social security number and home address, but also checks the credit score, representing the belief about the applicant, formed based on previous behavior. Such belief, which we call dynamic trusting belief, can be used to measure the possibility that

trust computation for the digital world closer to the evaluation of trust in the real world.

Unlike other trust models in the literature, the proposed model accounts for different types of trust. Specifically, it distinguishes trusting belief in integrity from that in competence.

The model takes into account the subjectivity of trust ratings by different entities, and introduces a mechanism to eliminate the impact of subjectivity in reputation aggregation.

Empirical evaluation supports that the distinction between competence and integrity trust is necessary in decision-making [15]. In many circumstances, these attributes are not equally important. Distinguishing between integrity and competence allows the model to make more informed and fine-grained authorization decisions in different contexts. Some real-world examples are as follows:

a user will not conduct harmful actions.

In this work, we propose a computational dynamic trust model for user authorization. Mechanisms for building trusting belief using the first-hand (direct experience) as well as second-hand information (recommendation and reputation) are integrated into the model. The contributions of the model to computational trust literature are:

The model is rooted in findings from social science, i.e., it provides automated trust management that mimics trusting behaviors in the society, bringing

- 1) On an online auction site, the competence trust of a seller can be determined by how quickly the seller ships an item, packaging/item quality etc., each being a different competence type. The integrity trust can be determined by whether he/she sells buyers' information to other parties without buyer consent. In the case of an urgent purchase, a seller with low integrity trust can be authorized if he/she has high competence trust.
- 2) For an online travel agency site, competence consists of elements such as finding the best car deals, the best hotel deals, the best flight deals etc., whereas integrity trust is based on factors like whether the site puts fraudulent charges on the customers' accounts. In a context where better deals are valued higher than the potential fraud risks, an agency with lower integrity trust could be preferred due to higher competence.
- 3) For a web service, the competence trust can include factors such as response time, quality of results etc., whereas integrity trust can depend on whether the service outsources requests to untrusted parties.

Y. Zhong and Y. Lu are with Microsoft Corporation, Seattle,

WA. E-mail: yuhuiz@hotmail.com, yilu.cn@gmail.com.

B. Bhargava and P. Angin are with the Department of Computer Science, Purdue University, West Lafayette, IN 47907.

E-mail: {bb, pangin}@cs.purdue.edu.

Manuscript received 6 June 2013; revised 8 Feb. 2014; accepted 18 Feb. 2014.

Date of publication 27 Feb. 2014; date of current version 16 Jan. 2015.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TDSC.2014.2309126

While government agencies would usually prefer high integrity in web services, high-competence services with low integrity could be authorized for real-time missions.

Experimental evaluation of the proposed integrity belief model in a simulated environment of entities with different behavior patterns suggests that the model is able to provide better estimations of integrity trust behavior than other major trust computation models, especially in the case of trustees with changing behavior.

## 2 RELATED WORK

### 2.1 McKnight's Trust Model

The social trust model, which guides the design of the computational model in this paper, was proposed by McKnight and Chervany [16] after surveying more than 60 papers across a wide range of disciplines. It has been validated via empirical study [15]. This model defines five conceptual trust types: trusting behavior, trusting intention, trusting belief, institution-based trust, and disposition to trust. Trusting behavior is an action that increases a truster's risk or makes the truster vulnerable to the trustee.

Trusting intention indicates that a truster is willing to engage in trusting behaviors with the trustee. A trusting intention implies a trust decision and leads to a trusting behavior. Two subtypes of trusting intention are:

- 1) Willingness to depend: the volitional preparedness to make oneself vulnerable to the trustee.
- 2) Subjective probability of depending: the likelihood that a truster will depend on a trustee.

Trusting belief is a truster's subjective belief in the fact that a trustee has attributes beneficial to the truster. The following are the four attributes used most often:

- 1) Competence: a trustee has the ability or expertise to perform certain tasks.
- 2) Benevolence: a trustee cares about a truster's interests.
- 3) Integrity: a trustee is honest and keeps commitments.
- 4) Predictability: a trustee's actions are sufficiently consistent.

Institution-based trust is the belief that proper structural conditions are in place to enhance the probability of achieving a successful outcome. Two subtypes of institution-based trust are:

- 1) Structural assurance: the belief that structures deployed promote positive outcomes. Structures include guarantees, regulations, promises etc.
- 2) Situational normality: the belief that the properly ordered environments facilitate success outcomes.

Disposition to trust characterizes a truster's general propensity to depend on others across a broad spectrum of situations. Two subtypes of disposition to trust are:

- 1) Faith in human: The assumptions about a general trustee's integrity, competence, and benevolence.
- 2) Trusting stance: A truster's strategy to depend on trustees despite his trusting belief about them.

Trust intention and trusting belief are situation and trustee specific. Institution-based trust is situation specific. Disposition to trust is independent of situation and trustee. Trusting belief positively relates to trusting intention, which in turn results in the trusting behavior. Institution-based trust positively affects trusting belief and trusting intention. Structural assurance is more related to trusting intention while situational normality affects both. Disposition to trust positively influences institution-based trust, trusting belief and trusting intention. Faith in humanity impacts trusting belief. Trusting stance influences trusting intention.

### 2.2 Computational Trust Models

The problem of establishing and maintaining dynamic trust has attracted many research efforts. One of the first attempts trying to formalize trust in computer science was made by Marsh [13]. The model introduced the concepts widely used by other researchers such as context and situational trust.

Many existing reputation models and security mechanisms rely on a social network structure [1]. Pujol et al. propose an approach to extract reputation from the social network topology that encodes reputation information [19]. Walter et al. [22] propose a dynamic trust model for social networks, based on the concept of feedback centrality. The model, which enables computing trust between two disconnected nodes in the network through their neighbor nodes, is suitable for application to recommender systems. Lang

[9] proposes a trust model for access control in P2P networks, based on the assumption of transitivity of trust in social networks, where a simple mathematical model based on fuzzy set membership is used to calculate the trustworthiness of each node in a trust graph symbolizing interactions between network nodes. Similarly, Long and Joshi [11] propose a Bayesian reputation calculation model for nodes in a P2P network, based on the history of interactions between nodes. Wang and Wang [23] propose a simple trust model for P2P networks, which combines the local trust from a node's experience with the recommendation of other nodes to calculate global trust. The model does not take the time of feedback into consideration, which causes the model to fail in the case of nodes with changing behavior. Reliance on a social network structure limits wide applicability of the mentioned approaches, especially for user authorization.

FCTrust [8] uses transaction density and similarity to calculate a measure of credibility of each recommender in a P2P network. Its main disadvantages are that it has to retrieve all transactions within a certain time period to calculate trust, which imposes a big performance penalty, and that it does not distinguish between recent and old transactions. SFTrust [25] is a double trust metric model for unstructured P2P networks, separating service trust from feedback trust. Its use of a static weight for combining local and recommendation trust fails to capture node specific behavior.

Das and Islam [3] propose a dynamic trust computation model for secure communication in multi-agent systems, integrating parameters like feedback credibility, agent similarity, and direct/indirect trust/recent/historical trust into trust computation. Matt et al. [14] introduce a method for modeling the trust of a given agent in a multi-agent

system by combining statistical information regarding the past behavior of the agent with the agent's expected future behavior.

A distributed personalized reputation management approach for e-commerce is proposed by Yu and Singh [24]. The authors adopt ideas from Dempster-Shafer theory of evidence to represent and evaluate reputation. If two principals "a" and "b" have direct interactions, b evaluates a's reputation based on the ratings of these interactions. Other-wise, b queries a TrustNet for other principals' local beliefs about a. The reputation of "a" is computed based on the gathered local beliefs using Dempster-Shafer theory.

Sabater and Sierra propose a reputation model called the Regret system [20] for gregarious societies. The authors assume that a principal owns a set of sociograms describing the social relations in the environment along individual, social and ontological dimensions. The performance highly depends on the underlying sociograms, although how to build sociograms is not discussed.

The above mentioned trust computation approaches do not consider "context" as a factor affecting the value of trust, which prevents an accurate representation for real life situations. Skopik et al. [21] propose a dynamic trust model for complex service-oriented architectures based on fuzzy logic. Zhu et al. [26] introduce a dynamic role based access control model for grid computing. The model determines authorization for a specific user based on its role, task and the context, where the authorization decision is updated dynamically by a monitoring module keeping track of user attributes, service attributes and the environment. Fan et al.

[5] propose a similar trust model for grid computing, which focuses on the dynamic change of roles of services. Liu and Liu [10] propose a Bayesian trust evaluation model for dynamic authorization in a federation environment, where the only context information is the domain from which authorization is requested. Ma and He [12] propose a genetic algorithm for evaluating trust in distributed applications. Nagarajan and Varadharajan [18] propose a security model for trusted platform based services based on evaluation of past evidence with an exponential time decay function. The model evaluates trust separately for each property of each component of a platform, similar to the consideration of competence trust in our proposed model. Although these approaches integrate context into trust computation, their application is limited to specific domains different from the one considered in our work.

### 3 OVERVIEW OF THE TRUST MODEL

The trust model we propose in this paper distinguishes integrity trust from competence trust. Competence trust is the trusting belief in a trustee's ability or expertise to perform certain tasks in a specific situation. Integrity trust is the belief that a trustee is honest and acts in favor of the truster. Integrity and benevolence in social trust models are combined together. Predictability is attached to a competence or integrity belief as a secondary measure.

The elements of the model environment, as seen in Fig. 1, include two main types of actors, namely trusters and trustees, a database of trust information, and different contexts, which depend on the concerns of a truster and the

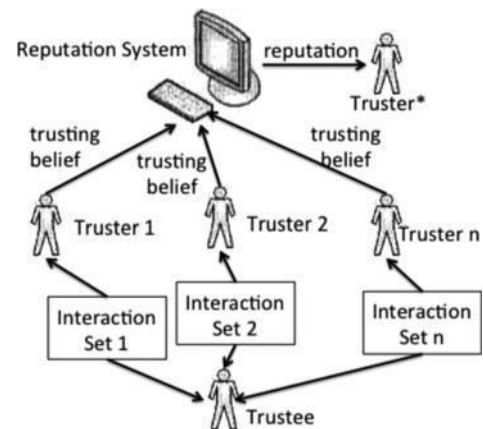


Fig. 1. Model elements.

competence of a trustee. For the online auction site example in Section 1, let us assume that buyer B needs to decide whether to authorize seller S to charge his credit card for an item I (authorize access to his credit card/contact information). The elements of the model in this case are:

Trusters are the buyers registered to the auction site.

Trustees are the sellers registered to the auction site. The context states how important for B the shipping, packaging and item quality competences of S for item I are. It also states how important for B the integrity of S is for this transaction.

B can gather trust information about S from a data-base maintained by the site or a trusted third party. This information includes the ratings that S received from buyers (including B's previous ratings, if any) for competence in shipping, packaging and quality of I as well as S's integrity. It also includes the ratings of buyers (including B) for sellers other than S in different contexts and ratings of S for different items. Trust evaluation is recorded in the database when a buyer rates a transaction with a seller on the site.

#### 3.1 Context and Trusting Belief

Context. Trust is environment-specific [13]. Both trusters' concern and trustees' behavior vary from one situation to another. These situations are called contexts. A truster can specify the minimum trusting belief needed for a specific context. Direct experience information is maintained for each individual context to hasten belief updating.

In this model, a truster has one integrity trust per trustee in all contexts. If a trustee disappoints a truster, the misbehavior lowers the truster's integrity belief in him. For integrity trust, contexts do not need to be distinguished. Competence trust is context-dependent. The fact that Bob is an excellent professor does not support to trust him as a chief. A representation is devised to identify the competence type and level needed in a context. Two functions that relate contexts are defined.

Let  $S_C$  denote the universe consisting of all types of competences of interest  $\{C_1; C_2; \dots; C_n\}$ , where each  $C_i$  is a different competence type. For example,  $S_C \not\sim \{\text{cooking, teaching, writing, } \dots\}$ . Let  $M_C$  denote the measurement of

a competence type  $c$ . For example,  $M_{\text{cooking}} \frac{1}{4} \{\text{very bad, bad, ok, good, excellent}\}$ . A partial order  $< : M \rightarrow M$  and a function  $\text{dis} : M \times M \rightarrow [0; 1]$  are defined on  $M_C$ . For two elements  $m_i$  and  $m_j$  in  $M_C$ ;  $m_j$  is the higher competence level if  $m_i < m_j$ . The  $\text{dis}$  function measures the numerical distance between two elements (its two arguments), outputting a value between 0 and 1.  $m_k$  is closer to  $m_i$  than  $m_j$  is, if  $\text{dis}(m_i, m_k) < \text{dis}(m_i, m_j)$ .

Two functions, `hCtxxt: contextId contextId ! {true, false}` and `simLCTX: contextId contextId  $R_{\frac{1}{20}; 1\&}$  ! {true, false}`, are defined as follows. Here,  $R_{\frac{1}{20}; 1\&}$  denotes a real number ranging from 0 and 1. `ct` and `cl` are abbreviations of `cType` and `cLevel`.

If  $\text{hContext}(\mathbf{c}_1; \mathbf{c}_2)$  is true,  $\mathbf{c}_2$  requires the same type of competence with higher level as  $\mathbf{c}_1$  does. SimLCTX specifies whether the levels required for two contexts with the same type are sufficiently close to each other.

TABLE 1  
Trust Model Notation

|                                                                        |                                                                                                            |
|------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------|
| $TC_{t_1 \rightarrow u_1}^v(c),$<br>$TC_{t_1 \rightarrow u_1}^p(c):$   | $t_1$ 's initial or continuous trusting belief in $u_1$ 's competence in context $c$ .                     |
| $DTC_{t_1 \rightarrow u_1}^v(c),$<br>$DTC_{t_1 \rightarrow u_1}^p(c):$ | $t_1$ 's competence belief about $u_1$ in $c$ based on direct experience (called direct competence trust). |
| $RC_{u_1}^v(c),$<br>$RC_{u_1}^p(c):$                                   | $u_1$ 's competence reputation in context $c$ .                                                            |
| $TI_{t_1 \rightarrow u_1}^v,$<br>$TI_{t_1 \rightarrow u_1}^p:$         | $t_1$ 's initial or continuous trusting belief in $u_1$ 's integrity.                                      |
| $DTI_{t_1 \rightarrow u_1}^v,$<br>$DTI_{t_1 \rightarrow u_1}^p:$       | $t_1$ 's integrity belief about $u_1$ based on direct experience (direct integrity trust).                 |
| $RI_{u_1}^v, RI_{u_1}^p:$                                              | $u_1$ 's integrity reputation.                                                                             |

This section presents the operations defined on the trust model. The notations in Table 1 are used for presentation. The notation with superscript  $v$  is the value of a belief. The one with superscript  $p$  is the associated predictability.

The trust model defines four types of operations:

[Operation 2.] Aggregate trusting beliefs about  $u_1$  from multiple trusters to his reputation  $\delta RC_{u_1}^V \delta \phi B$ ;  $RC_{u_1}^P \delta \phi B$  or  $\delta RI_{u_1}^V$ ;  $RI_{u_1}^P \triangleright$ . A method that evaluates reputation is called a reputation aggregation method. A reputation aggregation method takes trusting beliefs from different trusters as input.

[Operation 4.] Update and test continuous trust. This operation computes  $\delta TC^v$ !

$$t_1 \quad u_1 \quad \delta c p; TC_t^p \quad 1!u_1 \quad \delta c p \text{ or } \delta Tl_t^v \quad 1!u_1; Tl_t^p \quad 1!u_1 \text{ p.}$$

TABLE 2  
Test a Competence Trusting Belief

|                                               | $TC_{t_1 \rightarrow u_1}^p(c) \geq \delta_c$ | $TC_{t_1 \rightarrow u_1}^p(c) < \delta_c$ |
|-----------------------------------------------|-----------------------------------------------|--------------------------------------------|
| $TC_{t_1 \rightarrow u_1}^p(c) \leq \delta_p$ | true                                          | false                                      |
| $TC_{t_1 \rightarrow u_1}^p(c) > \delta_p$    | uncertain                                     | false                                      |

The results are compared with two constants  $\delta_c$  and  $\delta_p$  respectively. This is the same as Operation 3 except that it is used for updating beliefs, not initializing.

Belief formation and reputation aggregation are atomic operations. They are presented in the next two sections. The last two operations are needed for user authorization. They output a test result according to Table 2 or Table 3.

The methods that can be used to build a trusting belief are summarized below, followed by a discussion of using these methods to establish initial or continuous trust.

### 3.2.1 Methods to Build a Trusting Belief

Seven methods that can be used to build competence trust:

[M1.] Form trusting belief based on direct experience in a specific context.

Precondition:  $DTC_{t_1 \rightarrow u_1}^v \delta c p \neq \text{null}$

$$TC_{t_1 \rightarrow u_1}^v \delta c p : \frac{1}{4} DTC_{t_1 \rightarrow u_1}^v \delta c p; \quad (3a)$$

$$TC_{t_1 \rightarrow u_1}^p \delta c p : \frac{1}{4} DTC_{t_1 \rightarrow u_1}^p \delta c p; \quad (3b)$$

[M2.] Consider direct trust about  $u_1$  in contexts that require a higher competence level than  $c$ . Use the maximum value and minimum predictability.

Precondition:  $9c^0 \text{ hContxt } \delta c; c^0 \neq \text{null} \& DTC_{t_1 \rightarrow u_1}^v \delta c^0 \neq \text{null}$

$$TC_{t_1 \rightarrow u_1}^v c : \max_{\delta} DTC_{t_1 \rightarrow u_1}^v c^0 \text{ t1!u1 } \delta c^0 \neq \text{null}; \quad (4a)$$

$$TC_{t_1 \rightarrow u_1}^p c : \min_{\delta} DTC_{t_1 \rightarrow u_1}^p c^0 \text{ t1!u1 } \delta c^0 \neq \text{null}; \quad (4b)$$

[M3.] Consider direct trust about  $u_1$  in contexts that require a lower competence level than  $c$ . Use the minimum value and maximum predictability.

Precondition:  $9c^0 \text{ simLCTX } \delta c; c^0 \neq \text{null} \& DTC_{t_1 \rightarrow u_1}^v \delta c^0 \neq \text{null}$

$$TC_{t_1 \rightarrow u_1}^v c : \min_{\delta} DTC_{t_1 \rightarrow u_1}^v c^0 \text{ t1!u1 } \delta c^0 \neq \text{null}; \quad (5a)$$

$$TC_{t_1 \rightarrow u_1}^p c : \max_{\delta} DTC_{t_1 \rightarrow u_1}^p c^0 \text{ t1!u1 } \delta c^0 \neq \text{null}; \quad (5b)$$

TABLE 3  
Test Integrity Trusting Belief

|                                            | $TI_{t_1 \rightarrow u_1}^v \geq \delta_c$ | $TI_{t_1 \rightarrow u_1}^v < \delta_c$ |
|--------------------------------------------|--------------------------------------------|-----------------------------------------|
| $TI_{t_1 \rightarrow u_1}^p \leq \delta_p$ | true                                       | false                                   |
| $TI_{t_1 \rightarrow u_1}^p > \delta_p$    | uncertain                                  | false                                   |

[M4.] Request  $u_1$ 's competence reputation in context  $c$ .

Precondition:  $9t^0 DTC_{t_1 \rightarrow u_1}^v \delta c p \neq \text{null}$

$$TC_{t_1 \rightarrow u_1}^v \delta c p : \frac{1}{4} RC_{u_1}^v \delta c p; \quad (6a)$$

$$TC_{t_1 \rightarrow u_1}^p \delta c p : \frac{1}{4} RC_{u_1}^p \delta c p; \quad (6b)$$

[M5.] Use the most common belief value about trustees that  $t_1$  encountered in  $c$ . Suppose the belief values are in the range of  $(a, b)$ . Partition  $(a, b)$  into  $k$  (e.g., 10) intervals. Let  $(a', b')$  be the interval containing most values. If there are multiple such intervals (called multi-modal situation),

use the uncertainty handling policies in  $t_1$ 's profiles to choose one. mode  $DTC_{t_1 \rightarrow u_1}^v c$  is defined as  $\frac{a' + b'}{2}$ . mode

$\delta DTC_{t_1 \rightarrow u_1} \delta c p$  is computed in the same way.

Precondition:  $9u^0 DTC_{t_1 \rightarrow u_1}^v \delta c p \neq \text{null}$

$$TC_{t_1 \rightarrow u_1}^v \delta c p : \frac{1}{4} \text{ mode } DTC_{t_1 \rightarrow u_1}^v \delta c p \text{ jDTC}_{t_1 \rightarrow u_1}^v \delta c p \neq \text{null}; \quad (7a)$$

$$TC_{t_1 \rightarrow u_1}^p \delta c p : \frac{1}{4} \text{ mode } DTC_{t_1 \rightarrow u_1}^p \delta c p \text{ jDTC}_{t_1 \rightarrow u_1}^p \delta c p \neq \text{null}; \quad (7b)$$

[M6.] Use the most common belief about all trustees encountered by all trusters in  $c$ . This method is considered only if (1) both  $c$  and  $u_1$  are new to  $t_1$ ; and (2) no truster knows  $u_1$ . mode  $\delta DTC_{t_1 \rightarrow u_1} \delta c p$  and mode  $\delta DTC_{t_1 \rightarrow u_1}^p \delta c p$  are computed in the same way as mode  $\delta DTC_{t_1 \rightarrow u_1}^v \delta c p$ .

Precondition:  $9u^0; t^0 DTC_{t_1 \rightarrow u_1}^v \delta c p \neq \text{null}$

$$TC_{t_1 \rightarrow u_1}^v \delta c p : \frac{1}{4} \text{ mode } DTC_{t_1 \rightarrow u_1}^v \delta c p \text{ jDTC}_{t_1 \rightarrow u_1}^v \delta c p \neq \text{null}; \quad (8a)$$

$$TC_{t_1 \rightarrow u_1}^p \delta c p : \frac{1}{4} \text{ mode } DTC_{t_1 \rightarrow u_1}^p \delta c p \text{ jDTC}_{t_1 \rightarrow u_1}^p \delta c p \neq \text{null}; \quad (8b)$$

[M7.] Use priori competence trusting belief specified in  $t_1$ 's local or global profile (defined in Section 3.3). The priori belief in the local profile overrides the global one.

