# Outlier Detection in Secure Shell Honeypot using Particle Swarm Optimization Technique

**M.Sithara**
Department of Computer Science, SRMV College of Arts and Science, Coimbatore-20
Email: sitharamohammed@gmail.com
**M.Chandran**
Department of Computer Applications, SRMV College of Arts and Science, Coimbatore-20
Email: onchandran@gmail.com
**G.Padmavathi**
Department of Computer Science, Avinashilingam Institute for Home Science and Higher Education for Women
University, Coimbatore-43
Email: ganapathi.padmavathi@gmail.com

----------------------------------------------------------------------**ABSTRACT**----------------------------------------------------------------------
**With trends and technologies, developments and deployments, network communication has become vital and inevitable with human beings. On the other side, a network communication without security is powerless. There are so many technologies and developments have been rooted to provide a secure and an efficient means of communication through network. Parallel to this, network threats and attacks are also trendy and much technologized. In order to detect such a kind of threats and attacks, this research work proposes honeypot technology. Honeypot is a supplemented active defense system for network security. It traps attacks, records intrusion information about tools and activities of the hacking process, and prevents attacks outbound from the compromised system. This research work implements a kind of honeypot called Secure Shell (SSH) honeypot. SSH honeypot is a secure communication channel which allows users to remotely control computer systems. With the implementation of SSH honeypot, this research work collects the incoming and outgoing traffic data in a network. The collected traffic data can be then analyzed to detect outliers in order to find the abnormal or malicious traffic. This research work detects outliers from the collected SSH honeypot data using Particle Swarm Optimization technique which belongs to the category of cluster-based outlier detection method. With experiments and results, Particle Swarm Optimization shows best results in detecting outliers and has best cost function when compared to other cluster-based algorithms like Genetic Algorithm and Differential Evolution algorithm.**

Keywords - **Differential Evolution, Genetic Algorithm, Honeypots, Particle Swarm Optimization, Secure Shell**

## 1. INTRODUCTION

Honeypot is a supplemented active defense system for network security. It traps attacks, records intrusion information about tools and activities of the hacking process, and prevents attacks outbound from the compromised system [1]. It appears as an ordinary system doing work, but all the data and transactions are phony. Honeypots actively give way to attacker to gain information about new intrusions. Honeypots can be classified into three different types based on its level of interaction. They are low-interaction, medium-interaction and high-interaction honeypots [3]. This research work implements a honeypot called Secure Shell (SSH) honeypot which belongs to the category of both low and medium-interactions.

Secure Shell (SSH) is a protocol for secure remote login and secures network services over an insecure network using client/server architecture. It is a secure communication channel which allows users to remotely control computer systems [4]. Remote monitoring can be accomplished by deploying SSH in a remote server system. This remote server system can then be monitored using a SSH client system. The data traffic, data transfers, data logs, data failures, data speed, network performance and other criteria in the SSH server can be monitored remotely through SSH client system. Once the SSH protocol is deployed in SSH server system through the SSH client, it is possible to take over the control of SSH server. The idea of deploying honeypots is to lure the attackers and to log the network data. The SSH server logs all the incoming and outgoing network traffic, and this traffic can be monitored and controlled through the SSH client system. In order to authenticate the data log and data transmission between both SSH client and server there are three main authentication methods that are being followed. They are through passwords, keys and hybrid. The deployed scenario in this research work is authenticated through password and SSH server plays the role of honeypot by logging all the incoming and outgoing network traffic. Once the SSH client and server connection is established, the authentication is provided and the network traffic is logged in SSH server. The traffic data can now be analyzed for detection of outliers.

An outlier is a data point which is very different from the rest of the data points based on some measures. Such a data point often contains useful information on abnormal behavior of the system described by data. Abnormal behavior in the network includes malicious traffic, malicious data, extrusions and intrusions like malwares etc. The outlier detection methods try to find these anomalous types of traffic data among the normal traffic data. There are different types of outlier detection methods and they are statistical-based, distance-based, deviation-based, distribution-based, depth-based, density-based and clustering-based. This research work is about detecting outliers in Secure Shell (SSH) honeypot using clustering-based outlier detection methods such as Particle Swarm Optimization (PSO), Genetic Algorithm (GA) and Differential Evolution (DE) algorithms. The results of Particle Swarm Optimization are then compared with the results of Genetic Algorithm and Differential Evolution Algorithm.

