

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2021.Doi Number

# Intelligent Behavior-Based Malware Detection System on Cloud Computing Environment

Ömer Aslan<sup>1</sup>, Merve Ozkan-Okay<sup>2</sup>, and Deepti Gupta<sup>3</sup>

<sup>1</sup>Department of Computer Engineering, The University of Siirt, Siirt, Turkey

<sup>2</sup>Department of Computer Engineering, The University of Ankara, Ankara, Turkey

<sup>3</sup>Department of Computer Science, The University of Texas at San Antonio, San Antonio, TX 78249, USA

Corresponding author: Ömer Aslan (omer.aslan@siirt.edu.tr).

**ABSTRACT** These days, cloud computing is one of the most promising technologies to store information and provide services online efficiently. Using this rapidly developing technology to protect computer-based systems from cyber-related attacks can bring many advantages over traditional protection schemes. The protected assets can be any computer-based systems such as cyber-physical systems (CPS), critical systems, desktop and laptop computers, mobile devices, and Internet of Things (IoT). Malicious software (malware) is any software which targets the computer-based system to launch cyber-attacks to threaten the integrity, confidentiality and availability of the data. To detect the massively growing malware attacks surface, we propose an intelligent behavior-based detection system in the cloud environment. The proposed system first creates a malware dataset on different virtual machines which identify distinctive features efficiently. Then, selected features are given to the learning-based and rule-based detection agents to separate malware from benign samples. Totally, 10,000 program samples have been analyzed to evaluate the performance of the proposed system. The proposed system can detect both known and unknown malware efficiently with high detection and accuracy rate. Besides, the proposed method results have outperformed the leading methods' results in the literature. Our evaluation results show that the proposed algorithms along with machine learning (ML) classifiers achieve 99.8% detection rate, 0.4% false positive rate, and 99.7% accuracy. Our proposed system and algorithms may assist those who would like to develop a novel malware detection system in the cloud environment.

**INDEX TERMS** Cloud computing, virtualization, malware detection, behavioral detection, rule-based detection.

## I. INTRODUCTION

Nowadays, there is a tremendous increase in both the amount and severity of cyber-related attacks. In general, different malware variants are the main reason for cyber-attacks. Malware is any kind of software which is designed to exploit computer and network systems' vulnerabilities to perform malicious activities and gain financial benefits. Virus, worm, Trojan, backdoor, rootkits, and ransomware are well-known examples of malware. Each malicious code variant and its family are designed for different purposes. While some malware variants steal sensitive data, others initiate distributed denial of service (DDoS) attacks and allow remote code execution [1]. During sophisticated attacks, more than one malware type and family are used.

Over the years, the number of malware samples have been increasing rapidly. According to business and scientific reports, around 1 million malicious software variants are generated every day. Most of these malware variants are evolving versions of existing malware. Adding new devices to the computer networks every day such as IoT devices, the amount of applications created in a short period of time, and the amount of data created every day in social media also increase the malware-related attacks in the virtual world. On the other hand, new malware variants are using sophisticated concealing techniques such as obfuscation and

packing to hide from detection systems. This makes it almost impossible to identify and classify complex malware with a conventional detection method.

The sophistication of malware attacks, spread methods and economic damage to the world economy have hit the peak recently. According to the researchers, cyber-attacks cause trillion dollars damage to the world economy globally. The evolution of malware related attacks over the years is given in Table 1. It can be seen that back in the early days, viruses and worms were used to launch attacks, but over the years Trojans and ransomwares are mostly used. Attack spread methods that have been evolved, and damages that have been inflicted are changing over the years as well. Social engineering techniques which exploit user trust, software vulnerabilities, malicious emails, and phishing scams are used for attacks' spread methods. Most of the recent attacks steal information from individuals like credit card details on banking systems, encrypt computer data on hard drives to block victims' access to the system, and cause damage to millions of users around the globe.

Malware detection is the process of specifying whether a given program is malware or benign. There are a lot of different methods presented to detect malware which can be categorized as traditional and new detection approaches. Traditional approaches include signature-, heuristic-,

TABLE 1. The evolution of malware-related attacks over the years.

Malware Related Attack	Year	Attack Spread Method	Result
Melissa Virus	1999	It used social engineering techniques to persuade users to click on the email attachment.	It caused billions of dollars in losses across many countries.
ILOVEYOU Worm	2000	It used social engineering to entice users to open the attachment.	It stole users' credentials and infected more than 45 million computer users.
MyDoom worm	2004	It spread by email using attention-grabbing subjects, such as errors, tests, etc.	It launched DDos attacks and allowed remote control.
Zeus Trojan	2007	Malicious emails in the form of spam and drive-by downloads.	It stole login details for social networks, bank and email accounts.
Stuxnet Worm	2010	Attack on the programmable logic unit by stealing source codes.	It took control of industrial processes.
Mirai Malware	2016	It exploited the vulnerability of IoT devices.	It launched DDos attacks.
WannaCry Ransomware	2017	It exploited Windows vulnerability.	It encrypted computer hard drives and affected 150 countries.
Emotet Trojan	2018	Malicious emails in the form of spam and phishing campaigns.	It stole information from individuals, like credit card details on banking systems.
MyFitnessPal	2018	By exploiting software vulnerability.	It affected 150 million users.
LockerGoga Ransomware	2019	Malicious emails, phishing scams and credentials theft.	It completely blocked victims' access to the system and caused millions of dollars in damage.
CovidLock Ransomware	2020	It exploited users' trust by providing statistical information about COVID-19.	It encrypted data on Android devices and denied data access.

behavior-, and model checking-based while new approaches include cloud-, deep learning-, and mobile devices-based detection [2].

As it is known, the signature-based detection approach performs well for known and different versions of the same malware, but it fails to detect unknown malware which has a completely different signature. Behavior-, heuristic-, and model checking-based approaches may detect a significant portion of the zero-day malware. However, they cannot detect new malware which uses advanced packing techniques. Although deep learning- and mobile devices-based detection approaches improve the detection rate (DR) for mobile devices to a certain degree, they fail to detect malware which seems completely different from the previous version [2].

Cloud-based malware detection approach brings several benefits over the other approaches. The cloud computing environment provides easy access, on-request storage, more computational power and considerably bigger databases while decreasing the cost. Multiple execution traces of the same malware have been collected by using different virtual machines (VMs) and servers [3]. Cloud environment improves the DR for personal computers, mobile and IoT devices. In addition, various detection algorithms can be implemented on different servers. Using several algorithms improves the detection performance while decreasing the false positive and negative rates.