Four methods that can be used to build integrity belief:

[M8.] Form trusting belief based on direct experience if there are previous interactions.

Precondition:  $DTI_{t_1 \rightarrow u_1}^v \neq \text{null}$

$$TI_{t_1 \rightarrow u_1}^v : \frac{1}{4} DTI_{t_1 \rightarrow u_1}^v; \quad (9a)$$

$$TI_{t_1 \rightarrow u_1}^p : \frac{1}{4} DTI_{t_1 \rightarrow u_1}^p; \quad (9b)$$

TABLE 4  
Candidate Method Set to Build Initial Competence Trust

|                | c is new                  | c is known                    |
|----------------|---------------------------|-------------------------------|
| $u_1$ is new   | $\{M4\} > \{M6, M7\}$     | $\{M4\} > \{M5, M7\}$         |
| $u_1$ is known | $\{M2, M3, M4\} > \{M7\}$ | $\{M2, M3, M4\} > \{M5, M7\}$ |

[M9.] Request  $u_1$ 's integrity reputation.

Precondition:  $9t^0 DT_{t_1 \rightarrow u_1}^v \neq \text{null}$

$$T_{t_1 \rightarrow u_1}^v : \frac{1}{4} R_{u_1}^v ; \quad (10a)$$

$$T_{t_1 \rightarrow u_1}^p : \frac{1}{4} R_{u_1}^p ; \quad (10b)$$

[M10.] Use the most common beliefs about trustees that  $t_1$  encountered. mode  $\delta DT_{t_1 \rightarrow u_1}^v$  and mode  $\delta DT_{t_1 \rightarrow u_1}^p$  are computed in the same way as mode  $\delta DT_{t_1 \rightarrow u_1}^v$ . This method is always applicable except for the first trustee encountered by  $t_1$ .

Precondition:  $9u^0 DT_{t_1 \rightarrow u_1}^v \neq \text{null}$

$$T_{t_1 \rightarrow u_1}^v : \frac{1}{4} \text{ mode } DT_{t_1 \rightarrow u_1}^v ; DT_{t_1 \rightarrow u_1}^v \neq \text{null} ; \quad (11a)$$

$$T_{t_1 \rightarrow u_1}^p : \frac{1}{4} \text{ mode } DT_{t_1 \rightarrow u_1}^p ; DT_{t_1 \rightarrow u_1}^p \neq \text{null} ; \quad (11b)$$

[M11.] Use priori integrity trusting belief specified in  $t_1$ 's global profile.

### 3.2.2 Building and Testing Trusting Beliefs

Different methods are used under various situations for building and testing trusting beliefs. A candidate method set includes the methods considered in a specific situation. A method is applicable only if: (1) it is in the current candidate method set, and (2) its precondition holds.

Building and testing initial competence trust. There are four scenarios when  $t_1$  is about to establish initial trust about  $u_1$  in  $c$ : (1) both  $c$  and  $u_1$  are new; (2)  $c$  is known but  $u_1$  is new; (3)  $c$  is new but  $u_1$  is known; (4) both  $c$  and  $u_1$  are known. A context  $c$  is known if the truster has experience with some trustee in  $c$ . A trustee  $u_1$  is known if she interacted with  $t_1$  before. The candidate method set for all scenarios and the order of their priorities are summarized in Table 4. is a partial order defined on the method priority set. The relationship between two methods enclosed in one “{ }” is undefined by the model itself. This is an ambiguous priority set. is extended to a total order according to  $t_1$ 's method preference policies.

The algorithm to build and test an initial competence trusting belief is shown in Fig. 3. The algorithm initializes unusedMS using the appropriate candidate method set. It chooses the applicable method  $M$  with highest priority in unusedMS. The input threshold parameters  $d_c$  and  $d_p$  are compared with the trusting belief generated by  $M$ . If “true” or “false” is obtained, this result is output. Otherwise  $M$  is removed, trusting belief is saved and the process is repeated with the next  $M$ . In the case that the algorithm outputs no

```

Input:  $t_1, u_1, c, \delta_c, \delta_p$ 
Output : true/false
unusedMS := candidate method set defined in Table 4
i := 1
while unusedMS  $\neq \emptyset$ 
    M := the applicable method with highest priority
    result[i] := compute( $TC_{t_1 \rightarrow u_1}^v(c), TC_{t_1 \rightarrow u_1}^p(c)$ ) using M
    testResult := compare result[i] with  $\delta_c, \delta_p$  based on Table 2
    if (testResult = uncertain){
        i := i + 1;
        delete M from unusedMS
    }
    else{
        return testResult
    }
}
Choose r from {results[i] | 0} based on imprecision handling policy
return (r.value >  $\delta_c$ )

```

Fig. 3. Algorithm to build/test initial competence trusting belief. result after all methods are considered, one trusting belief is chosen (i.e.,  $r$  is chosen among all results) based on imprecision handling policies. The value of the belief is compared with  $d_c$ .

Building initial integrity trust. Truster  $t_1$  uses her priori integrity trusting belief for the first trustee she encountered. If  $u_1$  is not the first trustee, the candidate method set and the order of their priorities are: fM9g fM10; M11g. The algorithm to build and test initial integrity trusting belief is similar to that in Fig. 3.

Building continuous competence trust. The candidate method set and the order of their priorities are: f M1g fM2; M4g. The algorithm to build and test continuous competence trust is similar to that in Fig. 3.

Building continuous integrity trust. The candidate method set and the order of their priorities are: fM8g fM9; M10g. The algorithm to build and test continuous integrity trust is similar to that in Fig. 3.

### 3.3 Global and Local Profiles

Each truster  $t_1$  has one global profile. The profile contains:

- (1)  $t_1$ 's priori integrity and competence trusting belief;
  - (2) method preference policies;
  - (3) imprecision handling policies;
  - (4) uncertainty handling policies;
  - (5) parameters needed by trust-building methods.
- $t_1$  can have one local profile for each context. Local profiles have a similar structure as global profiles. The content in a local profile over-rides that in the global one. Fig. 4 shows the definition of global and local profiles.

As aforementioned, method preference policies, defined as PreferencePolicy, are to extend the partial order to a total order. Therefore, no two methods have the same priority. iCompetence and cCompetence are used when building initial and continuous competence trust respectively. iIntegrity and cIntegrity are for establishing integrity trusting belief. Relationships are separately defined on each ambiguous priority set. For

```

GlobalProfile ::= <truster, Priori, PolicySet, MethodPar>
Local Profile ::= <truster, contextId, Priori, PolicySet,
                                     MethodPar>

Priori ::= <IntegrityPriori, CompetencePriori>
IntegrityPriori ::= <value, predictability>
CompetencePriori ::= <value, predictability>
value ::=  $R^{[0,1]}$ 
predictability ::=  $R^{[0,1]}$ 
PolicySet ::= <PreferencePolicy, ImprecisionPolicy,
                                     UncertainPolicy>
PreferencePolicy ::= <iCompetence, iIntegrity, cCompetence,
                                     cIntegrity>
iCompetence ::= <<mld2>, <mld2>, <<mld3>, <mld2>
>>, <<mld3>, <mld2>>>
iIntegrity ::= <mld2>
cCompetence ::= <mld2>
cIntegrity ::= <mld2>
mld ::= string
ImprecisionPolicy ::= "false" | "priority" |
                                     MinPredictability | tValue
MinPredictability ::= "priority" | tValue
tValue ::= "min" | "max" | "median"
UncertainPolicy ::= tValue
MethodPar ::= <mld, parList>
    
```

Fig. 4. Global and local profile definitions.

example, the fourth scenario in building initial competence trust has two ambiguous priority sets {M2, M3, M4} and {M5, M7}. Hence, the third part of iCompetence is <<mld<sup>3</sup>>, <mld<sup>2</sup>>>. Here, mld is the identifier of a method.

<mld<sup>n</sup>> is the abbreviation of a string <mld; . . . ; mld> with n mlds. A method whose mld is in the *i*th place has the *i*th highest priority in that set.

Imprecision handling policies are used to choose a belief value when the tests on trusting beliefs generated by all applicable methods return "uncertain". There are three types of policies: If the "false" policy is specified, use "0" as the belief value. This implies that a test failed. If the "priority" policy is adopted, the belief generated by the method with the highest priority is chosen. If "MinPredictability" policies are used, the belief with the lowest predictability is selected. If multiple beliefs have the lowest predictability, they are distinguished by tValue, which is a constant specifying whether to choose the minimum, maximum or median belief value. The value of this constant is also set by the imprecision handling policy.

Uncertainty handling policies are used by three belief building methods, M5, M6, and M10, in the multimodal situation. Minimum, maximum or median values can be used based on the policy. They correspond to pessimistic, optimistic, and realistic attitudes as argued by Marsh [13].

MethodPar provides the values to the parameters a method needs. Currently, only the third method needs d that specifies the proximity threshold between two contexts.

## 4 BELIEF INFORMATION AND REPUTATION AGGREGATION METHODS

### 4.1 Competence Belief

Belief about a trustee's competence is context specific. A trustee's competence changes relatively slowly with time. Therefore, competence ratings assigned to her are viewed

as samples drawn from a distribution with a steady mean and variance. Competence belief formation is formulated as a parameter estimation problem. Statistic methods are applied on the rating sequence to estimate the steady mean and variance, which are used as the belief value about the trustee's competence and the associated predictability.

Let *R* denote the competence rating set  $R = \{r_1; \dots; r_n\}$  where  $r_i \in [0, 1]$ ;  $r_1, r_2, \dots, r_n$  are independently drawn from an underlying distribution. The mean and variance of the distribution need to be inferred based on the ratings. Like any statistical inference problem, the inference contains two parts: (1) estimated value (2) a measure of its goodness. Usually, a distribution from a restricted family is chosen to approximate the underlying distribution. We use the Normal distribution  $N(m, s^2)$ , where *m* corresponds to the mean and *s*<sup>2</sup> corresponds to the variance. *m* estimates the trusting belief about the trustee's competence denoted as *b<sub>C</sub>*; *s*<sup>2</sup> characterizes the variability of the underlying distribution and is positively correlated with predictability. The goodness of estimation is measured via a confidence interval. 90 percent confidence intervals for *m* and *s*<sup>2</sup> are constructed. The length of the confidence interval of *m* is used as the associated predictability *p<sub>C</sub>*. This method assumes that the observations are drawn from a normal distribution. This assumption may not hold and the result may be misleading. Goodness-of-fit tests can check whether the assumption is valid or not. In summary, the approach is: (1) estimate *m* and *s*<sup>2</sup>; (2) measure the goodness of the inferences; (3) test the normality of the distribution (4) Let *b<sub>C</sub>* = *m* and *p<sub>C</sub>* = length of the confidence interval if the inference is good enough and the data set approximately follows a normal distribution.

k-Statistic defined in (12) is used in computation. Let *m*<sup>^</sup> and *s*<sup>^</sup> denote the estimation values of *m* and *s*. The unbiased estimators of *m* and *s*<sup>2</sup> are *k*<sub>1</sub> and *k*<sub>2</sub>.

$$\begin{aligned}
 S_m &= \frac{1}{n} \sum_{i=1}^n r_i; \\
 k_1 &= S_1 = m; \\
 k_2 &= \frac{1}{n} S_2 = \frac{1}{n} \sum_{i=1}^n r_i^2;
 \end{aligned} \tag{12}$$

For *m*, if the number of ratings *n* is greater than 45, the length of the confidence interval can be approximated by (13a). Here, 1.645 is the z value such that 5 percent of the whole area lies to its right in a standard normal distribution. In this equation, *s* is replaced by *s*<sup>^</sup>. This substitution is applicable only when the size of the rating set is large. In the case there are few ratings, i.e., *n* < 45, (13a) is not suitable. Fortunately, the underlying population distribution is normal based on the assumption. t-distribution (i.e., student distribution) [17] is used to construct the confidence interval. Equation (13b) computes the length of the confidence interval in this case. In this equation, *t*<sub>0.05(*n*-1)</sub> denotes the critical value of t distribution with (*n*-1) degrees of freedom such that 5 percent of the area lies to its right:

$$\begin{aligned}
 \text{interval length for } m: \\
 & \frac{1}{n} \left( 1.645 \sqrt{\frac{s^2}{n}} \right) \quad n \geq 45 \quad \delta_{13a}; \\
 & \frac{1}{n} \left( t_{0.05(n-1)} \sqrt{\frac{s^2}{n}} \right) \quad n < 45 \quad \delta_{13b};
 \end{aligned}$$

$n \frac{1}{4} 45$  is chosen as the dividing line between large and small rating sets due to two reasons: (1) The critical value of  $t$  is always larger than the corresponding  $z$  value.  $t$  distribution approaches the normal distribution as  $n$  increases.

(2) The critical value of  $t_{0.05 \frac{1}{2} 1:6794}$  is quite close to  $z_{0.05}$ , (i.e., 1.645) when the degree of freedom is 44.

The length of the 90 percent confidence interval for  $s^2$  is shown in (14a). Here,  $x_{0.05 \frac{1}{2} 1:6794}$  is the critical value of  $\chi^2$  distribution with  $(n-1)$  degrees of freedom such that 5 percent of the area lies to the right.  $x_{0.95 \frac{1}{2} 1:6794}$  is defined similarly. Unlike  $z$  values,  $\chi^2$  is asymmetric and  $x_{0.05 \frac{1}{2} 1:6794} \neq x_{0.95 \frac{1}{2} 1:6794}$ . Equation (14a) is a straightforward application of Fisher's result (i.e.,  $x_{\alpha}^2$  is the value of  $\chi^2$  when  $n$  is large) for (14b).

$$\text{interval length for } s^2 = \frac{1}{4} 8 \frac{2\delta n \frac{1}{2} 1:6794}{1:645 \frac{1}{2} 1:645} \frac{2\delta n \frac{1}{2} 1:6794}{2n \frac{1}{2} 1:645} \frac{2\delta n \frac{1}{2} 1:6794}{2n \frac{1}{2} 1:645} \quad n \geq 45 \quad (14a)$$

In this model, we assume the existence of a reputation server that acts properly on behalf of trusters. It is assumed that trusters are honest in providing information. The attacks discussed in [4] do not exist. Trusters are subjective and utilize different evaluation criteria. Reputation aggregation methods shall eliminate the effect of subjectivity and output a result close to the trusting belief the reputation requester would have obtained if she had directly interacted with the trustee.

Let  $t$  denote the truster who requests reputation information about trustee  $u$ . Let  $t_1; t_2; \dots; t_k$  denote the trusters who submit a direct competence trusting belief about  $u$  to the reputation server. Suppose  $u$  follows a distribution with mean  $m$  and variance  $s^2$  from the perspective of  $t$ 's. Please note  $m$  and  $s^2$  are true values, not estimated values from existing ratings. Let  $m_i$  and  $s_i^2$  denote the mean and variance of  $u$  from  $t_i$ 's point of view. Because of the subjectivity,  $m$  and  $s^2$  are different from  $m_i$  and  $s_i^2$  even when  $u$  behaves consistently. Subjectivity between trusters is formulated in (15). Here,  $Dm_i$  and  $c_i$  are constants. Equation (15a) is interpreted as "An excellent behavior for Alice is just good for Bob". Equation (15b) can be explained by "the rating range of Bob is greater than that of Alice".

$$m_i \frac{1}{4} m \leq Dm_i; \quad (15a)$$

$$s_i^2 \frac{1}{4} c_i s^2; \quad (15b)$$

To eliminate the subjectivity of a truster from the perspective of a requestor is to calibrate such deviations.  $m$  and  $s^2$  are estimated based on estimated  $m_i$  and  $s_i^2$ , and interaction numbers submitted by  $k$  trusters. They are denoted as  $hm^{\wedge}_1; s^{\wedge}_1; n_1; hm^{\wedge}_2; s^{\wedge}_2; n_2; \dots; hm^{\wedge}_k; s^{\wedge}_k; n_k$ . From Section 4.1, we know  $m^{\wedge}_i$  can be viewed as the value of a random variable with mean  $m^{\wedge}_i$  and variance  $c_i s^{\wedge}_i$ .  $s^{\wedge}_i$  is the value of a random variable whose mean is  $c_i s$ .

Equation (16) defines an estimator for  $m$ . Two estimators for  $s^2$  are given in (17a) and (17b).

$$\text{estimator for } m = \frac{1}{k} \sum_{i=1}^k \frac{m^{\wedge}_i}{Dm_i}; \quad (16)$$

$$\text{estimator for } s^2 = \frac{1}{k} \sum_{i=1}^k \frac{\delta n_i \frac{1}{2} 1:6794}{c_i}; \quad (17a)$$

$$\text{estimator for } s^2 = \frac{1}{k} \sum_{i=1}^k \frac{\delta n_i \frac{1}{2} 1:6794}{c_i s^{\wedge}_i}; \quad (17b)$$

Estimation bound. Let  $X_i$  denote the random variable for  $m^{\wedge}_i$ . Let  $Y_1 = \frac{1}{k} \sum_{i=1}^k X_i$ ;  $Y_2; \dots; Y_k$  are independent. Let  $M_k$  denote  $\delta \frac{1}{k} \sum_{i=1}^k Y_i$ . Equation (18) gives the mean and variance of  $Y_i$ :

$$E \frac{1}{2} Y_i = \frac{1}{2} m; D \delta Y_i \leq \frac{1}{2} c_i s^2 = n_i; \quad (18)$$

According to Liapunov's central limit theorem, when  $k$  is large, we have the following result:

$$P \left( \frac{r}{1:645} \leq Y_i \leq \frac{r}{1:645} \right) < m < M_k \leq \frac{r}{1:645} \quad (19)$$

A threshold  $d_1$  on  $c_i = n_i$  is set with (19).  $i$ 's trusting belief is taken into consideration only if  $c_i = n_i \leq d_1$ . An interval enclosing  $m$  with at least 90 percent confidence coefficient can be constructed from (19). The intervals corresponding to different conditions are provided in (20a) and (20b):

$$P \left( \frac{r}{1:645} \leq Y_i \leq \frac{r}{1:645} \right) < m < M_k \leq \frac{r}{1:645} \quad (20a)$$

$$\text{Let } X_i \text{ denote the random variable for } s^{\wedge}_i. \text{ Let } Y_i = \frac{1}{k} \sum_{i=1}^k X_i. \text{ Here, } n_i \text{ and } c_i \text{ are constants. Let } S_k \text{ denote } \frac{1}{k} \sum_{i=1}^k Y_i. \text{ We have} \quad (20b)$$

$$P \left( \frac{r}{1:645} \leq Y_i \leq \frac{r}{1:645} \right) < s^2 < S_k \leq \frac{r}{1:645} \quad (21)$$

Let  $r_{\min}$  denote  $\frac{1}{2} \min \{n_1; n_2; \dots; n_k\}$ . We can get a simplified bound from (21):

$$P \left( \frac{r_{\min}}{1:645} \leq Y_i \leq \frac{r_{\min}}{1:645} \right) < s^2 < S_k \leq \frac{r_{\min}}{1:645} \quad (22)$$

with at least 90 percent confidence coefficient can be constructed from (22). Particularly,  $d_2 \geq 2$  leads to the following bound:

$$P \frac{S_k}{1} < s^2 < \frac{S_k}{1} \quad 0.9:(23)$$

The aforementioned estimators for  $m$  and  $s^2$  use  $Dm_i$  and  $c_i$  that are unknown. Two methods to estimate them are discussed in the rest of this section.

4.2 Estimation of  $Dm_i$  and  $c_i$  Based on Previous Knowledge

Two trusters become acquainted if they share a set of commonly rated trustees. It is assumed that a truster uses the consistent rating criteria for all trustees.  $Dm_i$  and  $c_i$  are estimated by comparing the trusting beliefs about trustees known by both  $t$  and  $t_i$ .  $Dm_i$  and  $c_i$  are computed using (15a) and (15b). This approach is named as competence reputation evaluation based on knowledge (CRE-K). The pre-requisite of CRE-K is that the reputation requester has a set of commonly rated trustees with each of the trusters who provide the trusting beliefs.

Suppose  $t$  is the truster who requests information. We want to evaluate  $Dm_i$  for truster  $t_i$ . Let  $u_1, u_2, \dots, u_n$  be the trustees about whom both  $t$  and  $t_i$  submit trusting beliefs.  $m_{t|u_1}, m_{t|u_2}, \dots, m_{t|u_n}$  and  $m_{t_i|u_1}, m_{t_i|u_2}, \dots, m_{t_i|u_n}$  denote the competence trusting beliefs from  $t$  and  $t_i$  respectively. Plugging  $m_{t|u_j}$  and  $m_{t_i|u_j}$  into (15a) yields  $n$  equations:

$$m_{t|u_1} \approx m_{t_i|u_1} + Dm_i; \dots; m_{t|u_n} \approx m_{t_i|u_n} + Dm_i \quad (24)$$

$Dm_i$  is the only unknown in above equation array. We find the  $Dm_i$  that minimizes the sum of square errors in (25):

$$Dm_i \approx \frac{1}{n} \sum_{j=1}^n (m_{t|u_j} - m_{t_i|u_j}) \quad (25)$$

Similarly, we construct  $n$  equations where  $c_i$  is the only unknown and find the  $c_i$  minimizing the sum of square errors in (26):

$$c_i \approx \frac{1}{n} \sum_{j=1}^n (s_{t|u_j}^2 - s_{t_i|u_j}^2) \quad (26)$$

If this method is used, we will use the first estimator for  $s^2$ .

estimator for  $m$  
$$\frac{1}{k} \sum_{i=1}^k \frac{1}{P} \sum_{j=1}^n \frac{\delta_{m_{t|u_j}}}{k} \frac{m_{t|u_j}}{P} \quad (27a)$$

The method discussed requires truster  $t$  to have a lot of acquaintances in the truster set. If this requirement is not

TABLE 5  
Hypothesis Test to Choose a Delegator

|             | Test statistic                                         | Rejection condition                         |
|-------------|--------------------------------------------------------|---------------------------------------------|
| $k \geq 45$ | $z = \frac{\bar{\mu}_{diff}}{s_{\mu_{diff}}/\sqrt{k}}$ | $z > z_{0.05}$ or $z < -z_{0.05}$           |
| $k < 45$    | $t = \frac{\bar{\mu}_{diff}}{s_{\mu_{diff}}/\sqrt{k}}$ | $t > z_{0.05}(k-1)$ or $t < -z_{0.05}(k-1)$ |

satisfied, we can enlarge  $t$ 's acquaintance set using the idea of delegation.  $t$  appoints some trusters  $t_d$  he knows as delegators and uses  $m_{t_d|t_i}$  and  $c_{t_d|t_i}$  as  $m_{t|t_i}$  and  $c_{t|t_i}$ . A delegator of  $t$  shall satisfy two constraints: (1)  $Dm_{t_d|t_d} \geq 0$ , and (2)  $c_{t_d|t_d} \geq 1$ . The method of hypothesis testing is used to check whether  $t$  shall choose a truster  $t_i$  as a delegator.