This section discussed about the concepts required for the study of research. Section 2 gives the review of literature about Honeypots, Honeypot Mechanisms, Outlier Detection techniques and Clustering techniques. Section 3 explains about the implementation of proposed methodology. Section 4 discusses about the results and outcomes of the experiments and Section 5 concludes the research with future work.

## 2. RELATED WORKS

There are number of security schemes that have been proposed and implemented in literature. Honeypots provide security to the network by logging traffic data and it acts as a normal user system which does not give any kind of information about its presence to the attackers. It is observed that some of the recent works in honeypot technology uses different schemes for implementation and different methods for outlier detection. Some of them are listed below.

Feng Zhang et al. [1] have proposed key components of data capture and data control in honeypot, and have given a classification for honeypot according to security and application goals. This paper have introduced honeypot and honeypot related technologies from the viewpoint of security management for network and have discussed about detection methods, reaction methods, data capture and data store methods. Deployment of virtual honeypots is also discussed here.

Ioannis Koniaris et al. [2] proposed details about two distinct types of honeypots. The first honeypot act as a malware collector, a device usually deployed in order to capture self-propagating malware and to monitor their activity. The second acts as a decoy server to log every malicious connection attempts. They have also shown the usage of honeypots for malware monitoring and attack logging.

Aaditya Jain et al. [3] discussed about the honeypot technology with its classification based on various factors. This paper classifies honeypots based on level of

interactions, purpose and physical presence. As two honeypots have been deployed here, one was used for inbound traffics and another one was for outbound traffics. This paper has also discussed about SSH Honeypots.

Shaik Bhanu et al. [4] presented the results of SSH honeypot operations in which it undertook the web trap of attackers who target SSH service in order to gain illegal services. They have collected the data and analyzed the information from SSH honeypots. Here the focus is on brute-force and dictionary attacks. They have analyzed the data collected from a large number of SSH attacks against a Virtual Private Server (VPS) which was set up as a honeypot to log all malicious activity.

Abdallah Ghourabi et al. [5] proposed a data analysis tool for honeypot router. The proposed tool is based on data mining clustering. They extracted useful features from the data captured by the honeypot router. These data will be then clustered by using the DBSCAN clustering algorithm in order to classify the captured packets and to extract the suspicious data. Suspicious packets will be then verified by a human expert.

Ren Hui Gong et al. [6] presented a genetic algorithm (GA) based approach to network intrusion detection. Genetic algorithm is employed to derive a set of classification rules from network audit data and is utilized as fitness function to judge the quality of each rule. The generated rules are then used to detect or classify network intrusions in a real-time environment.

**Table 1. Comparison of Methods**

| Authors | Techniques | Observations | Year Published |
|---|---|---|---|
| Feng Zhang et al. | Typical honeypot solutions are proposed and the deployment of virtual honeypots is discussed. | Honeypot is not a solution to network security but a good tool supplements other security technologies. | 2003 |
| Ioannis Koniaris et al. | Presented an open source visualization tool which was developed to help security professionals in analysis of data. | Using honeypots with other defense systems such as firewalls and IPS can be a good step | 2014 |
| Aaditya Jain et al. | Proposed honeypots that has multilayer data storing capacity with | Strong control mechanism like IPS is required | 2015 |

| | firewalls and helps in analyzing vulnerable activities. | | |
|---|---|---|---|
| Shaik Bhanu *et al.* | Presented SSH as a medium interaction honeypot and the focus was to detect brute-force and dictionary attacks | Results are very low in percentage. | 2014 |
| Abdallah Ghourab i *et al.* | Proposed a data analysis tool for honeypot router based on DBSCAN clustering algorithm. | False positives are generated with more in numbers. | 2015 |
| Ren Hui Gong *et al.* | Presented a genetic algorithm (GA) based approach to network intrusion detection. | Sometimes the generated rules can be biased to the training data set. | 2005 |

It is observed that Secure Shell (SSH) honeypot logs all the incoming and outgoing traffic and clustering algorithms are applied to detect outliers. This work implements clustering algorithms like PSO, GA and DE techniques. Then the results are compared based on the metrics such as cost function, number of clusters and number of iterations.

## 3. METHODOLOGY

There are different types of outlier detection mechanisms that have been used and studied by the researchers to detect abnormal traffic data. This proposed research work implements clustering-based outlier detection method.