In this paper, an intelligent behavior-based malware detection schema is proposed in the cloud computing environment. The cloud-based detection schema consists of two parts including feature extraction and detection phases. A client submits a suspicious file over the computer network and receives the analysis result from the server which shows whether the given suspicious file is malware or not. The suggested cloud-based system provides the following contributions:

- Suggested model creates a malware dataset with fewer features than known models do.
  - First, several system calls are mapped into relevant behaviors.
  - Second, relationships are determined among the behaviors.
  - Finally, features are extracted from behaviors which have semantic relationships between them.
- Learning-based detection engine is used to separate malware from benign.
- Rule-based detection engine is used to determine malware from benign as well.
- The proposed schema can detect both previously known and unknown malware.
- Proposed model detection and accuracy rate are measured higher than known models.

The rest of the paper is organized as follows: Section II explains the cloud computing environment. Related work is summarized in section III. Proposed methodology and case study are defined in section IV and V. The results and discussion are presented in section VI. The limitations and future works of the proposed model are given in section VII. Finally, the conclusion is given in section VIII.

## II. CLOUD COMPUTING

Cloud computing provides various computing services over the Internet [4]. By using the cloud, different users and businesses are storing their data remotely in the data centers instead of using their own local storage. This makes the data available from anywhere, anytime, and from any devices. Cloud environment provides data storage, servers, VMs, databases, networking, and software. Typical cloud computing environment and its infrastructure can be seen in Figure 1.

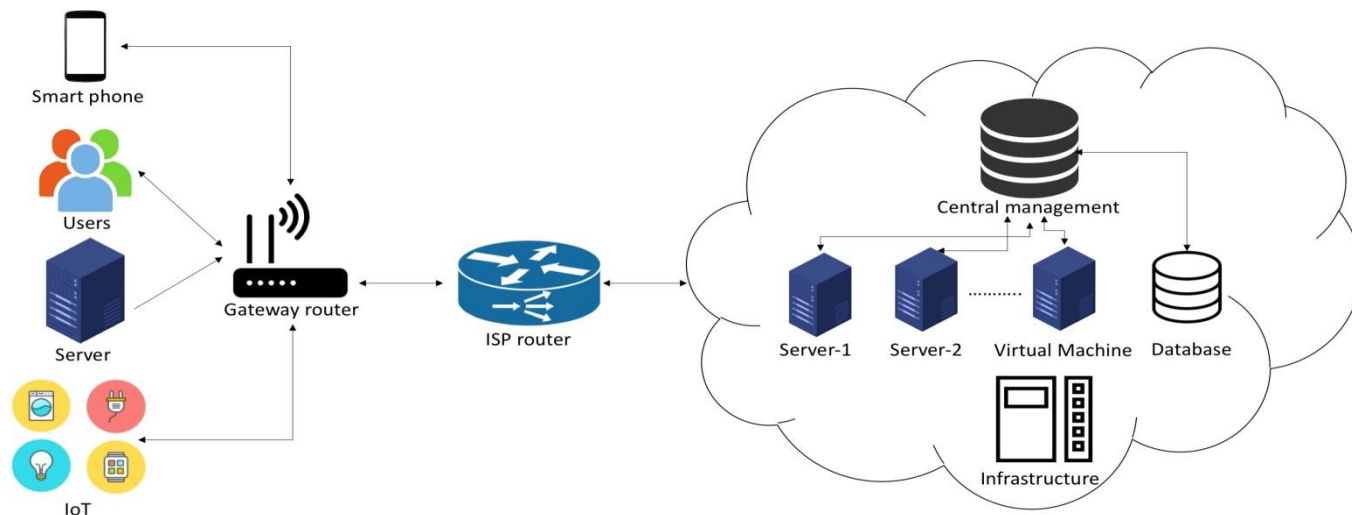


FIGURE 1. Cloud computing environment.

These days, cloud computing is used almost everywhere. Cloud computing brings several advantages over traditional kinds of information storage such as easy access, pay as you go, increase in speed, efficiency and performance, and decreasing cost. There are many cloud service providers including Amazon, Microsoft, Google, IBM, Rackspace, and Verizon.

Virtualization is the foundation of cloud computing. Without virtualization, cloud computing will be incomplete in many ways and will not be used as much as it is used today. Virtualization uses a hypervisor (Virtual machine monitor) to separate the operating system from the computer hardware. This allows us to use multiple operating systems that run on the same physical machine. Server virtualization on the cloud site brings many advantages including: less equipment, lower energy consumption, increase in server uptime, faster server provisioning, redundancy, and improvement in disaster recovery.

Cloud environment provides different types of services including infrastructure as a service (IaaS), platform as a service (PaaS) and software as a service (SaaS) [5]. For businesses' needs, one cloud service can have more advantages than others. There are also different kinds of cloud deployment models including public, private, community and hybrid clouds [6]. Public clouds are owned by third-party cloud providers and services, and available to the general public. Users can access the cloud services by using a web browser. In the public cloud, different organizations share the same infrastructure which may disclose sensitive data. Private clouds provide physical infrastructure and services for specific organizations. Infrastructure can be physically located in the company's data center or third-party's data center. Community clouds are created for exclusive use by a specific community. Hybrid clouds combine public and private clouds together. Hybrid clouds provide more deployment options, security and flexibility.

Cloud computing brings several advantages over traditional storing schema:

1. Users and organizations can store and back up their data in an efficient manner.

2. Regular users and organizations can access their data from any device, anywhere, and any time via browser or application.
3. Organizations can subscribe only for needed services.
4. Cloud environment provides cost-savings from small businesses to big organizations.
5. Cloud environment provides more storage space, computational power and considerably bigger databases.
6. It eliminates the need for onsite equipment, maintenance, and management issues.
7. It enables rapid response when increasing data volume requirements.
8. It reduces cost for physical resources, energy, and personnel training needs.

There are some issues which need to be addressed in the cloud computing environment:

1. Users lose control over their data.
2. Sensitive and top-secret data can be disclosed.
3. On public clouds the same physical resources are used for different organizations which also raises the security issues.
4. Data can be lost because of internal bugs, natural disasters and other reasons.
5. If the Internet is slow, it takes a lot of time to access the data.
6. Real time monitoring is not possible for all locations.

Using the cloud computing environment for malware detection brings many advantages. Cloud environment presents more computational power and much bigger databases [6]. Different methods and algorithms can be implemented in the cloud such as machine learning (ML), data mining, and deep learning. Multiple execution traces of the same malicious software can be collected. It enhances the detection performance for personal computers, mobile and IoT devices.

### III. RELATED WORK

Cloud computing technology has been rapidly developing recently owing to some advantages, including simple accessibility, lower costs and scalability. Due to the innovation and convenience of cloud technology, the interest and use of cloud computing have increased among users as well as researchers. Cloud computing plays an important role in the protection of computer systems such as smart cyber-physical systems (CPSs) [7], IoT devices [8] and personal computers from several cyber-attacks, especially from malware attacks. In literature review, several studies are presented to detect malware in the cloud environment. Different malware detection approaches were analyzed based on the main idea, algorithms that are used, and feature extraction methods. Well-known cloud-based malware detection methods are summarized in Table 2. The common goal of all these studies is to identify malware by increasing *DR* while decreasing misclassification rates. When these studies are examined, it is seen that although each detection method has its own superiorities and performs better for particular datasets in the cloud, none of them could detect all malware.