First, we will test the hypothesis related to  $Dm_{t_d|t_i}$  based on the data set  $m_{t_d|u_1}, m_{t_d|u_2}, \dots, m_{t_d|u_k}$  and  $m_{t_i|u_1}, m_{t_i|u_2}, \dots, m_{t_i|u_k}$ . Here,  $u_1, u_2, \dots, u_k$  are the trustees evaluated by both  $t$  and  $t_i$ . The mean and variance of the data set are computed as follows:

$$\bar{m}_{diff} = \frac{1}{k} \sum_{j=1}^k \frac{m_{t_d|u_j} - m_{t_i|u_j}}{diff} \quad (28a)$$

The following null and alternative hypothesis is used:

$$H_1: \text{Two tailed test : } Dm_{t_d|t_i} \geq 0$$
$$H_0: Dm_{t_d|t_i} < 0$$

The test statistic and rejection condition are summarized in Table 5. If  $t_i$  and  $t$  evaluate a lot of trustees together, we use  $Z$  value. Otherwise,  $t$  value is used.

Second, we will test the hypothesis related to  $c_{t_d|t_i} \geq 1$  in the same way. The data set is  $m_{t_d|u_1}, \dots, m_{t_d|u_k}$  and  $m_{t_i|u_1}, \dots, m_{t_i|u_k}$ . Those trusters who pass both tests are selected as delegators that are introduced to broaden the applicability of CRE-K.

4.3 Estimation Based on Priori Assumptions

The second method to estimate  $Dm_i$  and  $c_i$  is based on priori assumptions about the distribution of trusters. This method uses the second estimator of  $s^2$ , i.e., (17b). Instead of estimating each  $Dm_i$  and  $c_i$ , this method estimates  $E\{\bar{\delta}\}$  to substitute  $\delta_{m_{t|u_j}}$  and  $c_i$  in (16) and (17b). This approach is named as competence reputation evaluation based on assumption (CRE-A).

To estimate  $E\{\bar{\delta}\}$ , we assume:

- 1.  $Dm_i$ 's are independently drawn from the same distribution.
- 2. All states are equally preferable (i.e., the principle of insufficient reasoning).

The first assumption states that  $Dm_1; Dm_2; \dots; Dm_k$  are independent and identically distributed. According to the second assumption, given  $m \in [a, b]$ ,  $Dm_i$  has a uniform distributional expectation in Equation (29).

$$E[Dm_i] = \frac{a+b}{2} \quad (29)$$

Based on the same assumption,  $m$  has the same probability to assume any value in  $[0, 1]$ :

$$E[m] = \int_0^1 x \cdot 1 dx = \frac{1}{2} \quad (30)$$

To estimate  $c_i$ , we assume:

- $c_i$ 's are independently drawn from the same distribution.
- It is equally likely that  $s_i^c = a$  and  $s_i^c = b$ .

Let  $Y_k = \ln \prod_{i=1}^k c_i$ . The second assumption states

$c_i$  is symmetric and centered at 0, i.e.,  $E[\ln c_i] = 0$ . According to the strong law of large numbers theorem, we have

$$P(\lim_{k \rightarrow \infty} \frac{1}{k} \sum_{i=1}^k \ln c_i = 0) = 1 \quad (31a)$$

$$\text{estimator for } m = \frac{1}{k} \sum_{i=1}^k \frac{m_i^A}{2} \quad (31b)$$

$$\text{estimator for } s = \frac{1}{k} \sum_{i=1}^k \frac{s_i^A}{2}$$

## 5 INTEGRITY BELIEF

Integrity may change fast with time. Furthermore, it possesses a meaningful trend. Evaluation of integrity belief is based on two assumptions:

We assume integrity of a trustee is consistent in all contexts.

Integrity belief may vary largely with time. An example is a user behaving well until he reaches a high trust value and then starts committing fraud.

We used mean as an estimator for competence belief as it is relatively steady with time. For integrity belief, this assumption is excluded. When behavior patterns are present, the mean is no more a good estimator. The similarity between a rating sequence and a time series inspires us to adopt the method of double exponential smoothing [7] to predict the next rating based on a previous rating sequence.

Let  $r_i$  denote the  $i$ th rating and  $f_{ip1}$  denote the forecast value of  $r_{ip1}$  after observing the rating sequence  $r_1; \dots; r_i$ . Equation (32) shows how to compute  $f_{ip1}$ :

$$S_i = \frac{1}{2} (a + r_{i-2} + r_{i-1} + r_i) = \frac{3}{4} (r_{i-1} + r_i) \quad (32a)$$

$$b_i = \frac{1}{4} (b + r_{i-1} - r_i) \quad (32b)$$

$$f_{ip1} = \frac{1}{4} (S_i + b_i) \quad (32c)$$

In the above equation array, (32a) computes the overall exponential smoothing, (32b) is the trend smoothing, and (32c) is the prediction after observing  $r_1; \dots; r_i$ . The reason that we use the average of  $r_{i-2}; r_{i-1}; r_i$  is to make the model tolerant to random noise. We use the initial condition defined in (33):

$$S_1 = \frac{1}{4} r_1; b_1 = \frac{1}{4} (r_2 - r_1) \quad (33)$$

Equations in (32) may generate results out of the range  $(0, 1)$ . In this case, we resort to a single exponential

smoothing (SES) to bound double exponential smoothing. We call it BDES-S.

There are two parameters in (32),  $a$  and  $b$ . Both values are in the range  $(0, 1)$ . Their initial values are set by a trustor.  $a$  and  $b$  are updated after the ratings cumulate. They are

between the rating sequence and predictions are minimized. Given a rating sequence  $r_1; r_2; \dots; r_n$  we can determine  $a$  and  $b$  that minimize the MSE between two sequences using non-linear optimization algorithms such as

tional complexity, we approximate  $a$  and  $b$  using a simplified procedure:

The parameters are updated every time a new rating is added;

Only the latest sequence with length 4 is used to update the parameters;

We find the best parameters with a range between 0.1 and 0.9 precise in 1 decimal place.

We find  $a_i$  and  $b_i$  by minimizing by solving the optimization problem in (34):

$$MSE_i = \min_{a, b} \sum_{j=1}^4 (r_j - f_{jp1})^2 \quad (34)$$

Let  $t_i$  denote the trusting belief in integrity after observing  $i$  ratings. Let  $p_i$  denote the predictability associated with  $t_i$ . When  $i \geq 2$ ,  $t_i$  and  $p_i$  are defined in (35a)-(35b):

$$t_i = f_{ip1} \quad (35a)$$

$$p_i = \frac{1}{4} \sum_{j=1}^4 \frac{r_j - r_{j-1}}{r_j} \quad (35b)$$

The  $\delta_{ip1}$ th prediction computed using (32c) is used as the  $i$ th integrity trusting belief. Predictability characterizes the confidence in the belief (i.e., prediction). Mean of squared errors between predictions and ratings are used as  $p_i$ . The lower the value of  $p_i$ , the higher the confidence is. Like  $t_i$ ,  $p_i$  is computed using a recursive formula.

To evaluate  $t_i$  and  $p_i$ , the latest four ratings instead of the whole rating sequence are needed. In addition, we have to store two  $S$  values (i.e.,  $S_i$  and  $S_{i-1}$ ) and two  $b$  values (i.e.,  $b_i$  and  $b_{i-1}$ ). A time threshold is set by a trustor. Any  $S_i$ ,  $b_i$  and ratings before that time are discarded.

TABLE 6  
Algorithms to Build Integrity Belief

|         | Equation                                                       | Initial condition                           | Boundary                                                                                                  |
|---------|----------------------------------------------------------------|---------------------------------------------|-----------------------------------------------------------------------------------------------------------|
| Average | $t_i = \frac{\sum_{k=1}^i r_k}{i}$                             | $t_1 = c$                                   |                                                                                                           |
| SES     | $t_i = \alpha r_i + (1 - \alpha)t_{i-1}$<br>$\alpha \in (0,1)$ | $t_1 = c$                                   |                                                                                                           |
| Regret  | $t_i = \frac{\sum_{k=1}^i w(k,i)r_k}{\sum_{k=1}^i w(k,i)}$     | $t_1 = c$                                   |                                                                                                           |
| BDES    | Equation 32                                                    | $t_1 = c$<br>$S_1 = c$<br>$b_1 = r_2 - r_1$ | $t_i = \alpha \frac{r_{i-2} + r_{i-1} + r_i}{3} + (1 - \alpha)t_{i-1}$<br>if $t_i \geq 1$ or $t_i \leq 0$ |

Each time a new rating is added, the trusting belief and predictability are reevaluated. The procedure to evaluate  $t_j$  and  $p_j$  when  $r_1, r_2, \dots, r_i$  are available is outlined as follows:

1. Compute  $a_j$  and  $b_j$  using (34).
2. Compute  $S_j$ ,  $b_j$ ,  $f_{j|1}$  using (32).
3. Compute  $t_j$  and  $p_{j|1}$  using (35).

If  $i < 4$ , we ignore the step of updating  $a_j$  and  $b_j$  and use the initial parameters instead.

The aforementioned approach is extended to evaluate reputation. Let  $L_t$  denote the length of a time interval. Integrity trusting beliefs from different trusters are sorted in an ascending order based on time-stamp. The direct evaluation algorithm is applied on this sequence. The prediction generated is the evaluated integrity reputation.

## 6 EXPERIMENTAL STUDY OF TRUST MODEL

Experimental studies were conducted to evaluate the integrity belief model proposed in Section 5. The objective is to identify the suitable approaches for various scenarios (different types of trustees) and obtain guidelines to determine the appropriate values of parameters for the algorithms. Sections 6.1, 6.2 and 6.3 evaluate the approaches to build integrity belief based on direct experience.

We also conducted experiments to evaluate the competence belief model introduced in Section 4. The CRE-A and CRE-K methods were evaluated under different scenarios, with trustee behavior generated using a normal distribution. Experiments were conducted to compare the true mean and variance with the estimated mean and variance of competence reputation for different number of trusters. The relative error (RE) of CRE-A was found to be around 5 percent, and that of CRE-K was less than 3.5 percent, which are promising results. We omit detailed experiment results due to space constraints.

### 6.1 Study on Integrity Belief Building Methods

In this section, the BDES algorithm is compared with three other algorithms for five trustee behavior patterns.

Algorithms compared. The algorithms compared are BDES, simple average, single exponential smoothing, and the time-weighted average, called REGRET, proposed in [20]. Let  $t_i$  denote the trusting belief after observing rating

TABLE 7  
Parameters Used in Experiments

| SES                            | Regret                                                  | BDES                                                            |
|--------------------------------|---------------------------------------------------------|-----------------------------------------------------------------|
| $\alpha = 0.3$ (initial value) | $w(k,i)$ is a function linearly decreasing with $(i-k)$ | $\alpha = 0.3$ (initial value)<br>$\beta = 0.7$ (initial value) |

sequence  $r_1; r_2; \dots; r_i$ . Table 6 summarizes how  $t_j$  is evaluated under the four algorithms.  $w(k, i)$  in REGRET is a time dependent function giving higher values to ratings temporally close to  $r_i$ . Table 7 shows the initial values of the parameters of BDES and SES. A function linearly decreasing with  $(i - k)$  is used as  $w(k, i)$  in REGRET.

Experiment setup. For the experiments discussed in Sections 6.2 and 6.3 below, trustee behavior was simulated using the five different integrity rating generation functions detailed below. A rating for trustee  $u$  generated by a behavior pattern function at time  $i$  is considered to be the true integrity rating submitted for  $u$  by a truster  $t$  at time point  $i$ . For each behavior pattern experiment, a sequence of 100 ratings for each trustee were generated using the pattern function and the performances of the four integrity belief building methods listed above were evaluated by measuring the difference between the true rating and the rating output by the integrity belief method at each point in the sequence. Note that the identity of the trusters is not relevant in this case: The 100 ratings for a trustee could be submitted by a single truster or by 100 different trusters.

Generate ratings based on trustee behavior patterns. The true values of integrity trusting belief about a trustee can be viewed as the range of a time dependent function  $f(i)$ . A pattern is a family of  $f(i)$ s with the same form. It is impossible and unnecessary to enumerate all possible forms of  $f(i)$ s. We are interested in meaningful patterns revealing the trend and intention of a trustee's behavior. Five types of patterns, random trustee, stable trustee, trend trustee, jump-ing trustee and two-phase trustee, are identified and used in the experiments. The random pattern shows that the trustee's behavior is variable. Prediction based on previous knowledge may not lead to good results. On the other hand, we can expect to precisely predict the next performance of a trustee with a stable pattern. The trend pattern captures the improving or deteriorating behavior pattern. The jumping and two-phase patterns indicate a sudden shape change. Usually, they imply misbehaving of trust builders.

TABLE 8  
Trustee Behavior Patterns

|           | Form of $f(i)$                                                                                                                           | Figure |
|-----------|------------------------------------------------------------------------------------------------------------------------------------------|--------|
| Random    | $f(i) = U(0,1)$ for $\forall i$                                                                                                          | Fig. 5 |
| Stable    | $f(i) = c_1$ for $\forall i$                                                                                                             | Fig. 6 |
| Trend     | $f(i) = c_1 + ic_2$ for $\forall i$                                                                                                      | Fig. 7 |
| Jumping   | $f(i) = c_1$ for $i \leq n_0$<br>$f(i) = c_2$ otherwise                                                                                  | Fig. 8 |
| Two-phase | $f(i) = c_1$ if $i \leq n_0$<br>$f(i) = c_1 - \frac{(c_1 - c_2)(i - n_0)}{n_1 - n_0}$ if $n_0 < i \leq n_1$<br>$f(i) = c_2$ if $n_1 < i$ | Fig. 9 |

TABLE 9  
Instances of Trustee Behavior Patterns

|           | Definition of $f(i)$                                                                                             | Figure |
|-----------|------------------------------------------------------------------------------------------------------------------|--------|
| Random    | $f(i)=U(0,1)$ for $i \in [1,100]$                                                                                | Fig. 5 |
| Stable    | $f(i)=0.6$ for $i \in [1,100]$                                                                                   | Fig. 6 |
| Trend     | $f(i)=0.3+0.005i$ for $i \in [1,100]$                                                                            | Fig. 7 |
| Jumping   | $f(i)=0.8$ if $i \leq 50$<br>$f(i)=0.3$ if $50 < i \leq 100$                                                     | Fig. 8 |
| Two-phase | $f(i)=0.8$ if $i \leq 40$<br>$f(i) = 0.8 - 0.025(i - 40)$ if $40 < i \leq 60$<br>$f(i)=0.3$ if $60 < i \leq 100$ | Fig. 9 |

Table 8 shows the form of  $f(i)$  for each pattern. In this work, the independent variable of  $f(i)$  is the number of interactions.  $c_i$ s and  $n_i$ s are constants. Based on the behavior patterns, we can systematically evaluate a belief formation algorithm. The effectiveness of an algorithm in an environment is determined by (1) how the algorithm performs for each type of trustee, and (2) what is the distribution of trustees belonging to each?

In this section, we study the first issue. Algorithms are evaluated against the interaction sequences representing different trustee behaviors. Each interaction sequence is generated to reflect certain trustee behavior patterns. A trustee's behavior is determined by her trustworthiness and is influenced by some unpredictable factors. Therefore, the  $i$ th rating is generated using (36). The  $i$ th rating falls into  $[0; 1]$  with probability 90 percent. The interval is interpreted as the region where relative error is smaller than 10 percent:

$$N(f(i) \in [0; 1]) = 1 - 0.1 = 0.9 \quad (36)$$

## 6.2 Distribution of Errors

The first set of experiments compares absolute error (AE) and relative error, as defined in (37a) and (37b) respectively, of the four algorithms. We choose this measurement because the purpose of evaluating  $t_i$  is to forecast  $r_{ip1}$ , i.e., a good trust building algorithm shall output good predictions. Absolute and relative errors characterize how close one prediction is to the true value:

$$AE = |t_i - r_{ip1}| \quad (37a)$$

$$RE = |t_i - r_{ip1}| / r_{ip1} \quad (37b)$$

We generate 100 ratings for each type of behavior pattern. The parameters are summarized in Table 9. Four algorithms are applied on each trustee. The absolute and relative error for each prediction is computed. The distribution of errors generated by each algorithm is plotted using cumulative frequency figures.

### 6.2.1 Results and Observations

**A trustee with random behavior pattern.** For a trustee who has the random behavior pattern, the next behavior has no relation to the previous behaviors. The rating can increase or decrease sharply at any time. Because the behavior of the trustee is completely unpredictable, none of the evaluated algorithms is able to provide a good prediction of how the

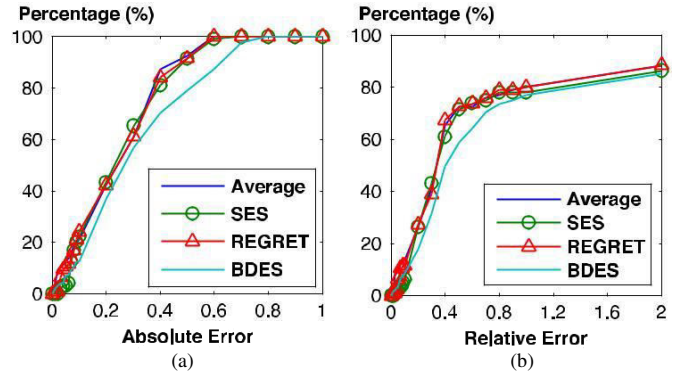


Fig. 5. (a) Absolute error, and (b) Relative error for a trustee with the random behavior pattern.

next behavior will be. The Average, SES, and REGRET algorithms have almost the same performance in terms of absolute error, as shown in Fig. 5a. The Average algorithm performs slightly better than the other two. About 88 percent of its results have an absolute error less than 0.4, while the percentages of the SES and REGRET algorithms are 85 and 81 percent respectively. Nearly all results of these three algorithms have an absolute error less than 0.6. The BDES algorithm fails to achieve low error rate in this experiment. Only 70 percent of its results have an absolute error less than 0.4. The upper bound of the error is 0.8 instead of 0.6. Fig. 5b shows that all algorithms generate large relative errors. For the Average, SES, REGRET, and BDES algorithms, the percentages of the results that have a relative error less than 100 percent are respectively, 80, 78, 80, and 77 percent. The percentages of the results that have a relative error greater than 200 percent are 12, 14, 12, and 15 percent respectively. The Average and REGRET algorithms perform the best.

**A trustee with stable behavior pattern.** For a trustee who has the stable behavior pattern, the next behavior tends to fluctuate around the mean of the previous behaviors. The rating has a greater probability of being closer to the mean of the previous ratings. All algorithms are able to produce very good results in terms of absolute error and relative error as shown in Figs. 6a and 6b. The REGRET algorithm performs the best, slightly better than the BDES algorithm. For these two algorithms, around 98 percent of the results have an absolute error that is less than 0.2. The corresponding

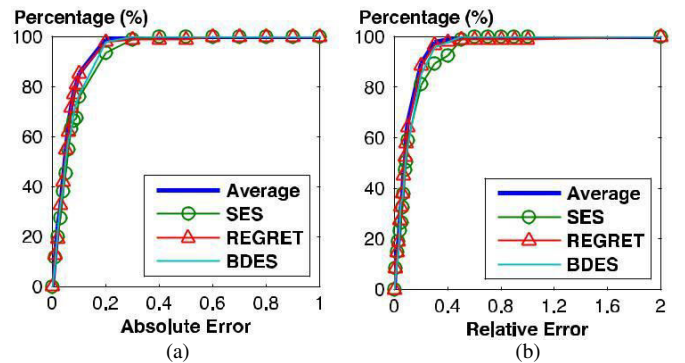


Fig. 6. (a) Absolute error, and (b) Relative error for a trustee with the stable behavior pattern.

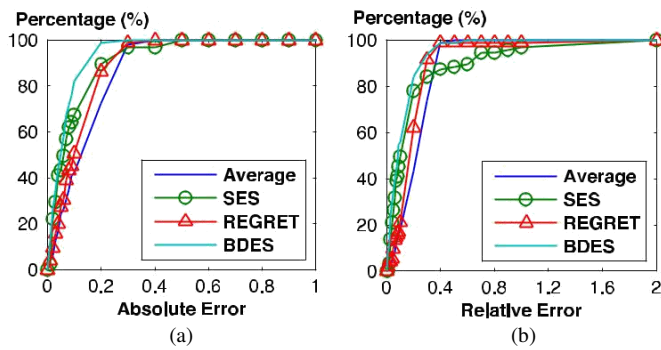


Fig. 7. (a) Absolute error, and (b) Relative error for a trustee with the trend behavior pattern.

percentage for the SES algorithm is 94 percent. The percentages of the ratings that have less than 0.1 absolute error are 86 percent for the Average and REGRET algorithms, and 78 percent for the SES and BDES algorithms. As shown in Fig. 6b, for every algorithm, almost all results have a relative error less than 60 percent. Ninety percent of the results generated by the Average and REGRET algorithms have a relative error less than 20 percent. The corresponding percentages for the SES and BDES algorithms are 82 and 87 percent respectively.