### 3.1. METHODOLOGY OVERVIEW

The methodology implementation starts by creating the network in Cisco Packet Tracer 7.0. Then the incoming and outgoing traffic data in SSH honeypot is collected. Once the data is collected, a subset of data is created. Now the subset of data is clustered using Genetic Algorithm (GA), Differential Evolution (DE) and Particle Swarm Optimization (PSO) algorithms using MATLAB. Finally, the performance of all the three algorithms is compared and the best algorithm is identified based on the experimental results. The methodology overview is given in the figure 1.

### 3.2. CREATION OF NETWORK

The network scenario is created and simulated using the tool Cisco Packet Tracer 7.0. The type of network used here is wired network. The network consists of DHCP server, user systems, routers and switches. DHCP server is deployed in the network to provide web connections between the user systems. It allows the user systems to establish communication between them by using webpages through internet.
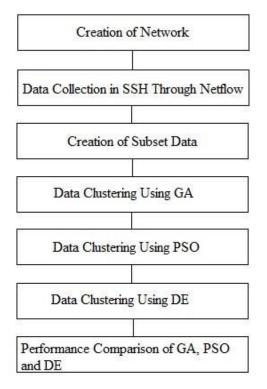


**Fig 1. Methodology Overview**

### 3.3. DATA COLLECTION IN SSH THROUGH NETFLOW

One of the major objectives of the proposed work is to collect data in the network using honeypots. With varying types of honeypots available, this work implements Secure Shell Honeypot (SSH). SSH honeypot can be implemented and monitored remotely. SSH honeypot is implemented in a router and this router acts as a SSH server. Now this SSH server collects traffic logs and this logs are monitored remotely by implementing SSH client in any one of the systems in the created network. Logging and monitoring in SSH client in Cisco Packet Tracer is achieved by applying Netflow. Netflow service provides with access of information about IP flows within the network. Now the traffic data logs in the SSH server is monitored and collected remotely by the SSH client. The collected data includes IPV4 Source Address, IPV4 Destination Address, IPV4 source mask, IPV4 destination mask, next hop

address, source interface , source IP address , destination IP address , destination interface,  protocol, source protocol, destination protocol, no of packets, flow sampler id, idflow direction, source mask, destination mask, counter packets and timestamp.

### 3.4. CREATION OF SUBSET DATA

The creation of subset is to extract only the needed data from the whole set of collected data. The collected data is stored and is accessed using Microsoft Access. The extracted data includes protocol, source and destination protocol, no of packets, flow sampler id, source mask, destination mask, counter bytes and counter packets. The collected data is now clustered to detect the outliers.

### 3.5. DATA CLUSTERING USING GENETIC ALGORITHM

Genetic Algorithm is one of the efficient evolutionary algorithms to detect outliers using clustering technique. Genetic Algorithm can detect misuse and outlier data in the network. This employs metaphor to iteratively evolve a population of high quality individual, where each individual represents a solution to the problem. The operation of Genetic Algorithm starts from an initial population of randomly generated individuals. Then the population is evolved for a number of generations and the qualities of the individuals are gradually improved. During each generation, three basic genetic operators are sequentially applied to each individual with certain probabilities (selection, crossover and mutation). First a number of best-fit individuals are selected based on a user-defined fitness function. The remaining individuals are discarded. Next, a number of individuals are selected and paired with each other. Each individual pair produces one offspring by partially exchanging their genes around one are more randomly selected crossing points. At the end, a certain number of individuals are selected and the mutation operations are applied [6].

### 3.6. DATA CLUSTERING USING DIFFERENTIAL EVOLUTION

The Differential Algorithm is a heuristic algorithm or an evolution strategy. This algorithm is a heuristic algorithm for global optimization and is operated by using decision variables in a real number form. The individuals occurring in this algorithm are represented by real number strings. Its searching space must be continuous. By computing the difference between two individuals chosen randomly from the population, the Differential Evolution algorithm determines a function gradient within a given area (not at a single point). Therefore, this algorithm prevents the solution of sticking at a local extreme of the optimized function. Another important property of this algorithm is a local limitation of the selection operator to only two individuals (parent $(x_i)$ and child $(u_i)$), and, owing to this property, the selection operator is more effective and faster. Also, to accelerate the convergence of the algorithm, it is assumed that the index r1 points to the best individual in the population.  [7].