Martignoni *et al.* [9] introduced a new framework to support dynamic behavior-based malware analysis based on cloud computing. The proposed framework is based on two assumptions. First, the security lab has no limit on available computing resources and can take advantage of hardware features. Second, end users' environments are more literal and nonhomogeneous than synthetic environments and are therefore more suitable for analyzing malware. They performed an empirical prototype to approve their ideas and integrated it into their existing behavior-based malware detection system. The evaluation results showed that the proposed framework enables security labs to advance the integrity of the analysis while performing a detailed analysis of the program's behavior without computational costs for end-users. On the other hand, the proposed framework increased the security issues and inclined to several detection and hijacking attacks. Solving security-related problems and applying a framework resistant to evasion attacks will improve framework performance.

Cha *et al.* proposed a new malware detection system named SplitScreen [10]. It is a distributed malware detection system that uses a supplemental screening step before the signature matching stage. SplitScreen's two-stage screening step is separated into client-server processes. The suggested method was implemented as an extension of ClamAV, which increases scanning throughput with more than 2x the signature set using half of the memory. As the authors mentioned that the acceleration and memory savings of SplitScreen improves when the number of signatures increases. The proposed method is scalable with a wide range of low end consumer and handheld devices. Since only one server is used on the cloud side, it would be better to optimize server efficiency and load some work on the client side.

Win *et al.* studied cyber-attacks targeting the virtualization infrastructure underlying cloud computing services [11]. They proposed a malware and rootkit detection system that defends guests from several attacks. The system was combined with Support Vector Machines (SVM) based external monitoring on the host, with system call monitoring

and system call hashing in the guest kernel. The design of the proposed approach is to perform a system that detects the entity of attacks against guests in real time without the demand for a signature database. They indicated the efficiency of the proposed approach by appreciating it against well-known user-level malware and kernel-level rootkit attacks. According to the authors, the implemented solution eliminated the demand to use a signature database for malware classification.

Gupta *et al.* proposed a novel model for malware detection in the cloud [12]. The aim of this study is to detect malicious activities with some techniques and warn guest VMs about it. In this paper, DNA sequence detection process, the symbolic detection process and the behavioral detection process are combined. During the DNA sequence detection process, they extracted the DNA sequence from a file to detect malware. In the symbolic detection process, they clustered files according to file formats and used symbols to detect malware files. During the behavioral detection operation, they observed the behavior of the file and determined whether it was a malicious program using the Anubis sandbox. A prototype of the proposed approach (PMDM) is partially implemented on the Eucalyptus. According to the authors, PMDM is inexpensive, needs less runtime, and ensures well performance for large numbers of files compared to other known systems. However, this study can be improved further by using a bigger dataset.

Rakotondravony *et al.* categorized attacks in the IaaS cloud that can be analyzed using VMI-based mechanisms [13]. They focused on attacks that directly scramble VMs deployed in the IaaS cloud. The classification methodology takes into account the target, source and direction of attacks. They provided an overview of attacks where each actor could be threatened in the environment. They defined a common IaaS cloud scenario as a range of three different elements: cloud provider, external entity, and VMs. First, they summarized the distinct properties of attacks classified in the literature in respect to attack complexity, security effect, and suggested defense metrics. They then analyzed statistics on virtualization vulnerabilities misused by attacks, noticed them in public databases, and highlighted their evolution over time. Finally, they presented the economic impact of attacks on business processes. This study allowed several actors in a cloud scenario to evaluate different malware attacks and as a result design sufficient detection and mitigation mechanisms based on virtual machine introspection. Paper can be further enhanced by focusing not only on attacks involving direct VMs, but also on other types of attacks.

Sun *et al.* [14] explained a cloud-based malware detection system called CloudEyes. CloudEyes ensures effective security and data privacy for limited resource devices. Suspect bucket cross-filtering, a novel signature-based detection system for the cloud server, has been proposed based on reversible structure. It can provide retroactive and correct processing of malicious signature fragments. A scanning tool is applied to quickly define the file content suspicion with respect to the summary of the reversible sketch for the client. An interaction mechanism has been

TABLE 2. Summary of cloud-based malware detection methods.

Paper	Proposed Method	Goal/Success	Year
Martignoni <i>et al.</i> [9]	Presented a new framework based on cloud computing for dynamic behavior-based analysis.	It provides security labs to enhance the accuracy of the analysis.	2009
Cha <i>et al.</i> [10]	Anti-malware system called SplitScreen.	It increases detection while decreasing memory usages.	2011
Win <i>et al.</i> [11]	A malware and rootkit detection system.	It removes the necessity of using a signature database.	2015
Gupta <i>et al.</i> [12]	A novel malware detection model on cloud architecture.	PMDM is inexpensive, takes less working time and presents well performance for large numbers of files.	2016
Rakotondravony <i>et al.</i> [13]	Attack classification in the IaaS cloud that can be examined using VMI-based mechanisms.	It lets distinct actors in a cloud scenario evaluate different malware attacks and design sufficient detection and mitigation mechanisms based on VMI.	2017
Sun <i>et al.</i> [14]	CloudEyes, which presents effective and confident security services for limited resource devices.	It is effective, practical, and saves time and data storage when detecting malware.	2017
Babu and Murali [15]	Improved and designed an intermediary malware protection in cloud environments.	It protects the cloud from malware transportations, and decreases time and cost.	2017
Xiao <i>et al.</i> [16]	Malware detection scheme with Q-learning.	It increases the accuracy, while reducing the latency.	2017
Abdelsalam <i>et al.</i> [17]	Malware detection approach in cloud infrastructure.	The 2-D CNN model achieves the 79% accuracy rate, and 3-D model notably enhances to 90% the accuracy rate.	2018
Mirza <i>et al.</i> [18]	An energy effective hosting model in the cloud environment.	It shows important energy efficiency with regard to CPU usage by the hosting model.	2018
Mirza <i>et al.</i> [19]	Cloud-based energy effective hosting model for an intelligent malware detection	It performs better than the conventional antiviruses.	2018
Shen <i>et al.</i> [20]	Malware detection system implemented by an intrusion detection system with cloud and fog computing.	It decreases delay of data traffic as well as data transfer overhead.	2018
Zhou and Yu [21]	A cloud-assisted model for malware detection and the dynamic system against malware propagation.	It can prevent the spreading of malicious codes obviously and efficiently and is convenient to the resource limited WMS.	2018
Yadav [22]	Consolidated WFCM-AANN malware detection technique.	It successfully determines the malicious software with high detection precision thereby outperforming existing classifiers.	2019
Indirapriyadarsini <i>et al.</i> [23]	Random and some other modeling like KNN, Logistic Regression (LR), etc.	It has come up with the unique solution by working with ML and cloud computing simultaneously to determine the legitimacy of the file.	2020
Deyannis <i>et al.</i> [24]	Cloud-based malware detection solution called TrustAV.	It can protect the transmission and processing of user data even in distrusted networks.	2020

designed to protect the data privacy and decrease consumption of communication. The client transmits the coordinates of the suspicious file segments rather than the entire file content. They evaluated the performance of CloudEyes using both suspicious and normal traffic. According to the authors, the test results showed that CloudEyes is effective, practical and outperforms other existing systems in terms of time usage and consumption of communication. However, *DR* and accuracy can be further improved. In addition, some methods can be applied to reduce the data size to optimize storage and matching performances.