A trustee with a trend behavior pattern. For a trustee who has the trend behavior pattern, the behavior becomes better and better (or worse and worse depending on the trend) as time passes, i.e., the number of interactions increases. The BDES algorithm outperforms the other algorithms in terms of absolute and relative error when the trustee has a trend behavior pattern.

As shown in Fig. 7a, 88 percent of its results have an absolute error less than 0.2 and 83 percent of its results have an absolute error less than 0.1. The corresponding percentages are 72 and 42 percent for the Average algorithm, 89 and 67 percent for the SES algorithm, and 87 and 50 percent for the REGRET algorithm. As shown in Fig. 7b, although the percentage of the results that have less than 40 percent relative error is around 98 percent for the BDES, Average and REGRET algorithms, only the BDES algorithm is able to make 85 percent of its results having a relative error less than 20 percent. The Average and REGRET algorithms can achieve 42 and 61 percent respectively. 87 percent of the results obtained using the SES algorithm have less than 40 percent relative error, 78 percent of the results have less than 20 percent relative error.

A trustee with jumping behavior pattern. A trustee with the jumping behavior pattern behaves as if he had the stable behavior pattern, and suddenly changes his behaviors. Comparing the results of this experiment with those of the previous two experiments, we can see that the performance downgrades for all, especially the Average and REGRET algorithms. As shown in Fig. 8a, the BDES and SES algorithms still make, respectively, 93 and 88 percent of the results have less than 0.2 absolute error. The corresponding percentage is 48 percent for the Average algorithm and 61 percent for the REGRET algorithm. The upper bound of the absolute error is 0.6 for the BDES and SES, and 0.7 and 0.9 for the Average and REGRET algorithms respectively. Fig. 8b shows that the BDES algorithm has the highest

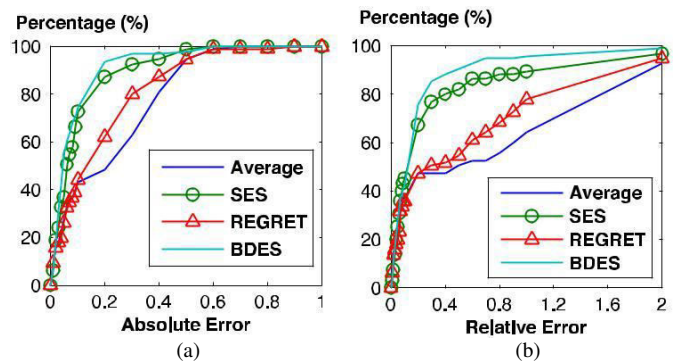


Fig. 8. (a) Absolute error, and (b) Relative error for a trustee with the jumping behavior pattern.

percentage of the results with less than 100 percent relative error, which is 96 percent. For the Average, REGRET, and SES algorithms, the percentages are, respectively, 63, 78 and 90 percent. From another perspective, 90 percent of the results obtained using the BDES algorithm have less than 47 percent relative error. The same percentage of results obtained using the Average, REGRET, and SES algorithms have a relative error less than 190, 170, and 100 percent respectively. The BDES algorithm has the best performance among the evaluated algorithms.

A trustee with two-phase behavior pattern. A trustee who has the two-phase behavior pattern has similar behaviors as compared to the trustee with the jumping behavior pattern, except that he changes his behaviors gradually instead of suddenly. In terms of absolute errors, the BDES and SES algorithms perform a little better, while the Average and REGRET algorithms perform slightly worse as compared to the jumping behavior pattern. As shown in Fig. 9a, the percentages of the results with less than 0.2 absolute errors are 85, 90, 62 and 47 percent for the BDES, SES, REGRET, and Average algorithms, respectively. The percentages of the results with less than 0.1 absolute error are 82, 69, 37, and 37 percent correspondingly. Fig. 9b shows that all algorithms have a better performance in terms of relative errors compared to what they achieve in the previous experiment. All the results obtained using the BDES algorithm have less than 100 percent relative error. For the SES, REGRET, and Average algorithms, 98, 83, and 71 percent of the results, respectively, have a relative error less than 100 percent. Ninety percent of the results obtained from the BDES, SES,

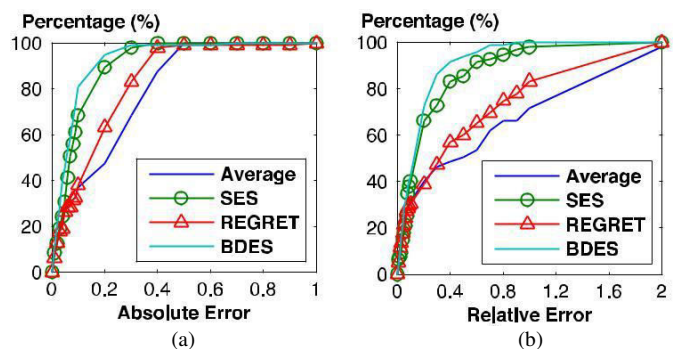


Fig. 9. (a) Absolute error, and (b) Relative error for a trustee with the two-phase behavior pattern.

TABLE 10  
Average MSE for Each Behavior Pattern

|         | Random   | Stable    | Trend     | Jumping  | Two-Phase |
|---------|----------|-----------|-----------|----------|-----------|
| Average | 0.086069 | 0.0037669 | 0.026811  | 0.072545 | 0.063554  |
| SES     | 0.093301 | 0.007613  | 0.014744  | 0.021522 | 0.014272  |
| REGRET  | 0.08932  | 0.0075901 | 0.016907  | 0.055423 | 0.046255  |
| BDES    | 0.12558  | 0.0056795 | 0.0062433 | 0.012282 | 0.0074293 |

Around 99 percent of them are less than 0.005. The REGRET and SES algorithms have almost the same performance, which is worse than that of the BDES algorithm. Ninety percent of the MSEs of the REGRET and SES algorithms are less than 0.009, while the same percentage of MSEs of the BDES algorithm are less than 0.0078. As shown in Fig. 10b the BDES algorithm introduces larger MSE than the other three algorithms when the trustee has the random behavior pattern. The MSEs range from 0.07 to 0.20. Ninety percent of them are less than 0.16. The MSEs of the other three algorithms are very close. All of them are in the range of 0.06 to 0.12. Fig. 10c shows that the BDES performs better than the other algorithms in terms of introducing less MSE when the trustee has the trend behavior pattern. Its smallest MSE is about 0.005. Ninety nine percent of its MSEs are less than 0.012, which is the smallest one among all the MSEs introduced by the other algorithms. The Average algorithm has the worst performance. Its MSEs are in the range of 0.02 to 0.04, 94 percent of them are less than 0.03. The SES algorithm performs slightly better than the REGRET algorithm. Its MSEs range from 0.012 to 0.018, while 99 percent of the MSEs of REGRET are in the range of 0.014 to 0.02.

As shown in Figs. 10d and 10e, when the trustee has the jumping or two-phase behavior pattern, the BDES algorithm has much better performance than the other algorithms. Even its largest MSE is smaller than the smallest one introduced by the other algorithms. For a trustee with the jumping behavior pattern, the ranges of the MSEs are 0.009 to 0.017 for the BDES algorithm, 0.018 to 0.03 for the SES algorithm, 0.04 to 0.07 for the REGRET algorithm, and 0.06 to 0.09 for the Average algorithm. For a trustee with the two-phase behavior pattern, the corresponding ranges are 0.004 to 0.001, 0.012 to 0.017, 0.04 to 0.06, and 0.05 to 0.08, respectively.

Table 10 shows the average MSE for each behavior pattern obtained in the experiment. For a completely unpredictable trustee, i.e., the one with random behavior, no algorithm is able to provide practically useful integrity trusting belief. For a stable trustee, all algorithms can provide satisfactory information, with the Average algorithm being the best. When a trustee has the trend to change his behavior, e.g., the trend, jumping, and two-phase behavior pattern, only the BDES algorithm is able to catch this trend. The accuracy of the integrity trusting belief computed using the BDES algorithm is not affected much by the change of behavior, as seen in the last row.

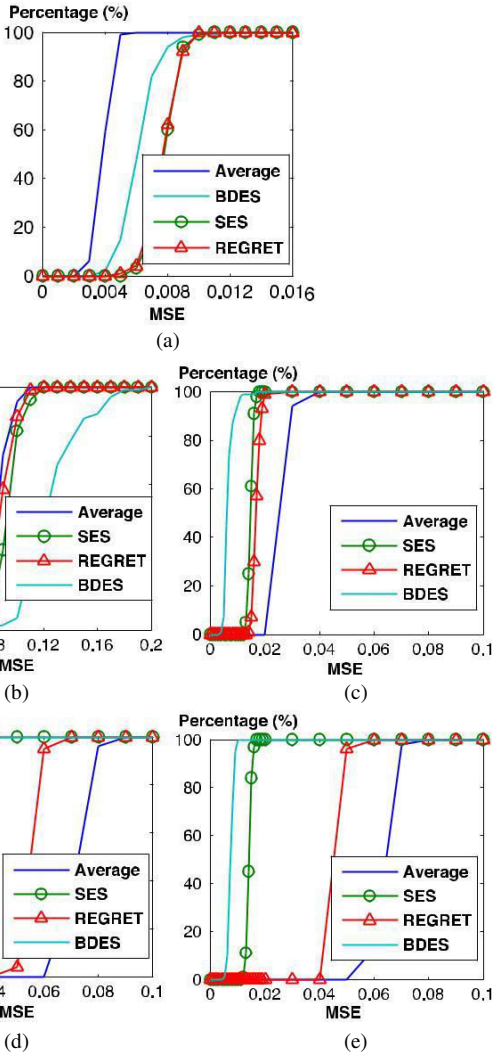


Fig. 10. Distribution of mean squared error for (a) Stable, (b) Random, (c) Trend, (d) Jumping, and (e) Two-phase behavior patterns.

REGRET, and Average algorithms have less than 38, 58, 140, and 170 percent relative errors respectively.

### 6.3 Distribution of Mean Squared Errors

Previous experiments studied the errors generated by a single user per type. The second set of experiments explores the distribution of mean squared errors, as defined in (38). Here,  $n$  is the number of predictions:

$$MSE \approx \frac{1}{n} \sum_{i=1}^n r_i^2 \quad (38)$$

We generate the 100 trustees per behavior pattern using the parameters in Table 9. Please note the trustees are not the same. Four algorithms are applied on each trustee. For each trustee, the MSE generated by each algorithm is computed. For each type of trustee, the distribution of MSE generated by each algorithm is plotted using cumulative frequency figures. Also, average MSE is computed.

#### 6.3.1 Results and Observations

When the trustee has the stable behavior pattern, the Average algorithm outperforms the other algorithms in terms of MSE as shown in Fig. 10a. Its MSEs range from 0.002 to 0.01.

findings from social science, and is not limited to trusting belief as most computational methods are. We presented a representation of context and functions that relate different contexts, enabling building of trusting belief using cross-context information.

The proposed dynamic trust model enables automated trust management that mimics trusting behaviors in society, such as selecting a corporate partner, forming a coalition, or choosing negotiation protocols or strategies in e-commerce. The formalization of trust helps in designing algorithms to choose reliable resources in peer-to-peer systems, developing secure protocols for ad hoc networks and detecting deceptive agents in a virtual community. Experiments in a simulated trust environment show that the proposed integrity trust model performs better than other major trust models in predicting the behavior of users whose actions change based on certain patterns over time.

## REFERENCES

- [1] G.R. Barnes and P.B. Cerrito, "A Mathematical Model for Interpersonal Relationships in Social Networks," *Social Networks*, vol. 20, no. 2, pp. 179-196, 1998.
- [2] R. Brent, *Algorithms for Minimization without Derivatives*. Prentice-Hall, 1973.
- [3] A. Das and M.M. Islam, "SecuredTrust: A Dynamic Trust Computation Model for Secured Communication in Multiagent Systems," *IEEE Trans. Dependable and Secure Computing*, vol. 9, no. 2, pp. 261-274, Mar./Apr. 2012.
- [4] C. Dellarocas, "Immunizing Online Reputation Reporting Systems against Unfair Ratings and Discriminatory Behavior," *Proc. Second ACM Conf. Electronic Commerce*, pp. 150-157, 2000.
- [5] L. Fan, "A Grid Authorization Mechanism with Dynamic Role Based on Trust Model," *J. Computational Information Systems*, vol. 8, no. 12, pp. 5077-5084, 2012.
- [6] T. Grandison and M. Sloman, "A Survey of Trust in Internet Applications," *IEEE Comm. Surveys*, vol. 3, no. 4, pp. 2-16, Fourth Quarter 2000.
- [7] J.D. Hamilton, *Time Series Analysis*. Princeton University Press, 1994.
- [8] J. Hu, Q. Wu, and B. Zhou, "FCTrust: A Robust and Efficient Feedback Credibility-Based Distributed P2P Trust Model," *Proc. IEEE Ninth Int'l Conf. Young Computer Scientists (ICYCS '08)*, pp. 1963-1968, 2008.
- [9] B. Lang, "A Computational Trust Model for Access Control in P2P," *Science China Information Sciences*, vol. 53, no. 5, pp. 896-910, May 2010.
- [10] C. Liu and L. Liu, "A Trust Evaluation Model for Dynamic Authorization," *Proc. Int'l Conf. Computational Intelligence and Software Eng. (CiSE)*, pp. 1-4, 2010.
- [11] X. Long and J. Joshi, "BaRMS: A Bayesian Reputation Management Approach for P2P Systems," *J. Information & Knowledge Management*, vol. 10, no. 3, pp. 341-349, 2011.
- [12] S. Ma and J. He, "A Multi-Dimension Dynamic Trust Evaluation Model Based on GA," *Proc. Second Int'l Workshop Intelligent Systems and Applications*, pp. 1-4, 2010.
- [13] S. Marsh, "Formalizing Trust as a Concept," PhD dissertation-Dept. of Computer Science and Math., Univ. of Stirling, 1994.
- [14] P. Matt, M. Morge, and F. Toni, "Combining Statistics and Arguments to Compute Trust," *Proc. Ninth Int'l Conf. Autonomous Agents and Multiagent Systems (AAMAS '10)*, pp. 209-216, 2010.
- [15] D. McKnight, V. Choudhury, and C. Kacmar, "Developing and Validating Trust Measures for E-Commerce: An Integrative Topology," *Information Systems Research*, vol. 13, no. 3, pp. 334-359, Sept. 2002.

# TIME ORIENT SPATIAL TRAFFIC PATTERN BASED MITIGATION OF DISTRIBUTED DENIAL OF SERVICE ATTACKS WITH DJN IN DISTRIBUTED WIRELESS NETWORKS

A. Saraswath <sup>1</sup>/Research Scholar & AssistantProfessor/P.G. and Research Department of Computer Science/

Government Arts College (Autonomous),/Karur, Tamil Nadu,India.

Dr.K.Thangadurai<sup>2</sup>/ Research Scholar & AssistantProfessor/P.G. and Research Department of Computer Science/

Government Arts College (Autonomous),/Karur, Tamil Nadu,India.

Dr.A.Mummoorthy<sup>3</sup>/Asscoiate professor/Department of Computer Science & Engineering

Malla Reddy College of engineering & Technology,Secunderabad-500 100. Telangana State/India

[asaraswathipgm@gmail.com](mailto:asaraswathipgm@gmail.com), [amummoorthy@gmail.com](mailto:amummoorthy@gmail.com)

**Abstract** - The research of identifying the presence of distributed jamming network has not been leveraged by the researcher's in modern wireless networks. This paper addresses the problem of distributed jamming network which intended to generate distributed denial of service attacks in the distributed wireless networks. The lower cost of tiny nodes of DJN has the capability of rapid deployment but has greater impact in the depletion of quality of service of distributed wireless networks. To mitigate the problem of DDOS attacks generated by distribute jamming network of nodes, an time orient spatial pattern based traffic analysis and verification approach has been proposed. The base station maintains the spatial pattern which represents the details of network snapshot which has been generated at each time window and maintains the location details of all the nodes present under the coverage. The details of nodes can be exchanged between neighbor base stations and the base station computes the traffic rate at each time window which is available with the time orient spatial traffic pattern. From the available patterns of different time window, the deployment of DJN with large number's or small scale can be identified in efficient manner. The proposed scheme, helps

improving the performance of DWN and reduces the threat of DDOS attacks in all levels.

## **Index Terms:**

**DWN, DDOS attack, DJN, Time Orient Spatial Traffic Patttern, QoS.**

## **Introduction:**

Distributed wireless nework (DWN) has broader sense and collection of wireless nodes with fixed and mobile nodes. There will be base station, which maintains different information about the wireless nodes under the coverage of its own. The wireless nodes comes with the radio to perform transmissioin and reception of packets but has fixed transmission range. The wireless nodes can here the transmission from the nodes which has located within the coverage of the node. Similarly, the nodes under the coverage of a base station, can here the data packets transmitted from the base station.

The purpose of wireless network is to collect data from different location of the geographic region and the wireless network has devices installed in different location of the network and perform data collection through them. Sometimes, the nodes of the network are allowed to move between different locations. The data are collected

thorough various wireless nodes under the coverage of the network and they can perform cooperative transmission of packets to deliver to the control unit.

In any distributed network, there is a huge chance of denial of service attacks where there exist number of services provided by the wireless network. There will be set of malicious node which can perform such attack. In case of denial of service attack, the detection of such malicious node will be very difficult and the detection approach introduce many overhead into the network and spoils the quality of service of the network. The Distributed jamming network is such one which composed of number of malicious nodes, with the intension to degrade the performance of the wireless network or the service being provided.

The distributed jamming network is formed by group of nodes which sends thousands of malicious packets towards the service point or in case of data collection, the nodes replies malicious data to the data collection request. Also the jamming nodes creates links with the other nodes and performs variety of denial of service attacks and destroy the sensor node by depleting the energy of the wireless node.

The time orient spatial traffic pattern can be used in mitigating the denial of service attacks in distributed wireless network, where the pattern represents the different factors of wireless network. The single instance of time orient spatial traffic pattern has the following details namely the time factor  $T_\alpha$  which represents the single time window,  $T_\beta$  – represents the number of nodes present at any time window,  $T_\mu$  shows the number of packets has been received in any time window,  $T_l$  represents the location of any node at particular time window. Using these details the trustworthy of any node can be computed by the nodes of wireless network.

### **Related Works:**

There are many approaches has been discussed earlier for the mitigation of denial of service attack and we explore few of the methods here in this section.

Using channel hopping to increase 802.11 resilience to jamming attacks [1], explore how to protect 802.11 networks from jamming attacks by having the legitimate transmission hop among channels to hide the transmission from the jammer. Using a combination of mathematical analysis and prototype experimentation in an 802.11a environment, we explore how much throughput can be maintained in comparison to the maintainable throughput in a cooperative, jam-free environment. Our experimental and analytical results show that in today's conventional 802.11a networks, we can achieve up to 60% of the original throughput. Our mathematical analysis allows us to extrapolate the throughput that can be maintained when the constraint on the number of orthogonal channels used for both legitimate communication and for jamming is relaxed.

**ARES:** An anti-jamming reinforcement system for 802.11 networks [4], conduct extensive experiments on an indoor 802.11 network to assess the ability of two physical layer functions, rate adaptation and power control, in mitigating jamming. In the presence of a jammer we find that: (a) the use of popular rate adaptation algorithms can significantly degrade network performance and, (b) appropriate tuning of the carrier sensing threshold allows a transmitter to send packets even when being jammed and enables a receiver capture the desired signal. Based on our findings, we build ARES, an Anti-jamming REinforcement System, which tunes the parameters of rate adaptation and power control to improve the performance in the presence of jammers.

ARES ensures that operations under benign conditions are unaffected.

A Measurement-Driven Anti-Jamming System for 802.11 Networks [10], conduct extensive experiments on an indoor 802.11 network to assess the ability of two physical-layer functions, rate adaptation and power control, in mitigating jamming. In the presence of a jammer, we find that: 1) the use of popular rate adaptation algorithms can significantly degrade network performance; and 2) appropriate tuning of the carrier sensing threshold allows a transmitter to send packets even when being jammed and enables a receiver to capture the desired signal. Based on our findings, we build ARES, an Anti-jamming REinforcement System, which tunes the parameters of rate adaptation and power control to improve the performance in the presence of jammers.

Detection Approach for Denial of Service Attack in Dynamic Wireless Networks [11], proposed architecture and mechanism for detection and control of DDOS attacks over reputation and score based MANET. To studied a novel DoS attack perpetrated by JellyFish: relay nodes that stealthily misorder, delay, or periodically drop packets that they are expected to forward, in a way that leads astray end-to-end congestion control protocols. This attack is protocol-compliant and yet has a devastating impact on the throughput of closed-loop flows, such as TCP flows and congestion-controlled UDP flows. For completeness, we have also considered a well-known attack, the Black Hole attack, as its impact on open-loop flows is similar to the effect of JellyFish on closed-loop flows. We studied these attacks in a variety of settings and have provided a quantification of the damage they can inflict. We showed that, perhaps surprisingly, such attacks can actually increase the capacity of

ad hoc networks as they will starve all multihop flows and provide all resources to one-hop flows that cannot be intercepted by JellyFish or Black Holes.