### 3.7. DATA CLUSTERING USING PARTICLE SWARM OPTIMIZATION

Particle Swarm Optimization algorithm is one among the evolutionary algorithm to detect outliers using clustering technique. Particle Swarm Optimization is similar to Genetic Algorithm in that system is initialized with the population of random solutions. This is unlike to Genetic Algorithm with each of the potential solution is also assigned a random velocity and the potential solutions called particles are then flown through the problem space. Each particle keeps track its coordinates in the problem space which are associated with the best solution (fitness) it has achieved and the fitness value is stored. This value is called pbest. Another best value that is tracked by the global version of the particle swarm optimizer is the overall best value, and its location, obtained so far by any particle in the population. This location is called gbest. The Particle Swarm Optimization concept consists of each time step, changing the velocity (accelerating) of each particle towards its pbest and gbest locations. Acceleration is weighted by a random term, with separate random numbers being generated for acceleration towards pbest and gbest locations. There is also a local version of Particle Swarm Optimization in which, in addition to pbest, each particle keeps track of the best solution, called lbest, attained within a local topological neighborhood of particles [8].

**Table 2. Pseudo Code of Particle Swarm Optimization**

```
begin
Initialize population;
While stopping condition not satisfied do
for j = 1 to no of particles
Evaluate fitness of particle;
Set i=1
for g=1, number of groups
for k=1, number of agents in the group
for n=1, number of dimensions
Random initialization of x^{k,g}_n
Random initialization of v^{k,g}_n
next n
next k
next g
do while (Sub boundary case)
Flag set global best = FALSE
for g=1, number of groups
for k=1, number of agents in the group
Evaluate fitness^{k,g} (i), the fitness of agent k in group g at
instant i
next k
next g
for g=1, number of groups
Rank the fitness values of all agents included only in group
g
next g
for g=1, number of groups
for k=1, number of agents in the group
if fitness^{k,g} (i) is the best value ever found by agent k in
group g then
```

```
p^{k,g}_{best,n}(i) = x^{k,g}_n(i)
end if
if fitness^{k,g} (i) is the best value ever found by all agents then
  g^{g}_{best,n}(i) = x^{k,g}_n(i)
end if
next k
next g
i=i+1
if (i >= sub boundaries iterations) then
Sub boundary case= FALSE
end do
if (Flag set global best = FALSE) then
Flag set global best = TRUE
Rank all the g^{g}_{best,n}(i) and set the actual g^{g}_{best,n}(i)
else
end if
end
```

The overall performance of Genetic Algorithm, Differential Evolution and Particle Swarm Optimization algorithms are measured. The performance for these algorithms is measured in terms of cost function, number of clusters and number of iterations.

# 4. RESULTS AND DISCUSSIONS
## 4.1. EXPERIMENTAL SETUP

The experimental setup of the proposed research work is carried out using the tool Cisco Packet Tracer 7.0. Cisco Packet Tracer is a powerful network simulation program that allows users to experiment with network behavior. Cisco Packet Tracer provides simulation, visualization, authoring, assessment and collaboration capabilities. The proposed work starts experimental setup by placing user nodes, server, routers and switches. Wired channel is chosen as a medium of transmission between devices. Protocols like TCP/IP, UDP and ICMP are used to perform the network simulation. The other simulation parameters used in the experimental setup is given in the below table 3.

**Table 3. Simulation Parameters**

| Simulation Parameters | Values |
|---|---|
| Channel Type | Wired Channel |
| Medium Type | Fast Ethernet, Serial DCE |
| IP Address Type | IPv4 |
| Number of Nodes and Type | 100, Generic-PC-PT |
| Number of Server and Type | 1, Dynamic Host Configuration Protocol (DHCP) |
| Number of Routers and Type | 3, Generic-PT-Empty, 829 |
| Number of Switches | 4, 2960-24TT |
| Protocols | TCP/IP, UDP, ICMP |
| IP Addresses | 192.168.1.1,10 to 34 192.168.2.1,10 to 34 192.168.3.1,2 192.168.4.1,10 to 59 192.168.5.1,2 |

| Server Address | 192.168.5.2 |
|---|---|
| Gateway Address | 192.168.5.1 |
| Honeypot Node (SSH Server) | 192.168.1.1 |
| SSH Client | 192.168.1.34 |
| SSH Password | 123 |
| Netflow Node | 192.168.1.33,34 |

In order to perform clustering and to detect outliers from the collected honeypot data, MATLAB is used. MATLAB (matrix laboratory) is a fourth-generation high-level programming language and interactive environment for numerical computation, visualization and programming.

## 4.2. PERFORMANCE EVALUATION

Performance of algorithms is measured based on the metrics such as Cost Function, Number of Clusters and Number of Iterations.