Babu and Murali designed a protection system against malware spreading in cloud environments [15]. This investigation presents several layered protections to address the problem and creates a two-layered epidemic model for preventing spread of malware from network-to-network. In the proposed system, they designed the malware detection system for various cloud servers using a middle monitoring server, allowing scanning, detection and removal of malware before transferring to cloud servers. According to the authors, this study secures malware transfers to the clouds and saves time and cost.

Xiao *et al.* analyzed the malware detection game based on cloud in which mobile devices upload the traces of their application to security servers over access points or base stations in dynamic networks [16]. Q-learned malware

detection system was designed for a mobile device. The aim of this study is achieving the optimal payload transfer ratio without knowing the trace creation and radio bandwidth model of different mobile devices. They used the Dyna architecture to enhance performance and a post-decision learning method to speed up the reinforcement learning phase.

Abdelsalam *et al.* [17] presented a malware detection method based upon Convolutional Neural Network in cloud computing environments. They used a standard 2-D CNN, training on data existing for each of the processes in a VM acquired through the hypervisor. They improved CNN classifier accuracy rate by using a new 3-D CNN, which considerably helps decrease mislabeled samples while training and data collecting. They performed experiments on collected data by working varied malware on VMs. The 2-D CNN model achieves the 79% accuracy rate, and 3-D model notably enhances to 90% the accuracy rate. This study could be improved with increasing the experiments scale by examining more malware binaries.

Mirza *et al.* [18] proposed a combination of ML techniques applied on large dataset. The paper mainly focused on two important goals including higher *DR* and low resource consumption. They extracted a group of features from the dataset including malicious and normal files, and implemented a SVM, boosting, and decision tree on the

decision tree to obtain the highest possible detection rate. Boosting of the decision tree classifier showed a better performance in the assessment of CloudIntell. They also introduced a scalable cloud based architecture hosted on Amazon Web Services (AWS). They tested proposed methodology on different scenarios. According to the authors, their methodology produced high results with lowest energy consumption. Besides, implementing the boosting algorithms on a real-time platform is difficult and training the classifier with large amounts of data takes a lot of time and computation. In another study, Mirza *et al.* [19] suggested an energy efficient hosting model which consists of distinct components of Amazon's cloud services to improve a unique and scalable model. This research examined the set benchmarking numbers and known antiviruses for the cloud based hosting model. According to the paper, the proposed approach not only was successful for the hosted detection framework, but also performed optimally better than traditional antiviruses. However, the malware detection framework and hosting model can be improved further by integrating the intrusion detection mechanism to be assisted by the cloud based engine.

Shen *et al.* [20] explained a malware detection structure implemented by a cloud and fog computational intrusion detection system (IDS) to accomplish the IDS spreading problem in smart objects. There are three main contributions of the proposed study. First, they suggested an intrusion detection approach to detect malicious software in fog cloud based IoT networks. Second, they introduced a multistage privacy-preserved game which is based on confidentially leakage evaluation of smart objects to detect malware in IoT networks. Finally, they explained a framework to integrate the presented game into fog cloud based IoT networks using the right detection strategies. According to the authors, the proposed model fulfilled the large data processing requirement caused by the greatly increasing number of smart objects and the reduced data traffic latency as well as the data transfer overhead.

Zhou and Yu suggested a cloud assisted model for the dynamic differential game against malware spread and malware detection [21]. In the suggested model, first, a malware detection model based on SVM is created by sharing data on the security platform in the cloud. Second, the number of malware infected nodes that physically infect sensitive nodes is calculated according to attributes of wireless multimedia system (WMS). Finally, the transition of states between WMS devices is described by the changed epidemic model and Hamilton function has been presented to simplify the saddle point solution. Also, a target cost function and dynamic differential game has been sequentially derived for the Nash equilibrium between the WMS system and malware. According to the paper, obtained results demonstrated that the proposed algorithm is capable of suppressing the spread of malicious code clearly and efficiently and is suitable for resource-constrained WMS.

Yadav explained a unified WFCM-AANN malware detection approach to identify malware on the system [22]. The presented study consists of 2 modules, including classification and clustering. In the clustering module, the input data set is obtained in clusters by applying the WFCM (Weighted Fuzzy C-mean) algorithm. In the classification module, the centroid from the clusters is given to the discontinuous Auto-Associative Neural Network, which is applied to characterize whether information is intruded or not. The author claims that the proposed classifier successfully determines malware with high detection rate and therefore outperforms the existing classifiers.

Indirapriyadarshini *et al.* [23] proposed a machine learning-based detection technique on the cloud environment. They first used random modeling to get the worst log loss and then used some modelling such as KNN, LR etc. They then looked at the log loss of each algorithm and determined whether it was a perfect model. Finally, they deployed the ML model with the user interface on the cloud AWS. According to the authors, they had found a unique solution by working simultaneously with ML and cloud computing to determine the legitimacy of the file. However, this study can be enhanced by applying different data mining techniques for feature selection or by implementing new learning models.

Deyannis *et al.* [24] presented a cloud based malware detection solution named TrustAV. This solution is based on a pattern matching technique to determine contaminated data. TrustAV transmits the processing of malware analysis to a remote server and it is proposed as a cloud based solution. According to the paper, TrustAV can protect the transmission and processing of user data even in distrusted environments. In addition, TrustAV also uses a variety of techniques offered by Intel SGX technology to overcome general performance loads and limit the risk. However, there is no real data to evaluate the proposed cloud-based TrustAV solution.