Distributed Detection of DoS Using Clock Values in Wireless Broadband Networks [12], proposes a novel flooding attack, the most severe denial-of-service attack that occurs at the transport layer of the internet. The main objective of this approach is to install local and global monitoring agents at various points in order to monitor and filter real-time TCP traffic and UDP traffic thereby allowing legitimate traffic to flow in the network during attack traffic filtration process and to avoid buffer overflow at the monitoring agents. Also, a novel algorithm has been proposed by taking the clock values of each node into account for effective detection of the attack. This distributed defense mechanism reduces the burden on a single global monitoring agent thereby introducing local monitoring agents at various points in the network.

Explicit Query based Detection and Prevention Techniques for DDOS in MANET [13], study how various detection parameters together work as a single and efficient method to detect various DDOS attacks in Manet. Later in this paper a technique to prevent DDOS attacks in Manet is also presented which help in preventing the attacks to communicate in the network and did not allow them in the network.

Mitigating Denial of Service Attacks in Wireless Networks [16], prevent the cyberspace from DOS attack. Zombie is a computer that has been compromised by the hackers for performing malicious activities. Group of zombie computers involved in jamming activities is called botnets. Denial of Service attack either injects malicious packet or drops legitimate packets from the network. The objective of DOS attack is to

prevent the receiver from reception of legitimate data and to get control over the entire network. This paper proposes a survey on three types of DOS attack such as selective forwarding attack, pollution attack and jamming attack and its detection techniques.

All the above discussed approaches has the problem of identifying denial of service attacks and jamming attacks in efficient manner.

#### Proposed Method:

The proposed time variant spatial traffic pattern approach logs the time variant location of the nodes and for each not it generates the time variant pattern about the traffic it has generated at each time window. Based on the generated traffic pattern, the method identifies the malicious nodes and mitigate the jamming attacks. It has the following functional components as mentioned in the architecture diagram.

#### Time variant Snapshot Generation:

Each node in the network identifies the set of neighbors of the network, and for each of the neighbor identified, the location parameter for each of the node is identified and stored in the trace. This is performed in distributed manner by all the wireless nodes present in the network and by the base station also. The generated trace will be used by both base station as well as the wireless nodes of the network.

#### Procedure:

Input: Network Trace Nt.

Output: Null.

Start

If Node==Base Station

Identify all the nodes present in the coverage of the network.

Node List Nl =  $\sum \text{Nodes@Coverage}(\text{Base Station})$

Else

Node List Nl =  $\sum \text{Nodes@Coverage}(\text{Nodes})$

End.

For each Node Ni

Identify the location of Ni

NLoc = Ni(Loc).

Add to the trace NT =

$\sum \text{Trace}(\text{Nt}) + \{\text{Ni}, \text{NLoc}\}$

End

Stop.

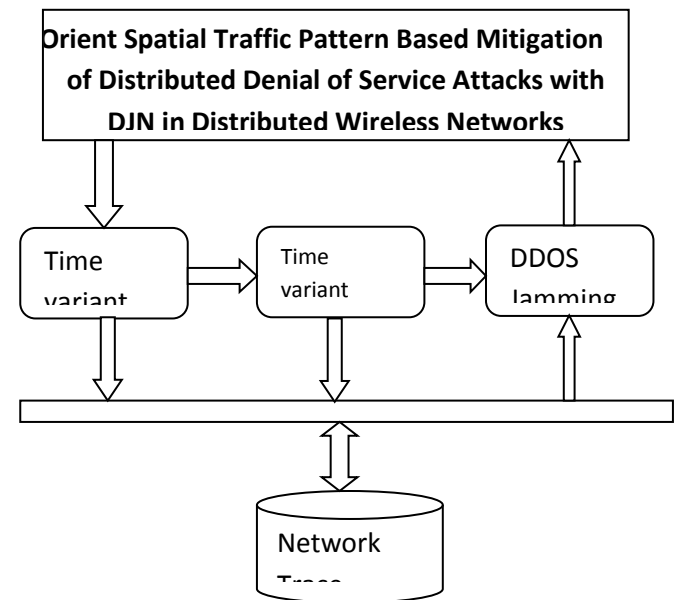


Figure 1: Proposed System Architecture

The Figure 1, shows the functional components of the proposed approach and we will explain each of the functional component in detail in this section.

#### Time Variant Spatial Traffic Pattern Generation:

At this stage, each node or base station generates the traffic pattern for the identified neighbors at the time window. At each time window, the node reads the network trace which has the snapshot at the time window and for each of the neighbor or member being identified, the method computes the number of packets has been received from the neighbor, the payload being sent by the node, the location of the neighbor node and so on. Using all these information, the method generates the pattern which will be used to perform the DDOS attack detection.

Procedure:

Input: Network Trace Nt.

Output: Pattern set Ps.

Start

Read Network trace Nt.

$Nt = \sum \text{Logs}(Nt)@T_{\infty}$

Identify the set of neighbors at the time window  $T_{\infty}$ .

Neighbor List NI =  $\sum \text{Neighbors}(Ni)@T_{\infty}$

For each neighbor Ni from N  
Compute number of packets has been received.

$Tpr = \sum \text{Packets} \in (Nt@T_{\alpha})$

Compute payload PI =  $(\sum \text{payload}(\text{Packets} \in (Nt@T_{\alpha}))/Tpr$

Identify location of the node  
 $NLoc = Loc(Ni)$

Generate Pattern  $Pi = \{ Ni, NLoc, Tpr, PI \}$ .

Add to pattern set  $Ps = \sum \text{patterns}(Ps)+Pi$ .

End

Stop.

DDOS Jamming Attack Detection:

The method performs the mitigation of denial of service attacks performed by distributed jamming attack. Whenever the node has a packet to transmit or received, the node retrieves the traffic pattern belongs to previous time window  $T_{\infty-1}$ , and computes the pattern for the neighbor node. Using all the information the method computes the trust weight for each of the neighbor or particular neighbor. If the trust weight is more than the threshold then the packet will be forwarded otherwise the node will be considered as jammer and ignored. The nodes identity is verified using one step verification with the base station where the location details of the all nodes are stored.

Procedure:

Input: Network Trace Nt, Pattern Set Ps

Output: Null

Start

Receive Packet P.

Identify Source of packet P.

Saddr = Source-Address(P)

Retrieve the Location from network trace.

$NLoc = Nt(Ni(Loc))$ .

Verify the location with the base station Bs.

If true then

Retrieve previous time windows pattern pi.

$Pi = Ps(Ni)@T_{\alpha-1}$

Compute current windows packet details.

$Tpr = \sum \text{packets}(Ni) \in Nt(\text{current Time window})$

Compute average payload  
 $Apl = (\sum \text{payload}(\text{Packets} \in (Nt@T_{\alpha}))/Tpr$

Compute Trust Weight  $Tw = Tpr \times Apl$ .

If  $Tw > TTh$  then //TTH-Trust threshold

Forward packet.

end

Else

Drop the packet or look for another neighbor.

End.

Stop.

Results and Discussion:

The proposed time variant spatial traffic pattern based denial of service attack with distributed jamming network has been implemented using Network simulator NS2. The protocol has been simulated using 100 nodes with 1000 meters range and with number of scenarios. The method has produced efficient results in all the factors of mitigating the denial of service attacks in wireless networks.

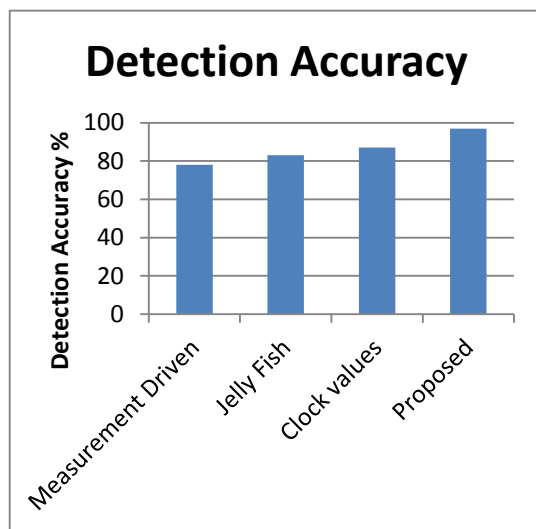
| Parameter       | Value                 |
|-----------------|-----------------------|
| Too             | Network Simulator NS2 |
| Simulation Area | 1000 Meters           |
| Number of nodes | 100                   |

|                          |            |
|--------------------------|------------|
| Nodes Transmission Range | 100 meters |
| Protocol                 | TSTP       |

Table 1: Simulation Details

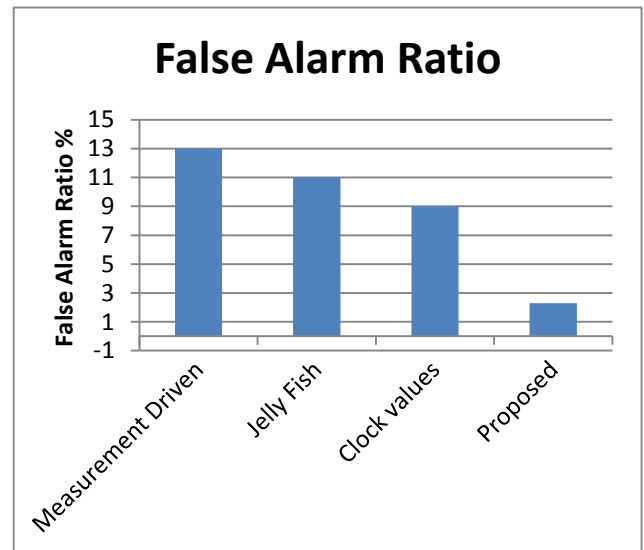
The Table 1, shows the simulation details being used to simulate and evaluate the proposed protocol.

The protocol has been simulated using the network simulator and has produced efficient results. The protocol has been evaluated for its efficiency using number of scenarios with varying number of nodes and varying number of jammer nodes.



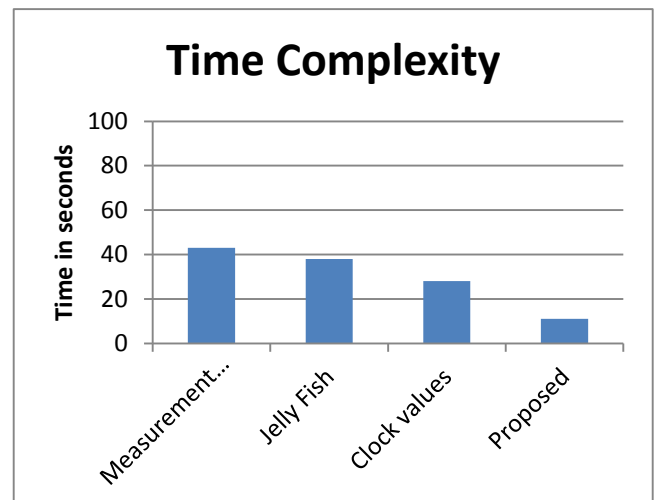
Graph 1: Comparison of detection accuracy

The Graph 1, shows the comparison of detection accuracy produced by different methods and it shows clearly that the proposed method has produced more accuracy than others.



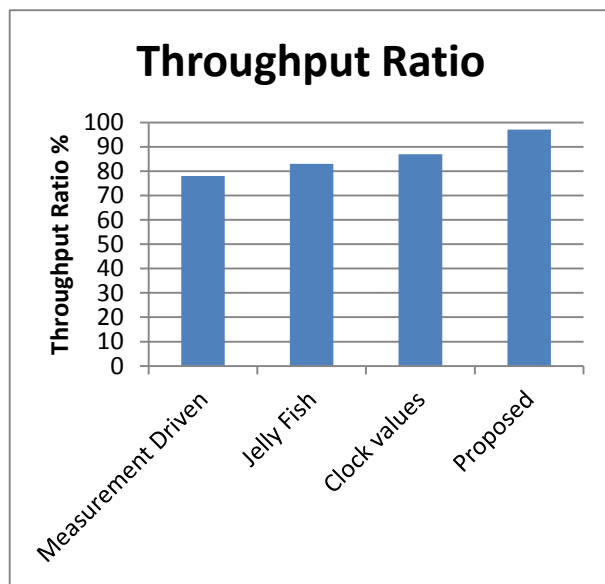
Graph 2: Comparison of false alarm ratio

The Graph 2, shows the comparative results of false alarm ratio produced by different methods and it shows clearly that the proposed method has produced efficient less false alarm than other methods.



Graph 3: Comparison of time complexity

The Graph 3, shows the time complexity produced by different methods in performing mitigation of DDOS jamming attack and it shows clearly that the proposed method has produced less time complexity than other methods.



Graph 4: Comparison of throughput performance

The Graph 4, shows the comparative result of throughput ratio produced by different methods and it shows clearly that the proposed method has produced more throughput than other methods.

#### Conclusion:

We proposed an time variant spatial traffic pattern based denial of service attack detection in a distributed jamming network. The method has been adapted to function in a distributed manner and will get to active mode in each time window and at the packet reception and transmission. The method generates traffic pattern at each time window and computes the trust weight for each of the neighbor who present in the current time window. Also the nodes verifies the location of the nodes at each time window or at the packet handling time to compute the trust weight. If the node has trust weight more than specific threshold then it will be forwarded to the node being identified otherwise it will be considered as jammer and ignored. The method has produced efficient results and increase the

throughput of the network by producing more accurate detection of DDOS attacks.

#### References:

1. V. Navda, A. Bohra, S. Ganguly, and D. Rubenstein, "Using channel hopping to increase 802.11 resilience to jamming attacks," in IEEE INFOCOM, Mini-Conf., 2007.
2. K. Pelechrinis, C. Koufogiannakis and S.V. Krishnamurthy, "Gamming the jammer: Is frequency hopping effective?," in WiOpt, 2009.
3. A. Sampath, H. Dai, H. Zheng, and B. Y. Zhao, "Multi-channel jamming attaches using cognitive radios," in IEEE ICCCN, 2007.
4. K. Pelechrinis, I. Broustis, S.V. Krishnamurthy and C. Gkantsidis, "ARES: An anti-jamming reinforcement system for 802.11 networks," in ACM CoNEXT, 2009.
5. W. Xu, W. Trappe, and Y. Zhang, "Anti-jamming timing channels for wireless networks," in ACM WiSec 2008.
6. K. Pelechrinis, I. Koutsopoulos, I. Broustis and S.V. Krishnamurthy, "Lightweight jammer localization in wireless networks: System design and implementation," Globecom, 2009.
7. I. Martinovic, P. Pichota, and J. B. Schmitt, "Jamming for good: A fresh approach to authentic communication in WSNs," ACM WiSec, 2009.
8. A. Mpitziopoulos, D. Gavalas, C. Konstantopoulos, and G. Pantziou, "A survey on jamming attacks and countermeasures in WSNs," IEEE Commun. Surveys Tuts., vol. 11, no. 4, 2009.
9. Xiangqian Chen, Kia Makki, Kang Yen, and N. Pissinou, "Sensor

- network security: A survey,” IEEE Commun. Surveys Tuts., vol. 11 no. 2, 2009.
10. Pelechrinis, K. A Measurement-Driven Anti-Jamming System for 802.11 Networks, IEEE Transacction on Networking, Vol. 19, Issue:4, pp:1208-1222, 2011.
  11. Deepesh Namdev<sup>1</sup>, Monika Mehra<sup>2</sup>, Detection Approach for Denial of Service Attack in Dynamic Wireless Networks,, Journal of Electronics and Communication Engineering Research Volume 2 ~ Issue 6(2014) pp: 01-06.
  12. I.Diana Jeba Jingle, Elijah Blessing Rajsingh, P.Mano Paul, Distributed Detection of DoS Using Clock Values in Wireless Broadband Networks, International Journal of Engineering and Advanced Technology (IJEAT) ISSN: 2249 – 8958, Volume-1, Issue-5, June 2012.
  13. Neha Singh, Sumit Chaudhary, Kapil Kumar Verma and A K Vatsa. Article: Explicit Query based Detection and Prevention Techniques for DDOS in MANET. *International Journal of Computer Applications* 53(2):19-24, September 2012.
  14. Deepak Vishwakarma and D s Rao. Article: Detection Mechanism for Distributed Denial of Service (DDoS) Attack in Mobile Ad-hoc Networks. *International Journal of Computer Applications* 102(9):23-26, September 2014.
  15. R.Ragupathy<sup>1</sup> and Rajendra Sharma, Detecting Denial of Service Attacks by Analysing Network Traffic in Wireless Networks, International Journal of Grid Distribution Computing Vol.7, no.3 (2014), pp.103-112.
  16. S. Raja Ratna, Mitigating Denial of Service Attacks in Wireless Networks, International Journal of Advanced Research in Computer Engineering & Technology (IJARCET) Volume 2, No 5, May 2013.

# Using Data Mining Techniques Analysing About Road Accidents

J.Sofia<sup>1</sup>

<sup>1</sup>Assistant Professor, Dept of CSE,  
Malla Reddy College Of Engineering,  
Maisammaguda,Dhulapally,Secunderabad

P.Sandeep<sup>2</sup>

<sup>2</sup>Assistant Professor, Dept of CSE,  
Malla Reddy College Of Engineering,  
Maisammaguda,Dhulapally,Secunderabad

**Abstract:** Globalization has influenced numerous nations. There has been an uncommon increment in the financial exercises and utilization level, prompting extension of travel and transportation. The expansion in the vehicles, activity prompt street mishaps. Considering the significance of the street security, government is endeavoring to distinguish the reasons for street mischances to decrease the mishaps level. The exponential increment in the mischances information is making it hard to investigate the limitations causing the street mishaps. The paper depicts how to mine incessant examples causing street mischances from gathered informational index. We discover relationship among street mischances and anticipate the kind of mishaps for existing as wel with respect to new streets. We make utilization of affiliation and grouping standards to find the examples between street mishaps and also foresee street mischances for new streets.

**Key Words:**Data mining, Association rule, Classification rule, Apriori algorithm, Naïve Bayes algorithm.

## 1.INTRODUCTION

India has second largest road network in the world. Road accidents happen quite frequently and they claim too many lives every year. It is necessary to find the root cause for road accidents in order to avoid them. Suitable data mining approach has to be applied on collected datasets representing occurred road accidents to identify possible hidden relationships and connections between various factors affecting road accidents with fatal consequences. The results obtained from data mining approach can help understand the most significant factors or often repeating patterns. The generated pattern identifies the most dangerous roads in terms of road accidents and necessary measures can be taken to avoid accidents in those roads.

## 2. METHODOLOGY

Descriptive or predictive mining applied on previous road accidents data in combination with other important information as weather, speed limit or road conditions creates an interesting alternative with potentially useful and helpful outcome for all involved stakeholders.

Association rule mining is used to analyse the previous data and obtain the patterns between road accidents. The two criterion used for association rule mining are support and confidence. Apriori algorithm is one of the techniques to implement association rule mining. In the proposed system, we use apriori algorithm

to predict the patterns of road accidents by analyzing previous road accidents data.

The steps for the apriori algorithm:

- Scan the data set and find the support(s) of each item.
- Generate L1 (Frequent one item set).Use Lk-1, join Lk-1 to generate the set of candidate k - item set.
- Scan the candidate k item set and generate the support of each candidate k – item set.
- Add to frequent item set, until C=Null Set.
- For each item in the frequent item set generate all non empty subsets.
- For each non empty subset determine the confidence. If confidence is greater than or equal to this specified confidence .Then add to Strong Association Rule.

INPUT DATASET ( A,B,C,D and E are accident types):

| TID | ITEMS   |
|-----|---------|
| 1   | A,C,D   |
| 2   | A,C,E   |
| 3   | A,B,C,E |
| 4   | B,E     |

**Minimum Support = 50%**

**Minimum Confidence = 80%**

Item set: A, B, C, D, and E

**STRONG ASSOCIATION RULE:**

This is the result obtained.

1. {B}->{E}
2. {CE}-> {A}
3. {AE}->{C}
4. {A}-> {C}
5. {C}->{A}

Classification rule is used to predict road accidents for new roads. In the proposed system, Naïve Bayes algorithm is used to implement classification rule.

Naïve Bayes algorithm steps:

**Step 1:** Scan the dataset (storage servers)

**Step 2:** Calculate the probability of each attribute value.  
[n, n\_c, m, p]

**Step 3:** Apply the formulae

$$P(\text{attributevalue}(a_i)/\text{subjectvalue}(v_j)) = (n_c + mp)/(n+m)$$

Where,

n = the number of training examples for which v = v<sub>j</sub>

n<sub>c</sub> = number of examples for which v = v<sub>j</sub> and a = a<sub>i</sub>

p = a priori estimate for P(a<sub>i</sub>|v<sub>j</sub>)

m = the equivalent sample size

**Step 4:** Multiply the probabilities by p

**Step 5:** Compare the values and classify the attribute values to one of the predefined set of class.

**Sample Example:**

Attributes(Constraints) – SpeedLimit, Wheather, PedestrianDistance [m=3]

Subject (Accident Type) – A1, A2 [p=1/2=0.5]

**Training Dataset**

| Road  | SpeedLimit(X,Y,Z) | Wheather(A,B,C) | Pedestrian Distance(P,Q,R) | Accident Type |
|-------|-------------------|-----------------|----------------------------|---------------|
| Road1 | X                 | A               | P                          | A1            |
| Road2 | X                 | B               | Q                          | A1            |
| Road3 | Y                 | B               | P                          | A2            |
| Road4 | Z                 | A               | R                          | A1            |
| Road5 | Z                 | C               | R                          | A2            |

**New Road6 Features – SpeedLimit - X,Wheather**

**A,PedestrianDistance - R Which Accident Type**

**A1/A2?**

$$P = [n_c + (m*p)]/(n+m)$$

| A1                                                                                                        | A2                                                                                                        |
|-----------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------|
| X<br>$P = [n_c + (m*p)]/(n+m)$<br>n=2,<br>n <sub>c</sub> =2,m=3,p=0.5<br>$p = [2+(3*0.5)]/(2+3)$<br>p=0.7 | X<br>$P = [n_c + (m*p)]/(n+m)$<br>n=2,<br>n <sub>c</sub> =0,m=3,p=0.5<br>$p = [0+(3*0.5)]/(2+3)$<br>p=0.3 |
| A<br>$P = [n_c + (m*p)]/(n+m)$<br>n=2,<br>n <sub>c</sub> =2,m=3,p=0.5<br>$p = [2+(3*0.5)]/(2+3)$<br>p=0.7 | A<br>$P = [n_c + (m*p)]/(n+m)$<br>n=2,<br>n <sub>c</sub> =2,m=3,p=0.5<br>$p = [2+(3*0.5)]/(2+3)$<br>p=0.3 |
| R<br>$P = [n_c + (m*p)]/(n+m)$<br>n=2,<br>n <sub>c</sub> =1,m=3,p=0.5<br>$p = [1+(3*0.5)]/(2+3)$<br>p=0.5 | R<br>$P = [n_c + (m*p)]/(n+m)$<br>n=2,<br>n <sub>c</sub> =1,m=3,p=0.5<br>$p = [1+(3*0.5)]/(2+3)$<br>p=0.5 |

$$A1 = 0.7 * 0.7 * 0.5 * 0.5 (p) = 0.1225$$

$$A2 = 0.3 * 0.3 * 0.5 * 0.5 (p) = 0.0225$$

Since **A1 > A2**

So this new road6 is classified to **A1**

### 3. CONCLUSIONS

Current system is manual where government sector make use of ledger data and analyze the data manually, based on the analysis they will take the precautionary measures to reduce the number of accidents. Proposed system uses road accidents data to mine frequent patterns and important factors causing different types of accidents. It discovers the associations among road accidents using apriori algorithm. It also predicts the common accidents

that may cause for new roads with the help of Naïve Bayes algorithm.