### 4.2.1. COST FUNCTION

Cost function of a clustering algorithm gives a low cost for the best. It depends on whether all the cluster have same size, content inside each cluster is the same and consecutive clusters do not have the same content.

**Table 4. Comparison of Cost Function**

| Algorithms | Cost Function |
|---|---|
| Genetic Algorithm | 3972.172 |
| Differential Evolution | 4027.1142 |
| Particle Swarm Optimization | **3957.6313** |

### 4.2.2. NUMBER OF CLUSTERS

Number of clusters in a clustering technique is the total number of different clusters with varying group of data obtained from the given data.

**Table 5. Comparison of Number of Clusters**

| Algorithms | Number of Clusters |
|---|---|
| Genetic Algorithm | 3 |
| Differential Evolution | 3 |
| Particle Swarm Optimization | 3 |

### 4.2.3. NUMBER OF ITERATIONS

Iteration is a process or function which is obtained by composing other function with it and process itself with a certain number of times. Iteration applies the same function repeatedly.

**Table 6. Comparison of Number of Iterations**

| Algorithms | Number of Iterations |
|---|---|
| Genetic Algorithm | 200 |
| Differential Evolution | 200 |
| Particle Swarm Optimization | 200 |

The overall comparisons of Genetic Algorithm, Differential Evolution and Particle Swarm Optimization are given in the table 7.

**Table 7. Overall Comparisons of GA, DE and PSO**

| Metrics | GA | DE | PSO |
|---|---|---|---|
| Cost Functions | 3972.172 | 4027.1142 | **3957.6313** |
| | 3985.2668 | 3970.3906 | **3957.6313** |
| | **3817.2913** | 3945.7321 | 3954.1534 |
| | 3993.4905 | 3983.8924 | **3959.6528** |
| | 3979.6956 | 3960.6388 | **3951.7585** |
| No. of Clusters | 3 | 3 | 3 |
| No. of Iterations | 200 | 200 | 200 |

From the comparisons it is clear that the Particle Swarm Optimization has lowest cost function values and it shows best costs when compared to Genetic Algorithm and Differential Evolution algorithms. Hence, Particle Swarm Optimization has the best cost and it is better than Genetic Algorithm and Differential Evolution algorithms.

## 4.3. RESULTS

Sample results for 100 nodes with Secure Shell (SSH) Honeypot and the graph performance of PSO is given in the following figures from 2 to 7.



**Fig 2. Scenario Setup for 100 Nodes in the Network**

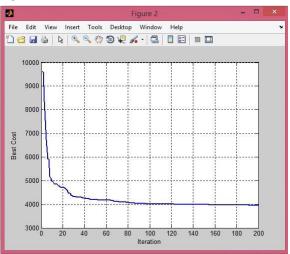

**Fig 3. Traffic Collected in SSH Client from SSH Server**
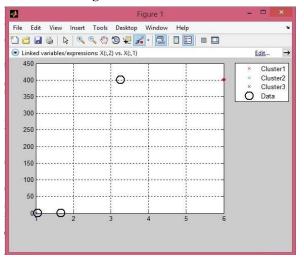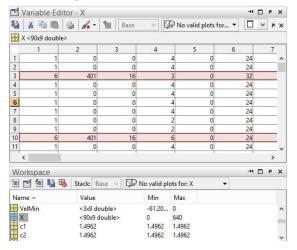


**Fig 4. Cost Function of PSO**



**Fig 5. Number of Clusters for PSO**

**Fig 6. Number of Iterations for PSO**



**Fig 7. Detected Outlier Data from the Collected Honeypot Data**

## 5. CONCLUSION AND SCOPE FOR FUTURE ENHANCEMENT

An efficient network is measured not only by its performance and working rather it is also measured by the terms of security which plays a key role for good, efficient and a secure network. So far, many security techniques have been implemented. This research work has proposed SSH honeypot technology. From the observation, SSH honeypot is efficient in collecting network traffic. From the collected traffic, the research work has implemented clustering algorithms like GA, DE and PSO to detect outliers. Once the outlier data is detected and analyzed, then it will be an easy task to detect the source of the abnormal traffic. The proposed clustering algorithms detect outliers from the collected SSH honeypot data. From the experiments conducted and results obtained, PSO has performed better than GA and DE. PSO has good detection mechanism and best cost function than other two algorithms.