When the existing studies are examined for the cloud environment, it is shown that various techniques such as preprocessing, feature reduction and extraction, and ML algorithms have been applied on the dataset to detect malware with high accuracy. When the proposed malware detection methods in these studies are evaluated, it is seen that preprocessing and feature selection stages before implementing ML algorithms improve the performance. In addition, some ML algorithms may perform better than other algorithms according to the size, distribution, and number of features used in the dataset. It can be concluded that the cloud based malware detection approach and its methods improve the detection performance for computers, mobile, and IoT devices with bigger malware databases, and heavy computing resources. Other benefits of cloud based detection are configurations, installations and regular updates. However, some portions of malware could not be detected by using a cloud based detection approach and its methods. To build a more effective detector on the cloud site, hybrid-based detection approach, which combines behavior-, model checking-, and using deep learning

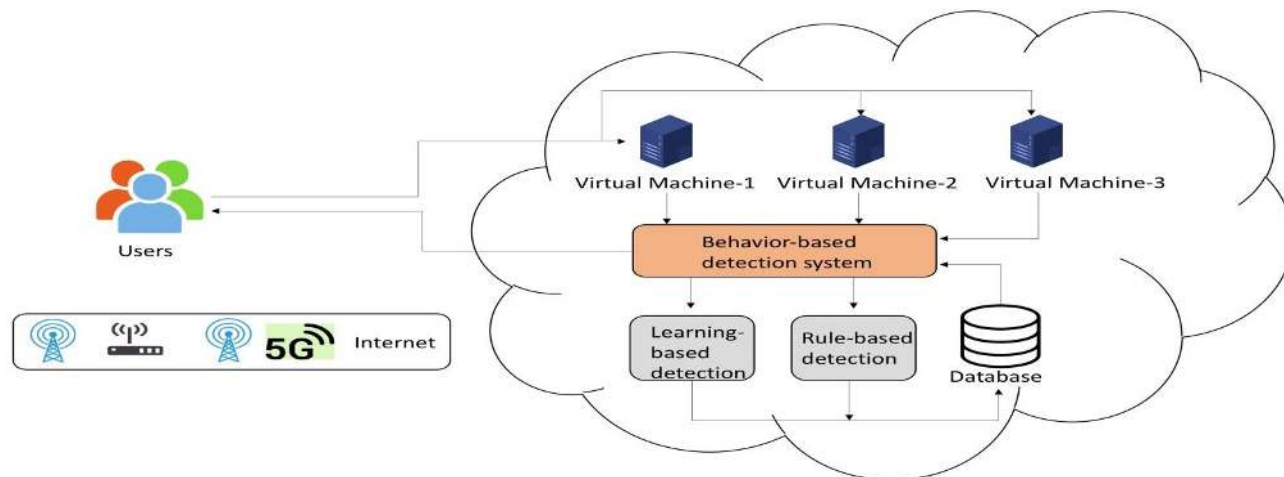


FIGURE 2. Proposed cloud-based malware detection architecture.

altogether on the cloud environment can be a promising method. We believe that cloud-based malware detection approach is still at the early stage, and there needs to be more studies in this area to see effectiveness of the cloud at detecting malware.

#### IV. PROPOSED SYSTEM

This section explains the proposed system including model architecture, dataset, features, and detection methods in detail.

According to our proposed system, the user submits a suspicious file to the cloud environment by using a computer network. Then, the submitted file is executed in different VMs and execution traces are gathered by using relevant dynamic tools. Generated execution traces are collected on behavior-based detection agent and behaviors are generated. Related behaviors are grouped according to the predefined rules in order to create features. When features are being created, a proposed cloud-based behavior centric model (CBCM) is used. After that, most discriminative features are selected by suggested algorithms and selected features are sent to the detection agents including learning-based detection and rule-based detection. In learning-based detection agent, selected features are trained by using machine learning algorithm such as logistic model trees (LMT), C4.5 (J48), random forest (RF), simple logistic regression (SLR), sequential minimal optimization (SMO), and k-nearest neighbor (KNN). On the other hand, in rule-based detection agent, features are evaluated based upon predefined features sets. Based on learning- and rule-based detection agents, each sample is marked as malware or benign and stored in the database. The analysis result is sent back to the user which shows whether the suspicious file is malware or not.

##### A. PROPOSED SYSTEM ARCHITECTURE

The system architecture of the cloud-based malware detection model is presented in Figure 2.

##### B. BEHAVIOR CREATION, FEATURE EXTRACTION AND SELECTION

Analyzing malware manually and extracting features require a lot of time and manpower. Therefore, there is an urgent need to build a system which can automatically analyze the malware and extract features. Although malware performs actions which are related to one another, it also carries out unrelated actions to hide its real behaviors. Because of that, it is vital to determine the interrelated actions and extract real features while creating a dataset. Automatic dataset creation models such as leading methods in the literature and the  $n$ -gram are lacking in this regard because those methods generate too many features as well as unrelated features. These deficiencies increase the detection time while decreasing the  $DR$ . For these reasons, the CBCM model is proposed in this study which creates features and selects features effectively.

Overview of malware analysis process can be seen in Figure 3 and feature creation and selection process can be seen in Figure 4. To create features for each suspicious file, an executable file is analyzed by using dynamic analysis tools such as Process Monitor, API Monitor, Process Explorer, Autoruns, and Debuggers in different VMs. Then, execution traces are collected and sent to the behavior-detection agent. In detection agent, behaviors and features are being created by using the CBCM model. Behavior creation, feature extraction and feature selection are intertwined in the proposed model.

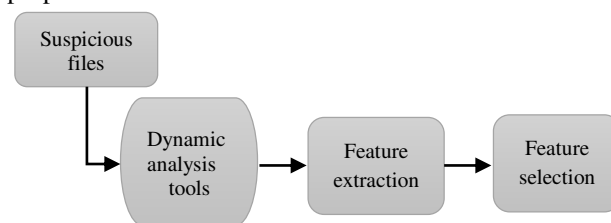


FIGURE 3. Malware analysis process.

The CBCM model is a modification of the subtractive center behavior model which was proposed in our previous work [25]. While creating features, malicious behavior patterns are determined. The malicious properties are those which can be frequently seen in malicious codes but rarely seen in non-malicious program samples. The goal is to gather the most important properties. To identify malicious properties: system calls, system call paths, system resource