## REFERENCES

- [1] R. Agrawal, T. Imieliński, A. Swami, "Mining Association Rules Between Sets of Items in Large Databases", Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data, ACM, New York, NY, USA, pp. 207–216, 1993.
- [2] R. Agrawal, R. Srikant, "Fast Algorithms for Mining Association Rules in Large Data-bases", Proceedings of the 20th International Conference on Very Large Data Bases, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp 487-499, 1994.
- [3] A Araar et al., "Mining road traffic accident data to improve safety in Dubai", Journal of Theoretical and Applied Information Technology, 47(3), pp. 911-927, 2013.
- [4] L. Breiman, "Random Forests", Machine Learning, Vol. 45, pp. 532, 2001.
- [5] Beshah, T. (2005). Application of data mining technology to support RTA severity analysis at Addis Ababa traffic office. Addis Ababa, Addis Ababa University.
- [6] Getnet, M. (2009). Applying data mining with decision tree and rule induction techniques to identify determinant factors of drivers and vehicles in support of reducing and controlling road traffic accidents: the case of Addis Ababa city. Addis Ababa Addis Ababa University.
- [7] Chang, L. and W. Chen (2005). "Data mining of treebased models to analyze freeway accident frequency." Journal of Safety Research 36: 365- 375
- [8] Data Mining: Bagging and Boosting available at: <http://www.icaen.uiowa.edu/~comp/Public/Bagging.pdf>
- [9] Evanco, W. M., The Potential Impact of Rural Mayday Systems on Vehicular Crash Fatalities. Accident Analysis and Prevention, Vol. 31, 1999, pp. 455-462.
- [10] E. Frank and I. H. Witten. Generating accurate rule sets without global optimization. In Proc. of the Int'l Conf. on Machine Learning, pages 144–151. Morgan Kaufmann Publishers Inc., 1998.
- [11] Gartner Group High Performance Computing Research Note 1/31/95
- [12].Gartner Group Advanced Technologies & Applications Research Note 2/1/95
- [13] Data Mining and Data Warehousing available at: <http://databases.about.com/od/datamining/g/Classification.html>.
- [14] Genetic algorithm available at: [http://en.wikipedia.org/wiki/Genetic\\_algorithm](http://en.wikipedia.org/wiki/Genetic_algorithm).
- [15] Road Traffic Accident Statistics available at: [http://www.td.gov.hk/en/road\\_safety/road\\_traffic\\_accident\\_statistics/2008/index.html](http://www.td.gov.hk/en/road_safety/road_traffic_accident_statistics/2008/index.html).
- [16] Statistical Analysis Software, Data Mining, Predictive Analytics available at: <http://www.statsoft.com/txtbook/stdatmin.html>



# Dynamic Access Control Policies in Multi Cloud Storage Based NCC Clouds.

K.Divya Bharathi, Asst.prof, Department of Computer Science, MRCE

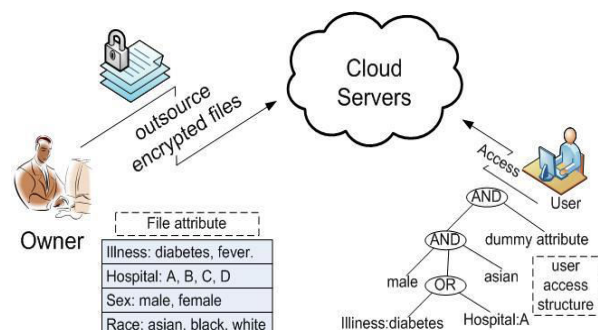
N.Keerthi, Asst.prof, Department of Computer Science, MRCE

**Abstract:** To provide mistake tolerance for reasoning space for storage, recent reports propose to stripe information across several reasoning vendors. However, if a reasoning suffers from a lasting failing and loses all its information, we need to fix the lost information with the help of the other surviving clouds to preserve information redundancy. We present a proxy-based space for storage system for fault-tolerant multiple-cloud space for storage called NCCloud, which accomplishes cost-effective fix for a lasting single-cloud failure. The protected transmitting of details among working together customers should be efficient as well as versatile in order to support accessibility management designs with different granularity levels for different kinds of programs such as protected team interaction, secure powerful conference meetings, and selective/hierarchical accessibility management published details. Accessibility management of short end users in cloud computing using Attribute-Set-Based Security (ASBE) with an requested structure of clients is not preferable for multi user access control in cloud computing. In this paper, we recommend the first provably protected Broadcast Group Key Management (BGKM) plan where each user in a team stocks a key with the reliable key server and the following re-keying for be a part of or leaving of customers needs only one transmitted concept. Our plan meets all the specifications set down for an effective GKM plan and needs no change to key stocks current customers have. We evaluate the security of our BGKM plan and evaluate it with the current BGKM techniques which are mostly ad-hoc.

**Index Terms:** Cloud computing, Attribute Based Encryption, Access Control, Security Model, Group Key Management, Trusted Authority for Key Sharing.

## I. INTRODUCTION

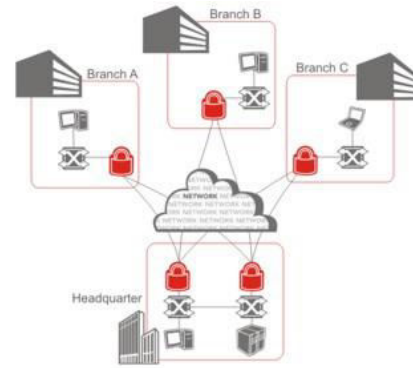
The fast advancement of the Internet and the Web in past decades has fundamentally changed the way individuals live, work, learn, think, shop, and impart everywhere throughout the globe. The open nature of the Internet makes it a twofold edged sword: On the one hand, telecom what's more, trade of data have never been speedier, less demanding, and more successful; on the other hand, new types of dangers like worms, infections, digital law violations have risen that bargain information/data security and client protection, and have postured numerous open difficulties to the world [1]. All sorts of client requests are actualized with great execution and association cost contains high. Clients may require any sort of assets to give the arrangements like pay per use way. Thinking handling gives the arrangements like unlimited wellsprings of subtle elements. We are going to take a shot at computation of time prerequisites, sources and asset necessities. Quality Based Encryption (EABE) permits just associations having a predefined arrangement of elements that can unscramble figure writings. EABE is suitable to openness administration, for example, the PC document talking about methods, in light of the fact that few associations can be accommodated the unscrambling of a figure content. We are recommending an improved EABE arrangement that is more viable than the previous one.



**Figure 1: Access control of data sharing in cloud.**

Through present sensitive computations we are going to devour the arrangements use with new security challenges in executing the system. In the storage room administration program, the thinking can let the client, data proprietor to shop his data, and talk about this data with different clients by means of the thinking, subsequent to the thinking can give the pay as you go air where individuals simply need to pay the cash for the storage room they utilize. For protecting the protection of the spared data, the data must be secured before presenting on the thinking. The security arrangement utilized here is quality based [4].

The EABE arrangement utilized a client's distinguishing proof as elements, and an arrangement of elements were utilized to secure and decode data. One of the primary weaknesses of the most current EABE method is that decoding is excessive for asset constrained contraptions because of coupling capacities, and the quantity of coupling capacities needed to unscramble a figure content creates with the many-sided quality in the availability arrangement [1][2][3]. The EABE arrangement can result the issue that data proprietor needs to utilize each sanction client's group key to secure data. Improved Attribute-Based Encryption (EEABE) which will be material for building adaptable, adaptable and fine grained access control of outsourcing information in distributed computing. EEABE grows the figure content approach quality set-based security (CP-ASBE, or ASBE for short) plot by (Bobbia et al., 2009) with requested structure of system clients, to perform adaptable, adaptable and fine-grained openness administration. All in all, the quality of information encryption with a symmetric-key calculation relies on upon the quality of the mystery key, which must be known by all taking an interest gatherings in correspondence. The procedure of selecting, circulating, putting away and upgrading mystery symmetric keys is called key administration. Solid, proficient and secure key administration is generally a testing issue in some genuine applications.



**Figure 2: Advanced key distribution in cloud server environment.**

Group key Management (GKM), as a particular instance of key administration, is identified with the taking after situation: Consider a server that sends information to a gathering of clients in a multicast/broadcast session through an open correspondence channel (As shown in figure 2). To guarantee information privacy, the server offers a mystery gathering key  $K$  with all gathering individuals and encodes the show information utilizing a symmetric encryption calculation with  $K$  as the encryption key [2]. Knowing the symmetric key  $K$ , any substantial gathering part can decode the scrambled telecast message. At the point when the gathering flow changes, i.e., when another client joins or a current client leaves the gathering, another gathering key must be produced and redistributed in a safe manner to all present gathering individuals, so that another gathering part can't recoup prior transmitted information (in reverse mystery), and a client who has left the gathering can't take in anything from future interchanges in the gathering (forward mystery). This procedure is called upgrade or re-keying. The procedure to keep up, circulate and upgrade the gathering keys is called gathering key administration.

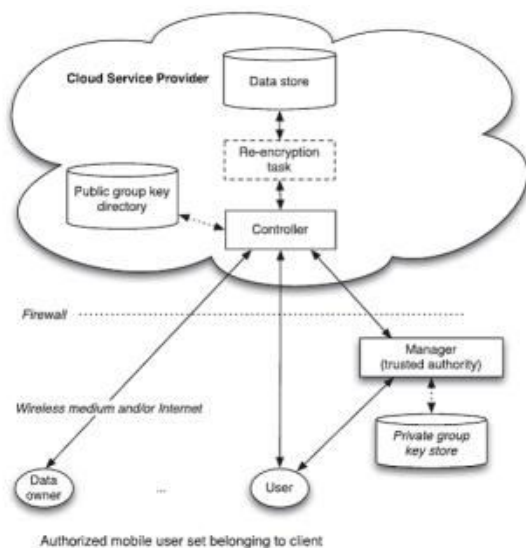
In this document, we recommend a new BGKM plan which, to the best of our information, is the first provably protected BGKM plan. Our new plan is versatile, effective and protected. It keeps the use of protected personal interaction programs little by not demanding any private communications when rekeying occurs either among the team associates or between the key server and a persisting team participant. The dimension the transmitted rekeying

information is linear with the count of team associates. In order to acquire a distributed team key, a team member need only execute effective hashing functions and an inner item of vectors over a limited area.

The rest of the paper organized as follows: Section II describes background/existing for access control services on cloud computing. Section III discusses related work in BGKM. Section IV formally defines BGKM with respect to security and efficiency. Section V presents experimental results with BGKM and EEABE in distributed key formation. Section VI concludes the paper.

## II. BACKGROUND APPROACH

Designing of multi cloud applications in NCC cloud may perform effective data assurance in real time cloud operations. Protecting regenerating rule qualities. We protect the fault patience and fix traffic preserving of FMSR requirements, with up to a small continuous expense. Thin-cloud storage space. Each server (or cloud-storage provider) only needs to give a primary interface for clients to create and read their saved data files.



**Figure 3: Network storage based data transmission in NCC Cloud.**

No computation capabilities are needed from the web servers to support our DIP plan.

Particularly, most cloud-storage providers nowadays give a RESTful interface, which contains the commands PUT and GET. PUT allows contacting data as a whole (no limited updates), and GET allows studying from a selected variety of bytes of information via a variety GET demand. Our DIP plan uses only the PUT and GET instructions to interact with each server. Our thin-cloud establishing allows our DIP plan to be portable to common types of storage space gadgets or solutions, since no execution changes are needed on the storage after sales. It is different from other “thick-”cloud-storage services where web servers have computational abilities and are capable of aggregating the evidence of several checks (e.g., [3], [4]). Nevertheless, we deal with how our approach can be prolonged to thick-cloud-storage of the additional data file, available online.

**Flexibility:** There should not be any boundaries on the number of possible difficulties that the consumer can make, since files can be kept for long-term archival. Also, the process size should be flexible with different parameter options, and this is useful when we want to lower the recognition rate when the saved data develop less important over time. Such flexibility should come without any additional charges.

## III. BGKM WITH SECURITY

In this area, we officially determine a transmitted team key control plan and its protection, and recommend a new team key control plan which allows any legitimate participant in the group which keeps an personal registration symbol (IST) to obtain a typical team key.

**Definition 1 (BGKM):** A transmitted team key control plan (BGKM) is consisting of two entities: 1) a key server (DIP), and 2) group members (Proposed Schemas), a chronic transmitted channel from DIP to all Proposed Schemas, an ephemeral personal channel<sup>3</sup> between DIP and each personal Proposed Schema, and the following phases:

**ParamGen** DIP requires as feedback a protection parameter  $k$  and results a set of community parameters Param, such as the sector KS of possible key principles.

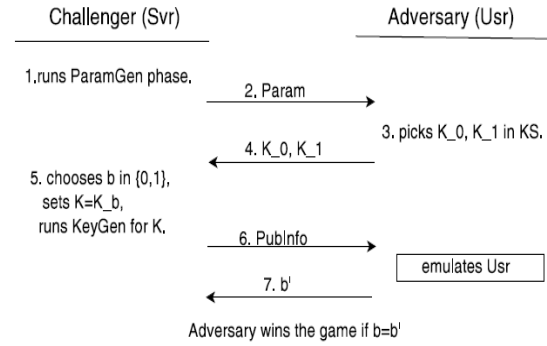
**TkDeliv** DIP delivers each Proposed Schema an personal registration symbol (IST) through a personal route.

**KeyGen** DIP selects a distributed team key  $K \in KS$ . In accordance with the ISTs of Proposed Schemas, DIP computes a set of principles PubInfo. DIP keeps  $K$  key, and shows through the transmitted channel PubInfo to all team associates Proposed Schema.

**KeyDer** Proposed Schema uses its IST and PubInfo to estimate the distributed team key  $K$ . Update When the distributed team  $K$  can no more be used (e.g., when there is a modify of group characteristics such as be a part of and leaving of team users), DIP produces new team key  $K'$  and PubInfo", then shows the new PubInfo to the team. Each Proposed Schema uses its IST and the new PubInfo" to estimate the new distributed team key  $K'$ . We contact the program after the Update phase a new "session". The Upgrade stage is also known as a rekeying stage.

**3.1. BGKM with Security:** A BGKM plan should allow a real team participant to obtain the distributed team key, and prevent anyone outside the team from doing so. Officially discussing, a BGKM plan should fulfill the following protection qualities. It must be appropriate, audio, key concealing, and forward/backward key defending.

1) **Correct:** Let Proposed Schema be a present team participant with an IST [14]. Let  $K$  and PubInfo be DIP's outcome of the KeyGen stage. Let  $K'$  be Proposed Schema's outcome of the KeyDer stage. A BGKM plan is appropriate if Proposed Schema can obtain the appropriate team key  $K$  with frustrating possibility, i.e.  $P_r[K = K'] \geq 1 - f(k)$  where  $f$  is negligible function for  $k$ .



**Figure 4: The attacker activity for BKGM's key concealing residence. With the information of PubInfo, the attacker is not able to distinguish one of its selected important factors from the other.**

**2) Sound:** Let Proposed Schema be a person without a legitimate IST. A BGKM is sound if the likelihood that Proposed Schema can get the right gathering key  $K$  by substituting the IST with a worth value that is definitely not one of the legitimate ISTs and afterward taking after the key induction stage KeyDer is immaterial.

**3) Key concealing:** A BGKM is key concealing if given PubInfo, any gathering which does not have a substantial IST can't recognize the genuine gathering key from an arbitrarily picked esteem in the key-space  $KS$  with non-negligible likelihood.

**4) Forward/in reverse key ensuring:** Suppose DIP runs an Update stage to produce Param for another shared gathering key  $K'$ , and a past part Proposed Schema is no more a gathering part after the Update stage. Let  $K$  be a past shared gathering key which can be inferred by Proposed Schema with token IST. A BGKM is forward key securing if a foe with learning of IST,  $K$ , and the new PubInfo can't recognize the new key  $K'$  from an arbitrary esteem in the key-space  $KS$  with non-negligible likelihood. Essentially, a BGKM plan is in reverse key securing if another bunch part Proposed Schema after the Update stage can't learn anything about the past gathering keys.

#### IV. EXPERIMENTAL EVALUATION

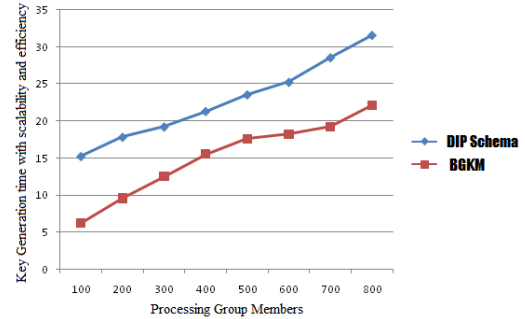
In this section we analyze the computational performance of ACV-BGKM. We imitate the KeyGen stage at DIP and the KeyDer stage at

Users. In the research, we differ both the dimension the actual primary area  $F_q$  and the dimension the number of Users, and evaluate the DIP-side and Proposed Schema-side calculations time [1]. To highlight on the mathematics functions, we do not depend plenty of here we are at hashing functions in the research. The rule is published in the Magma scripting language, and uses Magma's inner collection for limited area mathematics and fixing linear systems. Table 1 shows key generation with time for both attribute based encryption and broadcast group key management schemas in literal process.

| S.no | DIP     | PROPOSED SCHEMA in BGKM |
|------|---------|-------------------------|
| 1    | 15.2562 | 6.1245                  |
| 2    | 17.5891 | 9.2456                  |
| 3    | 21.5632 | 12.3256                 |
| 4    | 25.6785 | 15.5478                 |
| 5    | 32.4569 | 16.4563                 |
| 6    | 35.4587 | 23.2478                 |
| 7    | 39.5632 | 26.3547                 |

**Table 1: Generation of key values with respect to time**

As shown in above we construct efficient graphical representation of our sourcing data. The research was conducted on a machine running GNU/Linux kernel edition 2.6.9 with a Double Primary AMD Opteron(TM) Processor 2200 MHz and 16 G bytes storage. Only one processor was used for calculations. The following diagrams tell performance of broadcast group key management with access control opportunities.



**Figure 5: Comparison of DIP in ABE and PROPOSED SCHEMA in BGKM with respect to time in terms of key generation.**

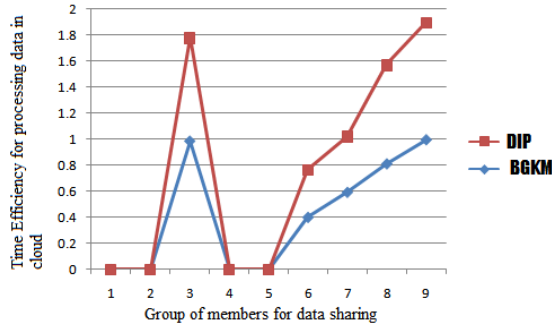
Fig 5 reviews the ACV-BGKM operating time at DIP and Proposed Schema for team dimensions 600, 800, 1000 and 1200, and with the dimension the primary area which range from 64 to 128 pieces. The operating time is averaged over 20 versions. As proven in the determine, the common calculations time improves in common as the dimension the primary area improves. The real operating time relies on the prime field that is selected and the way area mathematics is conducted in Magma. The following table 2 shows time efficiency for giving permissions to all the registered user in terms of access permissions regarding storage data in cloud.

| S.No | DIP    | PROPOSED SCHEMA BGKM |
|------|--------|----------------------|
| 1    | 0.9874 | 0.78965              |
| 2    | 0.4012 | 0.36581              |
| 3    | 0.5934 | 0.4263               |
| 4    | 0.8124 | 0.7569               |
| 5    | 0.9975 | 0.8965               |

**Table 2: User control access with storage of cloud data.**

An also we perform efficient performance evaluation in user granted permissions for accessing file from one to other users present in cloud. The performance evaluation of the user

access control in data storage in shown in figure 6.