In future, high-interaction honeypots can be used to get a more level of interactions with the attackers. Honeypots can be improved more effectively by combining it with other security mechanisms like firewalls, Intrusion Detection Systems (IDS) and Intrusion Prevention Systems (IPS). Other outlier detection methods like statistical, distance, deviation, distribution, depth and density based can be used.

## REFERENCES

[1] Feng Zhang, Shijie Zhou. Zhiguang Qin, Jinde Liu, Honeypot: A Supplemented Active Defense System for Network Security, *IEEE*, 2003, 231-235.

[2] Ioannis Koniaris, Georgios Papadimitriou, Petros Nicopolitidis, Mohammad Obaidat, Honeypots Deployment for the Analysis and Visualization of Malware Activity and Malicious Connections, *IEEE*, 2014, 1825-1830.

[3] Aaditya Jain, Dr. Bala Buksh, Advance Trends in Network Security with Honeypot and its Comparative Study with other Techniques, *International Journal of Engineering Trends and Technology, 29*, 2015, 304-312.

[4] Shaik Bhanu, Girish Khilari, Varun Kumar, Analysis of SSH attacks of Darknet using Honeypots, *International Journal of Engineering Development and Research, 3*, 2014, 348- 350.

[5] Abdallah Ghourabi, Adel Bouhoula, Data Analyzer Based on Data Mining for Honeypot Router, *IEEE*, 2015, 1-7.

[6] Ren Hui Gong, Mohammad Zulkernine, Purang Abolmaesumi, A Software Implementation of a Genetic Algorithm Based Approach to Network Intrusion Detection, *IEEE*, 2005, 1-5.

[7] Adam Slowik, Application of an Adaptive Differential Evolution Algorithm with Multiple Trial Vectors to Artificial Neural Network Training, *IEEE*, *58*, 2011, 3160-3167.

[8] Russell C. Eberhart, Yuhui Shi, Particle Swarm Optimization: Developments, Applications and Resources, *IEEE*, 2001, 81-86.

[9] Enrique Alba and Marco Tomassini, Parallelism and Evolutionary Algorithms, *IEEE*, *6*, 2002, 443-462.

[10] P. Garcı´a-Teodoro, J. Dı´az-Verdejo, G. Macia´-Ferna´ndez, E. Va´zquez, Anomaly-Based Network Intrusion Detection: Techniques, Systems and Challenges, *Elsevier*, 2009, 18-28.

[11] Roshan Chitrakar, Huang Chuanhe, Anomaly based Intrusion Detection using Hybrid Learning Approach of Combining k-Medoids Clustering and Naïve Bayes Classification, *IEEE*, 2015, 1-5.

[12] Robin Berthier, Michel Cukier, Profiling Attacker Behaviour Following SSH Compromises, *IEEE*, 2007, 1-7.

[13] A. M. Riad, Ibrahim Elhenawy, Ahmed Hassan and Nancy Awadallah, Visualize Network Anomaly Detection by Using K-Means Clustering Algorithm, *International Journal of Computer Networks & Communications (IJCNC), 5*, 2013, 195-208.

[14] Naila Belhadj Aissa, Mohamed Guerroumi, Semi-Supervised Statistical Approach for Network Anomaly Detection, *Elsevier*, 2016, 1090 – 1095.

[15] Amandeep Singh, Navdeep Singh, Review of Implementing a Working Honeypot System, *International Journal of Advanced Research in Computer Science and Software Engineering, 3*(6), 2013, 1007-1011.

**Biographies and Photographs**

*M. Sithara,* She received her M.Sc Computer Science in 2013 from Avinashilingam Institute for Home Science and Higher Education for Women University, Coimbatore. She is pursuing her M.Phil at Sri Ramakrishna Mission Vidyalaya College of Arts and Science, Coimbatore. Her areas of interest are Network Security, Cryptography and Mobile Technology.

*M. Chandran*, He is the Assistant Professor in the Department of Computer Applications at Sri Ramakrishna Mission Vidyalaya College of Arts and Science, Coimbatore. He has 11 years of teaching experience. His areas of interest include Java and Software Engineering.

*G. Padmavathi,* She is the Professor in the Department of Computer Science at Avinashilingam Institute for Home Science and Higher Education for Women University, Coimbatore. She has 29 years of teaching experience and one year of industrial experience. Her areas of interest include Real Time Communication, Wireless Communication, Network Security and Cryptography. She has significant number of publications in peer reviewed International and National Journals. Life member of CSI, ISTE, WSEAS, AACE and ACRS.