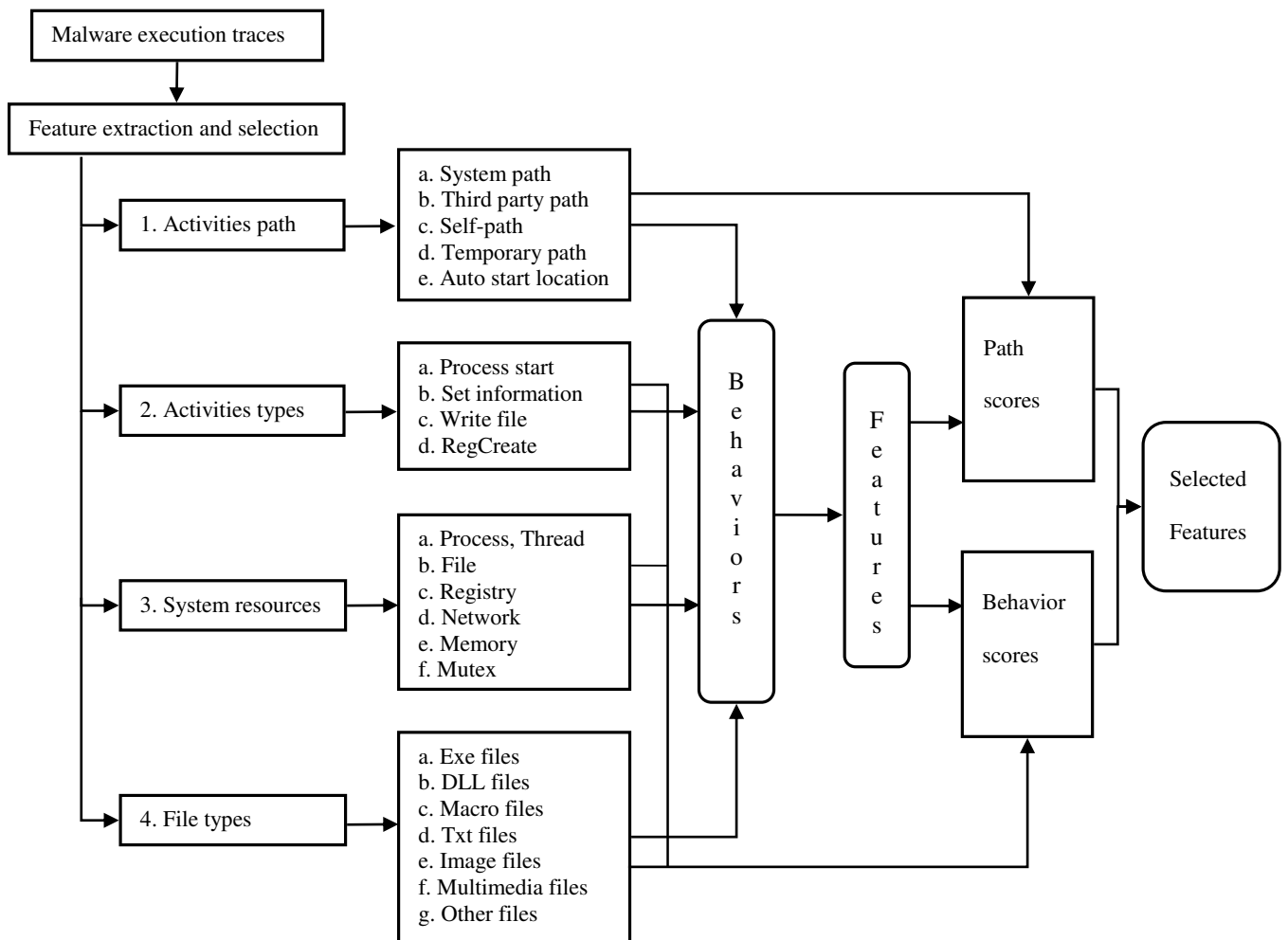


FIGURE 4. Malware feature creation and selection process.

types, and file types are taken into consideration (Figure 4). In this way properties that can distinguish malware from benign samples are obtained. To create behaviors, necessary relationships are established among the system calls. One or more system calls, which can represent meaningful activity, create a behavior. The behavior creation algorithm can be seen in Algorithm 1. It takes the list of activities (system calls)  $D_1$  as an input and generates a list of behaviors  $D_2$ . During the behavior creation, first based on the activities that are performed, action states ( $\psi$ : active (A) – passive (P)) are established. Then, action paths ( $\mu$ : self (SF), system (S), third party (TP), temporary (T), and auto start (AS)), where the system calls perform, are determined. After that, consecutive system calls, which can represent behaviors, are determined. The same consecutive system calls and ending system calls, which represent the actions, are excluded when behaviors are generated. Finally, obtained behaviors are written to do  $D_2$ .

Unlike  $n$ -gram, the CBCM uses twenty consecutive behaviors when creating properties. If there is a relationship between first and tenth or twentieth behaviors, it can create a property. The importance of the properties is determined as follows:

1. Paths that system calls are performed on.
2. Resources that system calls are performed on.
3. File types that are created.

### 1) Paths that system calls are performed on

We divided performed system calls locations into five categories, which is shown in Figure 4 including system, third party, self, temporary, and auto start locations. Each path is further divided into subfolders and path scores are established. These paths are used during the behavior and feature creation. Example list of system activities path can be seen in Table 3 and path score calculation for behaviors can be seen in Algorithm 2.

When the system call is performed in system folder, following criteria are taken into consideration (Algorithm 2):

- (i) If an analyzed program sample interacts with the operating system files and directories in order to work properly, these interactions are evaluated normal. Most of the time, these interactions are provided by system DLLs, background processes, and system services. These interactions are considered to be normal, so the risk level of these interactions will be low or moderate depending upon the other information
- (ii) If analyzed program tries to inject some codes to the system DLLs and exes including kernel32.dll, advapi.dll, svchost.exe, winlogon.exe, etc., those actions are considered to be malicious and the risk level of these interactions will be high.

When the system call is performed in third party folder, following criteria are taken into consideration:



- (i) Most programs need third-party software to run properly. If an analyzed program sample needs other programs to run properly, the risk level of these interactions will be low or moderate depending on the other information.
- (ii) However, if actions on the third party files and directories that are not related to the performed sample, those actions are considered malicious and the risk level of these interactions will be high.

**Algorithm 1** Behavior Creation

**Input (D<sub>1</sub> file):** List of activities/system calls

**Output (D<sub>2</sub> file):** List of behaviors

```

1: for each system call in D1 i do
2:   if D1[i][state] == 'active' then
3:     ψ = 'A'
4:   else
5:     ψ = 'P'
6:   end if
7:   if process.name==d1.filename then
8:     μ = 'SF'
9:   elif path=='system' then
10:    μ = 'S'
11:  elif path=='thirdParty' then
12:    μ = 'TP'
13:  elif path=='temporary' then
14:    μ = 'T'
15:  elif path=='autostart' then
16:    μ = 'AS'
17:  end if
18:  if D1[i][system call]! =D1[i+1][system call] then
19:    if D1[i][system call]! = 'ending system call' then
20:      write. D2 (D1[i][system call])
21:    end if
22:  end if
23:  if D1[i][system call] == D1[i+1][system call] then
24:    if D1[i][path]! =D1[i+1][path] then
25:      if D1[i][system call]! = 'ending system call' then
26:        write. D2 (D1[i][system call])
27:      end if
28:    end if
29:  end if
30: end for
    
```

generates normal actions that cannot be categorized as malicious. For those actions, the risk level will be low.

- (ii) However, if an analyzed program sample performs registry and network-related actions within some files or copies its own file content to other files, it is considered to be malicious and the risk level of these interactions will be high.

Temporary folder and auto start locations are other paths which need to be considered. This is because most of the malware types use temporary folders when performing malicious actions, and use auto start file-registry locations to become persistent in the system.

- (i) If an analyzed program sample is using temporary folder or auto start locations, these interactions are considered to be malicious and the risk level of these interactions will be fairly high.