**Figure 6: Computational time efficiency in application process of DIP and Proposed Schema in access control in usability.**

Fig. 6 reviews the ACV-BGKM operating time at DIP and Proposed Schema for set area measures (in bits) 64, 80, 96 and 112, with the dimension the team which range from 100 to 2000 associates. The running time is averaged over 20 versions. It reveals that the ACV-BGKM rekeying procedure operates fast on DIP when there are thousands of Proposed Schemas in the team. It requires less than two moments for DIP to generate new PubInfo when there are up to 2000 Proposed Schemas and when the primary area is huge enough. Both numbers display that it requires very little here we are at a Proposed Schema to obtain the distributed team key, and a essentially brief time frame for the DIP to produce the key and the transmitted rekeying details, when the real limited area and the team dimension are both significantly huge. Further performance gains can be carried out when the primary variety  $q$  is selected to be in a unique type, e.g., a common Mersenne primary (Solinas prime), for which quick area mathematics in  $F_q$  is available.

We lay the base towards a quicker method of ACV-BGKM, known as FACVBGKM, at the price of extra area, pre-computation and an catalog. It follows a baby step- giant-step (BSGS) rekey procedure where irregular massive actions are conducted analogous to the ACV-BGKM plan [2] . However, the amortized

computational and interaction price is reduced by the release of regular small actions. Due to area restrictions, we only describe the changes to the ACV-BGKM method below.

**Protocol (FACV-BGKM):** FACV-BGKM performs under identical circumstances as ACV-BGKM. ParamGen DIP chooses  $N'' = N + M$  where  $M \geq N$ . For the highest possible protection and minimum amortized price, it is suggested to set  $M = N$ .

**TkDeliv** DIP assigns an index  $i(1 \leq i \leq N)$  chosen consistently at unique, to each of the  $n$  current customers. DIP selects  $N$  ISTs and delivers an IST and corresponding catalog to each customer. The staying  $N - n$  precomputed ISTs are used for rekeying when new customers be a part of the team.

**KeyGen** DIP makes an  $N \times (N \times M) F_q - matrix$   $A$  ware for a given  $i(1 \leq i \leq N)$  .

$$a_{i,j} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } 1 \leq j \leq N \text{ and } i \neq j \\ H(ist_i \| Z_j) & \text{if } N < j \leq N + M \end{cases}$$

Like in ACV-BGKM method, DIP determines the zero area of  $A$  with a set of its  $M$  basis vectors, and chooses an accessibility management vector  $Y$  as one of the primary vectors. DIP caches these basic vectors and represents  $Y$  as “used.” DIP constructs an  $(N + M)$ -dimensional  $F_q$ -vector

$$X = \left( \sum_{i=1}^n K_i \cdot e_i^T \right) + Y \quad \text{where } e_i \text{ is the } i\text{th}$$

standard basis vector of  $F_q^{N+M}$  . Observe that, in contrast to ACVBGKM, the key is included to all the places corresponding to legitimate spiders. Like, ACVBGKM, DIP places  $PubInfo = \&X, (z_1, z_2, \dots, z_M)'$ , and shows PubInfo via the broadcast channel.

**KeyDer** Proposed Schemai, understanding the catalog  $i$  and  $isti$ , originates the  $(N + M)$ -

dimensional row Fq-vector  $v_i$  which matches to a row in  $A$ . Proposed Schemai originates the team key as  $K = v_i \cdot X$ .

Update Unlike ACV-BGKM, DIP does not run the finish KeyGen stage again. If a new Proposed Schemas connects the team, DIP chooses an rarely used catalog  $t$  and is it from the pre-computed ISTs and computes the new  $\hat{X}$  with a new key  $\hat{K}$ . If an present User results in the team, DIP chooses a new key  $\hat{K}$  and determines a new

$$\hat{X} = (\sum_{j=1}^n K.e_i^T) + Y$$

Where  $\hat{Y}$  is an “unused” foundation vector which is among the pre-computed set in KeyGen stage. DIP marks  $\hat{Y}$  as “used”, and shows only  $\hat{X}$  while maintaining the other community details the same. We contact these functions a “baby-step rekey” since it only needs time  $O(N)$  in comparison to  $O(N^3)$  in ACV-BGKM. A finish KeyGen (i.e. “giant-step rekey”) eventually  $O(N^3)$  needs to be performed every  $M$  Up-dates since otherwise a team participant who has been legitimate for the last  $M$  classes can restore the zero area of  $A$ , thus the matrix  $A$  itself. A giant-step rekey also needs to be conducted with a resized matrix  $A$  before  $M$  updates, if the variety of connects exceeds  $N - n$  after the present giant-step rekey to provide new customers. As described above, the KeyGen price is amortized to acquire a plan quicker than the ACVBGKM scheme. Due to area restrictions, we bypass the security/performance analysis of the FACV-BGKM plan from this document. We observe that it is an exciting start analysis problem to choose the maximum  $M$  and  $N$  principles based on the program situation.

## V. CONCLUSION

We have suggested a new BGKM plan ACV-BGKM which is managed by a trusted key server, and allows any legitimate customer in the team to acquire a distributed team key on its own from transmitted community details. The plan reduces the use of personal peer-to-peer communication programs, and

only uses a transmitted route to provide new rekeying messages when the team key needs to be modified. The interaction expense is straight line with the number of customers in the team. The plan uses only effective hash functions and straight line geometry over finite areas in calculations, and does not require any security plan. It is protected in that even a computationally unbounded attacker cannot acquire the distributed team key without a valid symbol from the key server. The key derivation is effective for any team participant. The experimental outcomes show that the creation of the rekeying details requires a few months on a laptop or computer for a number of a large number of associates. As upcoming work, we plan to empirically evaluate the efficiency of the FACV-BGKM plan under different parameters.

## VI. REFERENCES

- [1] “Enabling Data Integrity Protection in Regenerating-Coding-Based Cloud Storage: Theory and Implementation” by Henry C.H. Chen and Patrick P.C. Lee proceedings in IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL. 25, NO. 2, FEBRUARY 2014 by author and published journal.
- [2] M. Armbrust, A. Fox, R. Griffith, A.D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, “A View of Cloud Computing,” *Comm. ACM*, vol. 53, no. 4, pp 50-58, 2010.
- [3] G. Ateniese, R. Burns, R. Curtmola, J. Herring, O. Khan, L. Kissner, Z. Peterson, and D. Song, “Remote Data Checking Using Provable Data Possession,” *ACM Trans. Information and System Security*, vol. 14, article 12, May 2011.
- [4] K. Bowers, A. Juels, and A. Oprea, “HAIL: A High-Availability and Integrity Layer for Cloud Storage,” *Proc. 16th ACM Conf. Computer and Comm. Security (CCS '09)*, 2009.
- [5] K. Bowers, A. Juels, and A. Oprea, “Proofs of Retrievability: Theory and Implementation,” *Proc. ACM Workshop Cloud Computing Security (CCSW '09)*, 2009.
- [6] B. Chen, R. Curtmola, G. Ateniese, and R. Burns, “Remote Data Checking for Network Coding-Based Distributed Storage Systems,” *Proc. ACM Workshop Cloud Computing Security (CCSW '10)*, 2010.
- [7] H.C.H. Chen and P.P.C. Lee, “Enabling Data Integrity Protection in Regenerating-Coding-Based Cloud Storage,” *Proc. IEEE 31<sup>st</sup> Symp. Reliable Distributed Systems (SRDS '12)*, 2012.
- [8] L. Chen, “NIST Special Publication 800-108,” Recommendation for Key Derivation Using Pseudorandom Functions (Revised), <http://>

csrc.nist.gov/publications/nistpubs/800-108/sp800-108.pdf, Oct. 2009.

[9] R. Curtmola, O. Khan, and R. Burns, "Robust Remote Data Checking," *Proc. ACM Fourth Int'l Workshop Storage Security and Survivability (StorageSS '08)*, 2008.

[10] R. Curtmola, O. Khan, R. Burns, and G. Ateniese, "MR-PDP: Multiple-Replica Provable Data Possession," *Proc. IEEE 28th Int'l Conf. Distributed Computing Systems (ICDCS '08)*, 2008.

[11] Wan, Z., Liu, J. E., & Deng, R. H. (2012). HASBE: a hierarchical attribute-based solution for flexible and scalable access control in cloud computing. *Information Forensics and Security, IEEE Transactions on*, 7(2), 743-754.

[12] Yu, S., Wang, C., Ren, K., & Lou, W. (2010, March). Achieving secure, scalable, and fine-grained data access control in cloud computing. In *INFOCOM, 2010 Proceedings IEEE* (pp. 1-9)IEEE.

[13] N. Shang, M. Nabeel, F. Paci, and E. Bertino, "A privacy-preserving approach to policy-based content dissemination," in *ICDE '10: Proceedings of the 2010 IEEE 26th International Conference on Data Engineering*, 2010.

[14] X. Zou, Y. Dai, and E. Bertino, "A practical and flexible key management mechanism for trusted collaborative computing," in *INFOCOM. IEEE*, 2008, pp. 538–546. [Online]. Available: <http://dblp.uni-trier.de/db/conf/infocom/Infocom2008.html#ZouDB08>.

[15] J. Li, N. Li, and W. H. Winsborough, "Automated trust negotiation using cryptographic credentials," in *Proc. ACM Conf. Computer and Communications Security (CCS)*, Alexandria, VA, 2005.

[16] V. Goyal, O. Pandey, A. Sahai, and B. Waters, "Attribute-based encryption for fine-grained access control of encrypted data," in *Proc. ACM Conf. Computer and Communications Security (ACM CCS)*, Alexandri.

[17] Zhiguo Wan, Jun'e Liu, and Robert H. Deng, "HASBE: A Hierarchical Attribute-Based Solution for Flexible and Scalable Access Control in Cloud Computing".

[18] S. Yu, C. Wang, K. Ren, and W. Lou, "Achieving secure, scalable, and fine-grained data access control

in cloud computing," in *Proc. IEEE INFOCOM 2010*, 2010, pp. 534–542.

[19] J. Bethencourt, A. Sahai, and B. Waters, "Ciphertext-policy attributebased encryption," in *Proc. IEEE Symp. Security and Privacy*, Oakland,CA, 2007.

[20] G.Wang, Q. Liu, and J.Wu, "Hierarchical attribute-based encryption for fine-grained access control in cloud storage services," in *Proc. ACM Conf. Computer and Communications Security (ACM CCS)*, Chicago, IL, 2010.

# A Cloud Based System to Sense Security Vulnerabilities of Web Application in Open-Source Cloud IAAS

K.Himabindu<sup>1</sup>

G.mamatha<sup>2</sup>

S.Kavitha<sup>3</sup>

<sup>1,2,3</sup>Assistant Professor, Department of Computer Science, Malla Reddy institute of technology&Science,Hyderabad.  
khimabindu.2009@gmail.com, mamatha0503@gmail.com, kavi.shepuri34@gmail.com

**Abstract:-**In present times the use of web applications has money, reading news and such comparable exercises. It is watched that in online applications security aspect is disregarded by numerous turn into a need in our every day routine, for example, internet shopping, saving developers which may prompt significant loss of data and put an inquiry on secrecy of the client data. Mostly because of the need of security mindfulness, aptitude and now and then steady weight on the designer to finish the errand inside given dead line causes serious web vulnerabilities. In this paper, another imaginative administration demonstrate is proposed which brings new opportunity on request benefit pool. This is accomplished with the assistance of virtualization, cloud computing, and sharing different assets over a private cloud condition. A cloud based static examination framework for unmasking security vulnerabilities in PHP is proposed in this paper. The proposed framework will help a user to identify factors affecting the vulnerabilities on a client system and call attention to the polluted code inside the source code. Testing, private cloud, security, hazard investigation, vulnerabilities.

**Key words:** Automated Automated testing, private cloud, security, risk analysis, vulnerabilities.

## I.INTRODUCTION

Robert L. Grossman said about sorts of cloud, clarification about distributed computing[7], on-request distributed computing occurrences

and those that give on-request registering limit, points of interest, and impediments, with layered administrations. Following are the primary four territories of the cloud-based administration sector

- To start with, the administrators of Telecom industry will accelerate the usage of distributed computing system to accomplish business advancement with the guide of cloud computing.
- Second, little and medium-sized business enterprises are compelled to set-up open cloud benefit center.
- Third, government supports the change of electronic government through the distributed computing stage.
- Fourth, training industry enhances asset use by initiating sharing cloud stage Ironically security is disregarded by numerous associations which some of the time prompt loss of major financial and specialized asset of the association. The explanation for an absence of security mindfulness in little scale ranches is carelessness and here and there absence of assets. In programming testing different testing models can be incorporated and these are associated with each period of programming advancement life cycle. In conventional approach of programming testing re-quires heaps of time, labor and equipment assets. Here and there programming testing causes blunder because of human information mistake. In entire period of programming testing of web application or programming application security testing can be considered with an abbreviate the improvement life cycle.

To upgrade the product testing proficiency and testing adaptability distributed computing considered as great option way to deal with take care of the issue. Distributed computing transfer on designating asset sharing, for example, space memory, I/O, gadgets, handling speeds and related things to accomplish huge scale availability. The foundation of distributed computing is in consistence to all finished use framework and assets. Cloud can be classified into open, private, group, half breed cloud. [18] [22] Web administrations are intended to beg of association exercises, check monetary possibility, on-request resourcing, assignment planning, with that dispatching calculation to use of PC assets and versatility to perform amplification the testing fields. Distributed computing centers around business improvement work out. TaaS is another model in cloud testing[18]. Advantages of moving existing electronic security testing framework to cloud-based the web security testing framework

- Use of existing foundation: Cloud base security testing framework in distributed computing will empower the client to pay-per-utilize improvement display. according to the need of the client, a client can broaden its equipment assets whenever. any make framework anything as administration are takes after. [19]

- Quicker sending and re-ease of use: Applications created by one substance can be made accessible on the cloud. No requirements to recode the things simply convey on cloud and utilize it in versatile mode. Code can be shareable with alternate partners. because of the elite of cloud speedier advancement can be accomplished.

- Manageability and maintainability:

Combined visibility and control will be provided by the cloud under a single directory of services helping the requirement of lengthy

procurement and maintenance of cloud-based security testing infrastructure.

- Scalability: Applications and infrastructure deployed on the common Cloud platform can take advantage of the qemu virtualized kind of the cloud to scale as required. and cases, and profile the consumed resources. Experiments were carried out to help users figure out the bottleneck of existed testing program and cloud systems.

- Cost reduction: The major utility of cloud system and method for performance testing on a Cloud Computing architecture is that aims to providing easy access to costly software running on High performance processors to users at organizational startup considerable facilities. Cloud computing security testing system design and develop as very cost efficient.

- Reduced effort in managing technology: On Open Stack Cloud availability of services are more. Easy provisioning of computing resources will ensure more consistent technology upgrades and expedite fulfillment of IT resource requests. [20]

- Powerful computing and storage capacity: Cloud-based security testing server enables to process the data in the private cloud environment. Cloud computing frame-work supports powerful computing and storage capacity.

- Virtualization: It is managed, expended, migrated, and backup through virtualization platform. It put the underlying hardware, including servers, storage and networking equipment, comprehensive virtualization

- Resource pooling: The pool-based model result is that physical computing resources become invisible to consumers; Example of

resources includes storage, processing, memory, network bandwidth, and virtual machines.

Software testing is task intensive and often implicates substantial collaboration among end user, developer, and tester. This paper proposes “A cloud-based security vulnerabilities (CBSV) testing system”. which works on private cloud to fulfill the on-demand resource pool of end user is proposed.

## II. RELATED WORK

In cloud testing, Testing-as-a-Service has enough power to provide testing skillful to end users effectively. Designer has to design a model such as the private cloud should adapt the user resource and work on it. [12] Software are frequently disclosing the different kinds of threats and attacks such as SQLi, CSS, Code evaluation, Code execution, session fixation vulnerabilities. Sajjad Rafique,[2] Proposed solutions mapped against Firstly: for the software development stages for which the solution has been proposed and secondly: for the web application vulnerabilities mapping according to OWASP Top 10 security vulnerabilities [18]. Chorng-Shiuh [3] proposed an automatic software testing service with the benefits of real-time software testing and automatically computation scaling. It can also describe software testing in parallel, and can automatically pair the source code to perform statistics analysis, generate test drivers and cases, and profile the consumed resources. Experiments were carried out to help users figure out the bottleneck of existed testing program and cloud systems. Yue Zhou,[11] explained a Cloud Application performance testing, traditional performance testing tool and method limited by many characteristics which are different from Cloud Computing and traditional application. This article describes the limitation and weak point of traditional tool

and method for performance testing on a Cloud Computing application. Face to challenge, some new methodologies are introduced. Anh Nguyen-Tuong[17] proposed a fully automated approach for securely hardening web applications. It is based on precisely tracking tainted of data and checking specifically for dangerous content in only in parts of commands and output that came from untrustworthy sources. Unlike previous work in which everything that is derived from the tainted input is tainted, our approach precisely tracks tainted within data values.

### A. Security Vulnerabilities Detection approach

University Lisboa [21] Web application vulnerabilities detection can be done with two ways, either user white box testing or black box testing. White box testing also considered as static code analysis. white box testing can be done using various tools like Pixy, Frosty etc. Considering Black box testing, the internal details of testing can be hidden from tester. User fuzzy logic techniques over the web http request. So static scanning helps out to simulate numerous scenarios such as hacker's intentional attacks, or end user identified attacks. Security vulnerability testing can avoid hectic of recursive tasks. [15]

## III. PROPOSED CLOUD BASED WEB SECURITY TESTING SYSTEM(CBSV) ARCHITECTURE

Proposed system overview is given in fig 1, Here. Software-as-service model end-user can interact with services like testing of web application, analyses results, generation of result reports, add a network, add instances, allocate Disk size to instances. Overall management can be performed by the user as per the need. Proposed system can be also used as platform as service model to test security of web application. Authorized user can use the system as per requirement to test web applications. On infrastructure which supports

the storage capabilities, QEMU/KVM Virtualization, so we can utilize hardware resources. The advantage of this is we can pool resources like RAM, HDD, OS images, etc.

#### IV. IMPLEMENTATION AND EXPERIMENTS

##### A. System Implementation :

Implementation will consist of stages which define system's overall architecture

##### 1) Stage 1: Implement Cloud Architecture

In this stage, Construction of private cloud environment will be done. To support scalability, performance, support testing environment etc. We are using VMware machine to create private cloud environment and to achieve these we configure the network, hardware/resource with a help of Ubuntu 14.04 server on VMware tool. The following steps for creating the private cloud. Using standard procedure, Implement cloud network in existence.

- The configuration of Sources OS selection.
- The configuration of Network for controller node.
- Verification resource sharing (Virtual machine management)
- Controlling event management of Virtual machines.
- OpenStack version Icehouse.

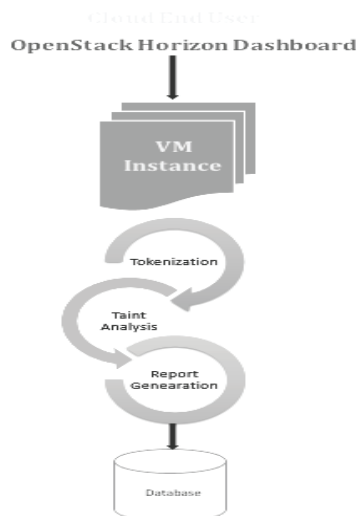


Fig. 1. Proposed System Architecture

2) Stage 2: Design Web service for Automated Testing Environment. In this phase, with the help of Security testing tools, the system will scan vulnerabilities of the given input application/ files.

- The tester should be able to specify the specific host name or IP address.

Setting testing capabilities.

- Client-side attack
- Command execution
- Information Disclosure
- cross side scripting
- session fixation
- Testing Customization.

3) Stage 3: Method for calculating Impact factor of Vulnerabilities

a) From given identified vulnerabilities, system will find out impact factor.

b) On the basis of impact factor risk analysis will be calculated.

c) System will predict vulnerabilities factor, for an end user.

d) To detect vulnerabilities author will use token tainted analysis in that treats every words, key- words as a token make them separate. The Tokenization: Analyze whole web application correctly; we need to split whole code into tokens to analyze the whole web application properly. This includes the additional character with function name or adding curly braces to program flow constructs where not a single brace is used. Tokenization concept is also helpful to reduce an overhead of comments, inline HTML. [4]

```
// securing LDAP injection
$_SECURING_LDAP = array(
);

// all specific securings
$_SECURES_ALL = array_merge(
    $_SECURING_XSS,
    $_SECURING_SQL,
    $_SECURING_PREG,
    $_SECURING_FILE,
    $_SECURING_SYSTEM,
    $_SECURING_XPATH
);

// securing functions that work only when embedded in quotes
$_QUOTE_ANALYSIS = $_SECURING_SQL;
```

Fig2: Snapshot of Token generation for PHP code4

4) Stage 4: Scanning Input and generating Vulnerability report. In this phase with the help of vulnerabilities, the end user will generate impact factor will be generated by each and every vulnerability. So it will helpful to end user to identify the effectiveness of detected vulnerability.

5) Stage 5: With the help of Impact Factor Generate Risk analysis report will be generated. Weiwei Wang[13] propose an algorithm framework for the empirical evaluation of classifier security at design phase that extends the model selection and shows evaluation steps of the classical design cycle. This classifier algorithm will be used in a proposed system for classification of vulnerabilities. In this phase, with the help of impact factor, calculate Risk with the vulnerability like in this phase, with the help of impact factor, calculate Risk with the vulnerability like

- The amount of memory it is accessing.
- Type of memory
- Type of Domain accessing in OS.

- The line of Code
- Generating report format in Human and machine readable format.– PDF

### *B. Challenges in Integrating Web Service System:*

The various challenges encountered while adopting in the web security testing system. The challenges in integrating CBSV web system pointed includes: Yue Zhou, [11] identified the anticipated challenges in adapting On Demand Testing System which includes more investment is needed in the area of infrastructure and human development, due to lack of a sufficient number of computers and room for installing them, appropriate software and adequate technology expertise are necessary for effective adaptation. [11] pointed out the key challenges in integrating security testing system which include: non availability of Cloud base security vulnerability (CBSV) infrastructure, resistance to adopting new technology, lack of motivational quality, frequent changes in administration which hamper the exploration and exploitation of CBSV opportunities. lack of IT infrastructure, inadequate electricity supply, inadequate IT tools. In the [10] authors pointed out the various problems associated with existing e-security testing system. These problems include human interference, security issues, inadequate training for the tester, developer and enduser and complexity of software.