**2) Resources that system calls are performed on**

In order to create behaviors and related properties, system resources are split into following categories: process, thread, file, registry, network, memory and mutex. During the determining behaviors and properties, usually the same types of resources are considered. When malware first runs, it creates some processes and threads to perform malicious actions. These processes and threads can make some changes on files, registry entries, memory and mutexes, or can connect other networks to exchange some sensitive data. Because of that each action which is carried out on those system resources is analyzed deeply during the feature creation.

**3) File types that are created**

Created file types are also taken into consideration during feature creation. We considered portable executable (exe, DLL) and macro files slightly more dangerous than other files including txt, image, multimedia files, etc. This is because several malware variants create exe extension files or inject malformed program codes into DLL files to launch attacks.

The feature extraction algorithm is presented in Algorithm 3. When features are generated from behaviors, twenty consecutive behaviors are considered. In this phase, features, feature action types, and path scores are calculated. The same types of system resources (file, registry, mutex, network, etc.) are considered when determining property relationships. In addition, different resources create features if relationships can be established among them. Path scores and action states (AA, AP/PA, PP) are used during the feature selection.

The feature selection algorithm is presented in Algorithm 4. First, the frequency of each property is calculated. During the feature frequency calculation, we try to reduce the number of different features as many as we can. The features of the same name, which occur on the same resource type and have the same path score but different locations, are combined with the same property and the frequency is increased. For instance, even though ReadFileWriteFile (\...\path1\, pathScore = 'x') and ReadFileWriteFile (\...\path2\, pathScore = 'x') have been performed in different locations and instances, they set to the same feature and frequency is increased.

After frequency calculation is finished, features are selected based upon path scores and action states. If the path score is moderate, high or very high, related property is chosen.

**TABLE 3.** Malware execution trace system calls path (The list is abbreviated).

Action System Path
'c:\windows', system folder
'hklm\system', system folder
'c:\windows\system32', system folder
'c:\program files', third party folder
'c:\program files (x86)', third party folder
'\...\suspicious file', self-folder
'c:\users\...\startmenu\...\startup', auto start location
' hklm\software \microsoft\active setup\installed components', auto start location
'hkcu\software\Microsoft\windows\currentversion\runonce\setu p', auto start location
'hklm\software\microsoft\windows\currentversion\run', auto start location
'c:\documents and settings\ user name\local settings\temp', temporary folder
'c:\ users\user name \appdata\local\temp', temporary folder

When the system call is performed in its own folder, following criteria are taken into consideration:

- (i) If an analyzed program needs some data from its own directory or file in order to run properly, it

---

**Algorithm 2** Behavior Path Score Calculation

---

**Input (D<sub>2</sub> file):** List of behaviors

**Output:** Behavior's path score

```
1: for each behavior i do
2:   if  $\mu$  == 'SF' then
3:     if process.name == D2. filename then
4:       pathScore = 'low'
5:     elif process.name! = D2. filename and D2. filename == 'system *.exe' then
6:       pathScore = 'high'
7:     else
8:       pathScore = 'moderate'
9:     end if
10:  elif  $\mu$  == 'TP' then
11:    if D2[i][path] == 'program files' then
12:      pathScore = 'moderate'
13:    elif D2[i][path] == 'AS' then
14:      pathScore = 'high'
15:    else
16:      pathScore = 'low'
17:    end if
18:  elif  $\mu$  == 'S' then
19:    if process.name == D2. filename then
20:      pathScore = 'low'
21:    elif D2. File.extension == 'exe' or DLL then
22:      pathScore = 'moderate'
23:    elif D2[i][path] == 'registry AS similar locations' then
24:      pathScore = 'high'
25:    elif D2[i][path] == 'system *.exe' then
26:      pathScore = 'very high'
27:    else
28:      pathScore = 'low'
29:    end if
30:  elif  $\mu$  == 'T' or  $\mu$  == 'AS' then
31:    pathScore = 'very high'
32:  end if
33: end for
```

---

---

**Algorithm 3** Feature Extraction from Behaviors

---

**Input (D<sub>2</sub> file):** List of behaviors

**Output (D<sub>3</sub> file):** List of Features

```
1: propertyName [] = ''
2: for each i in D2 do
3:   for j= i+1 to i+20 do
4:     if D2[i] [behaviorStatus] == D2[j] [behaviorStatus] and D2[i] [behaviorStatus] == 'A' then
5:        $\psi_i$  = 'AA'
6:     elif D2[i] [behaviorStatus] == D2[j] [behaviorStatus] and D2[i] [behaviorStatus] == 'P' then
7:        $\psi_i$  = 'PP'
8:     else
9:        $\psi_i$  = 'AP' = 'PA'
10:    end if
11:    if D2[i] [behaviorType] == D2[j] [behaviorType] and D2[i] [behaviorName] != D2[j] [behaviorName]
and (D2[i] [path] == D2[j] [path] or D2[i] [behaviorRead] before D2[j] [behaviorWrite]) then
12:      Algorithm 2 BehaviorPathScoreCalculation (D2[i] [behavior], D2[j] [behavior])
13:      propertyName[k] = D2[i] [behaviorName] + ' ' + D2[j] [behaviorName]
14:      k = k + 1
15:    end if
16:    write. D3 (propertyName[k],  $\psi_i$ , D2[i] [path], D2[j] [path], D2[i] [pathScore], D2[j] [pathScore])
17:    if D2[i] [behaviorType] != D2[j] [behaviorType] and D2[i] [behaviorRead] before D2[j] [behaviorWrite] then
18:      propertyName[k] = D2[i] [behaviorName] + ' ' + D2[j] [behaviorName]
19:      k = k + 1
20:      Algorithm 2 BehaviorPathScoreCalculation (D2[i] [behavior], D2[j] [behavior])
21:    end if
22:    write. D3 (propertyName[k],  $\psi_i$ , D2[i] [path], D2[j] [path], D2[i] [pathScore], D2[j] [pathScore])
23:  end for
24: write. D3 (D2[i] [behaviorName], D2[i] [ $\psi$ ], D2[i] [path], D2[i] [pathScore])
25: end for
```

---

Furthermore, even if the path score is low but the action state is AA or AP/PA, this property is also chosen. These properties are considered because we try to choose only malicious related patterns which differentiate malware from benign. That way, normal features which can be performed by malware and benign samples are removed from the dataset. Thus, our proposed algorithms create far fewer features than well-known algorithms and  $n$ -gram.