## V. CONCLUSION

In this paper, it is observed that combination of web security testing stage with cloud design can accomplish on-request asset for web application testing. It can tackle a noteworthy issue of clients which they look amid a period of testing that is asset inaccessibility. The paper

additionally features web security testing with Effect factor with insights results and bring up the spoiled code inside the source code. What's more, this framework empowers the client to pick a different kind of OS or Make possess volume space for testing and furthermore produce the report for investigation.

## REFERENCES

- [1]. Dan Tao, Zhaowen Lin, and Cheng Lu, Cloud Platform Based Automated Security Testing System for Mobile Internet, IEEE Transaction Tsinghua Science And Technology ISSN11007-021411011311pp537-544Volume 20, Number 6,December 2015.
- [2]. Sajjad Rafique, Mamoon Humayun, Bushra Hamid, Ansar Abbas,Muhammad Akhtar, Kamil Iqbal,Web Application Security Vulnerabilities Detection Approaches: a Systematic Mapping Study,IEEE SNPD 2015, June 1-3 2015, Takamatsu, Japan.
- [3]. Chong-Shiuh Koong, Tzu-I Yang, and Chang-Chung Wu,An Implementation of Automatic Software Testing on Cloud Computing Environment, Computer Science and Its Applications, Lecture Notes in Electrical Engineering, Springer-Verlag Berlin Heidelberg 2015
- [4]. Josip Bozic, Bernhard Garn, Dimitris E. Simos, Franz Wotawa, Evaluation of the IPO-Family Algorithms for Test Case Generation in Web Security Testing, Eighth International Conference on Software Testing, Verification and Validation Workshops (ICSTW) IEEE 2015.
- [5]. Yuan-Hsin Tung, Chen-Chiu Lin, Hwai-Ling Shan,Test as a Service: A framework for Web security TaaS service in cloud environment, 8th International Symposium on Service Oriented System Engineering IEEE 2015.
- [6]. Anna Thankachan, R. Ramakrishnan, M.Kalaiarasi,A Survey and Vital Analysis of Various State of the Art Solutions for Web Application Security, ICICES2014 - S. A. Engineering College, Chennai, TamilNadu, India IEEE 2015.
- [7]. Filip Holik, Josef Horalek, Ondrej Marik, Sona Neradova, Stanislav Zitta, Effective penetration testing with Metasploit framework and methodologies, 15th IEEE International Symposium on Computational Intelligence and Informatics,CINTI 2014.
- [8]. Dale D. Reitze, Using Commercial Web Services to Build Automated Test Equipment Cloud Based Applications, IEEE 2014.
- [9]. Sheng-Jen Hsieh, Guo-Heng Luo, Shyan-Ming Yuan, Hsiao-Wei Chen,A flexible public cloud based testing service for heterogeneous testing targets, IEICE - Asia-Pacific Network Operation and Management Symposium (APNOMS) 2014.
- [10]. Yao-Wen Huang, Shih-Kun Huang, Chung-Hung Tsai,Web Application Security Assessment by Fault Injection and Behavior Monitoring, ACM 1-58113-680-3/03/0005. ACM-2003
- [11]. Yue Zhou, Wenchuang Qin, Nafei Zhu, Challenges To Performance Testing Of The Cloud Application Developing, International Conference on Software Engineering and Computer Science (ICSECS2013) 2013.
- [12]. Yi-Lun Pan, Chang Hsing Wu, His-En Yu, Hui-Shan, Chen and We- icheng Huang,Creating Your Own Private Cloud: EzillaToolkit For Coordinated Storage, Computing, and NetworkingServices, 12th IEEE/ACM International Symposium on Cluster,Cloud and Grid Computing, 2012.
- [13]. K. V. Arunkumar & E. Samlinson,Testing as a service (TaaS)an enhanced security framework for taas in cloud environment,International Journal of Internet Computing ISSN No: 2231 ,2012.
- [14]. Weiwei Wang, Xiaosong Zhang, Ting Chen, Xueyang Wu, andXiaoshan Li,Cloud Computing Based Software TestingFramework Design and Implementation, Advances in Computer,Communication, Control & Automation, Springer-Verlag BerlinHeidelberg 2011.
- [15]. Jerry Gao, Xiaoying Bai, and Wei-Tek Tsai,Cloud Testing-Issues,Challenges, Needs and Practice, An International Journal(SEIJ), Vol. 1, No.1, SEPTEMBER 2011.
- [16]. Anh Nguyen-Tuong, Salvatore Guarnieri, Doug Greene, David Evans,Automatically

Hardening Web Applications Using Precise Tainting, University of Virginia Computer Science Technical Report CS-2004- 36 December 2004.

[17]. Robert L. Grossman, The case of cloud Computing, IEEE Computer Society 2009.

[18]. Open Web Application Security Project(OWASP) 2015, OWASP Top 10, Available FTP: <https://www.owasp.org/index.php>

[19]. Jerry Gao<sup>1,2</sup>, Xiaoying Bai<sup>2</sup>, and Wei-Tek Tsai<sup>2,3</sup> Testing as a Service (TaaS) on Clouds, 2013, 2013 IEEE Seventh International Symposium on Service-Oriented System Engineering

[20]. Openstack Installation Guide Organization, Single node machine. [www. docs.openstack.org](http://www.docs.openstack.org)

[21]. Realistic Vulnerability Injections in PHP Web Applications, Nuno Ferreira Neves, 2011.

[22]. Verification of multi-owner shared data with collusion resistant user revocation in cloud, D. S. Kasunde, A. A. Manjrekar, 2016 International Conference on Computational Techniques in Information and Communication Technologies (ICCTICT), New Delhi, 2016, pp. 182-185. doi: 10.1109/ICCTICT.2016.7514575

# Sentiment Analysis in Healthcare Using Social Media

Rajasekhar Nennuri

Department of CSE

Institute of Aeronautical Engineering

Hyderabad, India

rajasekharnennuri@gmail.com

Dr I Surya Prabha

Department of CSE

Institute of Aeronautical Engineering

Hyderabad, India

ipsurya17@gmail.com

## Abstract:

The global healthcare industry is under significant pressure to reduce costs and more efficiently manage resources while improving patient care. Healthcare organizations have operated based on imprecise or incomplete cost and care measurements and did not have the comprehensive view of clinical and operational processes they needed to identify areas for improvement. The healthcare industry has recently begun to turn to data and analytics in ways that are similar to other industries that rely on digital information to improve service and reduce costs. In this paper we discuss about opinions on health care using twitter data. Ailments like headache, flu, fever will be discussed. Mixed opinions are also discussed and observed.

## Keywords—

Data mining, Natural language processing, Rapid Miner

## I. Introduction

Sentiment analysis technique is an effective means of discovering public opinions various companies often use online or paper based surveys to collect customer comments. Due to the emergence of social networking sites and applications, people tend to comment on their facebook or tweet profile. Therefore the paper based approach is not an efficient approach. Only a very small customer base can be reached and there is no guarantee that their answers in the survey are honest or not. Here social media

comes into play. Facebook, Twitter and all other social media sites are full of people's opinions about products/services they use, comments about popular personalities and much more. Hence mining opinions about various subject matters from social media is a much more innovative approach for market analysis. A lot of research has been done on opinion mining from social media, most of which focuses on people's sentiment towards various topics. But analyzing social media data in this manner gives a much generalized idea. To make it more specific, sentiment analysis can be performed on social media data from explicit locations. Our approach is to find the sentiments in specific locations. we have analyzed a large data set from which we tried to determine the popularity of a given product in several locations. In order to do this we analyzed tweets from Twitter. Tweets are a reliable source of information mainly because people tweet about anything and everything they do. A number of research works has already been done on twitter data. Most of which mainly demonstrates how useful this information is to predict various outcomes. Our current research deals with outcome prediction and explores localized outcomes. We collected data using the Twitter public API which allows developers to extract tweets from twitter programmatically. The collected data, because of the random and casual nature of tweeting, need to be filtered to remove unnecessary information. Filtering out these and other

problematic tweets such as redundant ones, and ones with no proper sentences was done next. As the preprocessing phase was done in certain extent it was possible to guarantee that analyzing these filtered tweets will give reliable results. Twitter does not provide the gender as a query parameter so it is not possible to obtain the gender of a user from his or her tweets. It turned out that twitter does not ask for user gender while opening an account so that information is seemingly unavailable.

## II. Background Study

**Meenu Dave** et al [1] have done many researches on clustering of healthcare data. Examining huge volume of data, clustering algorithms aid in providing a powerful meta-learning tool. Numerous clustering techniques (including traditional and the recently developed) in reference to large data sets with their pros & cons are being discussed by him. He considered several clustering methods which are presently and widely used for big data analysis. This work delivered an all-inclusive study of the clustering procedures projected in the literature. Analyzing the online streamed data can be considered in the future work. Still there is a huge gap in examining the big data.

**Ramesh** et al [2] proposed a method for sentiment analysis of goods using twitter data in present trends. Their proposed approach is a dictionary based technique i.e. a dictionary of sentiment bearing words was used to classify the text into positive, negative or neutral opinion. Machine learning techniques [12] are not used because although they are more accurate than the dictionary based approaches, they take far too much time performing Sentiment Analysis as they have to be trained first and hence are not efficient in handling big sentiment data.

The main aim of this research by Kim et al. (2013) was to detect short period trends on twitter. Generally this refers to events or

holidays or anything similar which lasts for a while and then loses activity. A problem that was encountered was that simply counting the word frequency was not enough to discover a trending topic. This is because commonly used words such as 'love', 'like' are very common in all tweets and will obviously have a high frequency no matter what set of tweets is analyzed. The approach used by the authors involved plotting the tweet's frequency as a function of time. This resulted in a very helpful pattern where commonly used words had a very much constant frequency throughout the time period being considered but certain keywords showed spikes during certain times. The two events that were analyzed were Easter and weather patterns. The results were very clear as keywords related to Easter spiked on the day of Easter and slowly dropped down in the next couple of days. Similarly areas from weather related tweets were obtained showed sentimental consistency with the actual weather situation in that given area. The results were extremely consistent with the real world events as the output of the weather patterns was accurate as they matched with all the weather forecasts. So the way people are commenting on twitter about the weather was a good indicator of the actual weather turn out again justifying the accuracy of information prediction using twitter. This study shows people talk about certain events on twitter and it turns out that these discussions reflect actual events in those places. So if people in a given location are discussing a product on twitter, their opinions on twitter reflect the actual sentiment in that area.

### III. Methodology

**Data Extraction:** Twitter tweets were used as a data source. It is possible to extract tweets in a large scale from Twitter using the twitter public API that they provide. In our case we used the “twitteroauth” version of the public API by Williams (2012). This version has been implemented in PHP and can be run directly on the local host or on web servers. The query could contain several parameters. Twitter provides a large set of filtering parameters so that a well-defined set of tweets can be obtained. Once the query has been constructed it can be ran by the API and all relevant twitter data will be provided as output in the browser. This data was directly inserted into a MySQL database for the use later on. Each record or tweet that is obtained contains several types of information like user name, tweet id, text etc. But out of those only the text and tweet id were useful to us. Initially the twitter API allowed tweet locations in the form of latitude and longitude to be available with every tweet were the user has made his/her location public. But due to security issues and user complaints this was stopped in 2012. This means that the geographical location from where the tweet was created is not available with the tweet. What twitter does allow on the other hand is the use of location as a filtering parameter in the main query. So in compliance with this restriction we had to extract tweets based on a fixed set of locations. For our research we decided to focus on one nation, USA. We extracted tweets from seven major cities in the USA. The choice of location is very limited mainly due to data availability and language constraints. We decided to go with data from New York, Los Angeles, Boston, Chicago, Dallas, San Francisco and Philadelphia for the experiments. Each major city has a city center, the latitude and longitude that was used to define the city itself. The radius of coverage was chosen based on approximate measure

obtained from using free map tools by Viklund (2015). The radius was picked in such a way so that major parts of the city were covered. Even if a bit of excess was covered it does not really matter as those areas are generally very lightly populated and will not give results anyways. The latitude, longitude and radius are all values assigned to the ‘locations’ parameter in the query build. So now we have multiple data sets each obtained from a different city.

The keyword we choose here in Ailments and Disease. Even though it is possible to analyze any disease in twitter the major issue is that ailment is to be trending in twitter. Remembering that reasonable amount of data about the ailment is available. But our proposed model will obtain results of any required data given that good amount of data available on twitter. So only tweets which contain the ailment will be obtained.

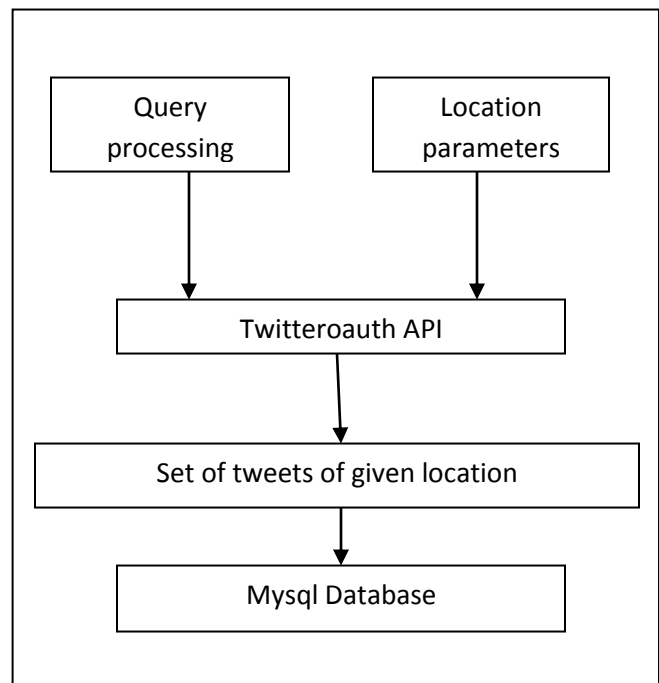


Figure. 1. Data extraction Procedure

**Data processing:**In order to filter out these useless data we mainly used the Stanford Natural Language Processing tool by The Stanford NLP Group (SNLP Group 2015) which is an open source natural language processing tool developed by Stanford University. This tool was used because it gives the grammatical relations between the words in a sentence as output. According to advanced linguistics several such relations are available in the English language. But not all relations are useful in general natural language research. So SNLP has 50 predefined relations which they call dependencies. These dependencies are listed and explained in the Stanford Type Dependencies Manual (SNLP Manual 2015). The reason 50 dependencies are defined in the SNLP is because these are the only word relations which are useful to information analysts, even though linguistics defines several other word relations within a sentence. Out of these 50 dependencies we chose three which will be useful to us.

**Implementation:**In order to assess the sentiment which is present in the tweet a numeric metric is required. This has been done using the tool SentiWordNet (2015), which comes bundled with the SNLP. What SentiWord does is that it takes a word and also the part of speech that a word has in a given sentence. Using the combination of part of speech and the word itself SentiWord gives it a numeric score between  $-1$  and  $1$  where lower value refers to more negative sentiment and higher value refers to higher sentiment. As a tweet text consists of a few words we can take the SentiWord score for each of those words and then sum them up to get a numeric score for each tweet. Another issue here is that SentiWord does not recognize sentences; it only takes words and their corresponding part of speech as input. The part of speech the word will have will depend completely on the sentence itself. So a way has to be devised to map each word in the sentence

to its corresponding part of speech. This was done using Parts of Speech tag extraction. This is also bundled with the SNLP and is used to identify the parts of speech a word has within a given sentence. So each tweet must first be analyzed using the POS tagger which will separate the tweet into individual words and assign a part of speech to it. This is required because by only assessing the word itself it is not possible to determine any sort of opinion, what part the given word plays within a sentence

is always defined by the part of speech it using. In order to map or normalize the POS tags assigned by the POS tagger we had to implement a custom program. Knowing that SentiWord only recognizes nouns, adjectives, adverbs and verbs, any parts of speech other than these three had to be mapped to any of these. An example of the mapping convention would be that if a word is assigned the VBZ tag, which stands for verb in present tense, it will be assigned the Verb tag by the mapper. This set of words along with their normalized POS tags are then sent to SentiWord and the sentiment for each word is calculated and then the individual numeric sentiments are added to obtain a final score for the tweet.

**Table.1 Normalization modeling:**

| Sentiment Score range            | Assigned sentiment         |
|----------------------------------|----------------------------|
| $\text{Score} \leq -0.5$         | Worst                      |
| $-0.5 < \text{Score} \leq 0$ Bad | Score = 0 Neutral          |
| $0 < \text{Score} \leq 0.5$ Good | Score $\geq 0.5$ Excellent |
| $-0.5 < \text{Score} \leq 0$ Bad | Score = 0 Neutral          |
| $0 < \text{Score} \leq 0.5$ Good | Score $\geq 0.5$ Excellent |

#### IV. Experimental Results

To properly understand the trends and variations in sentiments various comparisons were made. The comparisons started at a national level and then became more detailed by the introduction of cities and genders. All of the comparisons have been illustrated using graphs for easy understanding and comparability. The average scores taken are standard averages, not weighted. As the scores were normalized beforehand it was unnecessary to renormalize the average scores. The sentiment percentages were found using the proportion of tweets having a given sentiment among all the other tweets.

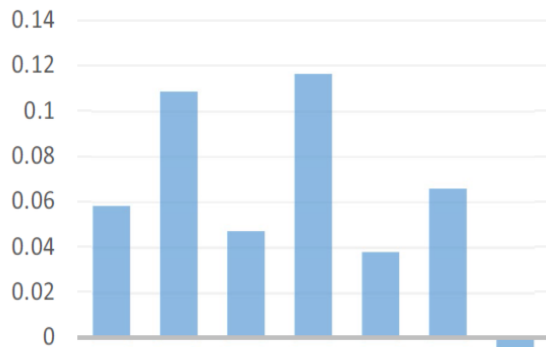


Figure.2. National Average of Ailments

#### V. Conclusion and Future Scope

We discussed a methodology by which it is possible to determine the popularity/ opinion/ sentiment of Ailment. The no of tweets must be significant of ailments. We have to give tags for pre processing of data but Sentiword which we gave for each tweet will be sentiment score.namsor mining tool was used to classify the gender of Diseased.

#### VI. References

1. Akhtar, N. (2014). Social Network Analysis Tools. In Fourth International Conference on Communication Systems and Network Technologies (pp 382–388).
2. Cho, S. W., Cha, M. S., Kim, S. Y., Song, J. C., Sohn, K.-A. (2014). Investigating Temporal and Spatial Trends of Brand

Images using Twitter Opinion Mining. In 2014 International Conference on Information Science and Applications (ICISA)

3. Namsor. (2015). <https://github.com/namsor/namsor-api>. Retrieved July 30, 2015.
4. Renzulli, M. (2015). Top Music Cities in the USA. <http://usatravel.about.com/od/Top-Destinations/ss/Top-Music-Cities-In-The-Usa.htm#showall>. Retrieved July 27, 2015.
5. SentiWordNet. (2015). <http://sentiwordnet.isti.cnr.it/>. Retrieved July 30, 2015.
6. SNLP Manual. (2015). Stanford Typed Dependencies Manual.
7. Viklund, A. (2015). Free Map Tools. <http://www.freemaptools.com/>. Retrieved August 23, 2014.
8. Ekram, T. (2015). Tahmid140/twitter-opinion - mining. <https://github.com/tahmid140/twitter-opinion-mining>. Retrieved July 31, 2015.
9. Hashtags for #election2016 in Instagram, Twitter, Facebook, Tumblr. (2015). <http://top-hashtags.com/hashtag/election2016/>. Retrieved July 30, 2015.
10. Williams, A. (2012). Abraham/twitteroauth. <https://github.com/abraham/twitteroauth>. Retrieved September 19, 2014.

## ABOUT MRGI

Malla Reddy group of Institutions is one of the biggest conglomerates of hi-tech professional educational institutions in the state of Telangana, established in 2001 sprawling over 200 acres of land. The group is dedicated to impart quality professional education like pharmacy, Engineering & Technology, MCA, MBA courses. Our sole objective is to turn out high caliber professionals from those students who join us.

## ABOUT MRCE

Malla Reddy group of Engineering (Formerly CM Engineering College) has been established under the aegis of the Malla Reddy Group of Institutions in the year 2005, a majestic empire, founded by chairman Sri Ch.Malla Reddy Garu. He has been in the field of education for the last 22 years with the intention of spearheading quality education among children from the school level itself. Malla Reddy College of Engineering has been laid upon a very strong foundation and has ever since been excelling in every aspect. The bricks of this able institute are certainly the adept management, the experienced faculty, the selfless non-teaching staff and of course the students.

## ABOUT ICTIMES

ICTIMES started long back with its banner to promote the vision of future technologies that change the trends of life on this planet earth. Under this banner, the Department of Humanities Sciences and Management at MRCE organizes the ICAHSM - International Conference on Advances in Humanities Sciences and Management to provide a scholarly platform to ignite the spirit of Research and bring out the latent potential in teaching fraternity and student community. ICAHSM accommodates major areas like, Automation, Data Analytics and Science, Cloud Computing, Neural Networks, Data Security, Big Data and Business Intelligence.

## ABOUT ICETCS

International Conference on Emerging Technologies in Computer Science (ICETCS-2017) will bring together innovative academicians, researchers and industrial experts in the field of Computer Science to a common forum. The idea of the conference is, for the scientists, scholars, engineers and students from the Universities across the world and the industry as well, to present ongoing research activities, and hence of foster research relations between the Universities and the industry with the rapid development of trends and studies in the fields concerned. ICETCS-2017 will provide a heartwarming platform to researchers, scholars, faculty and students to exchange their novel ideas face to face together.