**Algorithm 4** Frequency Calculation and Feature Selection

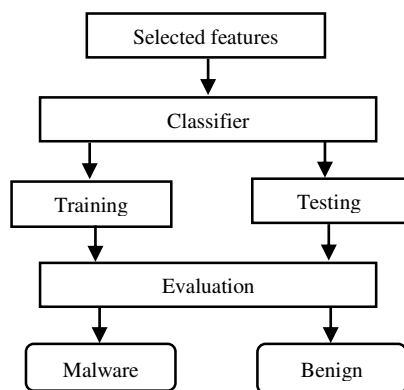
```

Input ( $D_3$  file): List of Features
Output ( $D_5$  file): Selected Features with Frequency
1: for each feature  $i$  in  $D_3$  do
2:   calculateFeatureFrequency ()
3:   write.  $D_4$  (propertyName, frequency)
4: end for
5: for each feature  $i$  in  $D_4$  do
6:   if (pathScore == 'moderate' or == 'high' or == 'very high') then
7:     write.  $D_5$  (propertyName, frequency)
8:   end if
9:   if (pathScore == 'low') and ( $\psi$  == 'AA' or == 'AP/PA') then
10:    write.  $D_5$  (propertyName, frequency)
11:   end if
12: end for

```

**C. LEARNING-BASED DETECTION**

After features are selected from the previous section, each program sample is represented by a row vector. For each property, frequency value is written. If property is repeated  $x$  times,  $x$  is written as a property value, if property is not repeated, 0 is written as a value. After the dataset is built based on feature vectors, learning algorithms are applied. In learning-based detection agent, selected features are trained by using machine learning algorithms (classifiers) including C4.5, LMT, RF, KNN, SLR and SMO. Several classifiers are used for classification to measure the proposed method efficiency. Learning-based detection agent in the cloud can be seen in Figure 2, training and testing phase can be seen in Figure 5.



**FIGURE 5.** Learning-based malware detection agent.

During the training phase, cross-validation and holdout methods are used to measure the performance. Decision trees such as C4.5, LMT and RF are used for training and testing as a classifier because they return scalable and highly accurate results in the cloud environment. Besides, decision

trees are suitable classifiers for our dataset features distribution to separate malware from benign. C4.5 uses gain ratio for feature placement. In gain ratio, the feature with the maximum gain is selected recursively for splitting criteria when features are placed on the tree. The gain ratio is prone to unbalanced partitioning and hence can create uneven trees. The information gain ratio is measured as follows.

$$Gain\ Ratio(A) = Gain(A) / SplitInformation_A(D) \quad (1)$$

$$Gain(A) = Information(D) - Information_A(D) \quad (2)$$

$$SplitInformation_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left( \frac{|D_j|}{|D|} \right) \quad (3)$$

$Gain(A)$  shows how much information will be gained when branching using the property  $A$  and  $SplitInformation_A(D)$  shows the intrinsic information which measures the entropy of the sub-dataset. C4.5 is appropriate for our dataset because it works with continuous data, eliminating data with noise, and prunes decision trees effectively.

RF is a combination of many trees which classifies using attributes that each tree is sampled independently. This classifier produces satisfying results in a dataset with low variance and many interrelated features. It uses the CART (Classification and regression tree) algorithm to generate RF trees. On the other hand, LMT is a classifier which uses a supervised learning algorithm that combines LR and decision tree learning, and produces results with high accuracy. LMT classifier creates LR functions on each node using the LogitBoost algorithm [26] and prunes the tree using the CART algorithm. The CART algorithm works according to the depth priority search and uses the Gini index as a criterion for splitting features. The Gini index is used to measure the differences between the probability distributions of target feature values. The feature with the minimum Gini index is selected as the splitting attribute. The Gini index does not work well when the number of classes and the value of properties are very large. The Gini index is calculated as follows.

$$Gini(A) = Gini(D) - Gini_A(D) \quad (4)$$

$$Gini(D) = 1 - \sum_{i=1}^m (P_i)^2 \quad (5)$$

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2) \quad (6)$$

KNN is a statistical model classifier which uses example-based learning. It is a classifier that produces good results when there is no prior knowledge about data distribution. Even though KNN classifier needs a lot of storage space during the learning phase, it performs well in the cloud environment for our dataset. Even if the SLR algorithm is not adequate to solve non-linear problems and comprise high bias which reduces the efficiency of the classifier, it is suitable and fast when combined with the proposed model. Since SMO works well for non-linear boundary situations and performs well on high-dimensional data, it performs well on our dataset on the cloud. After training and testing phases are performed by using C4.5, LMT, RF, KNN, SLR and SMO, the results are sent to the behavior-based detection agent.

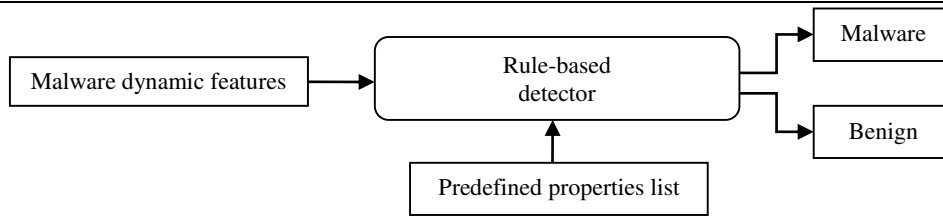


FIGURE 6. Rule-based malware detection agent.

**Algorithm 5** Rule-Based Detection

```

Input: List of Features
Output: List of Marked Program Samples
1: pathScore {'very low', 'low', 'moderate', 'high', 'very high'}
2: list {...} ← predefinedList {...}, fC ← frequencyCategory {'few', 'average', 'many', 'excessive'}
3: for each feature i do
4:   if ipathScore=='low' and i ∈ list {} and (fC == 'many' or 'excessive') then
5:     suspiciousFile = 'malware'
6:   elif ipathScore=='moderate' and i ∈ list {} and (fC == 'average' or 'many' or 'excessive') then
7:     suspiciousFile = 'malware'
8:   elif ipathScore=='high' and i ∈ list {} and (fC == 'average' or 'many' or 'excessive') then
9:     suspiciousFile = 'malware'
10:  elif ipathScore=='very hig' and i ∈ list {} and (fC == 'few' or 'average' or 'many' or 'excessive') then
11:    suspiciousFile = 'malware'
12:  elif ipathScore=='hig' and (fC == 'excessive') then
13:    suspiciousFile = 'malware'
14:    list = list + 'ifeature'
15:  elif ipathScore=='very hig' and (fC == 'many' or 'excessive') then
16:    suspiciousFile = 'malware'
17:    list = list + 'ifeature'
18:  else
19:    suspiciousFile = 'benign'
20:  end if
21: end for
  
```

**D. RULE-BASED DETECTION**

Rule-based behavior malware detection agent is running on different machines in the cloud. For detection there is no training or learning phase, instead detection is performed based on the predefined property list (Figure 6). We use malware behaviors when creating predefined properties. This list consists of features which differentiate malware from benign based on malicious behavior patterns. The malicious behavior patterns are the features which can be frequently performed by malware while rarely performed by benign samples. The malicious behavior pattern list is dynamically updated when new malware features are determined. After features are created and selected in section IV.B, the feature values are categorized based on repeated frequencies into four categories. These categories are {few}, {average}, {many}, and {excessive}. If the analyzed program features are substantially similar to the features in the list, the program is marked as malware (Algorithm 5). Otherwise, the program is marked as benign. For instance, analyzed program features are somehow in the predefined properties list but without enough repeated frequency, the analyzed program is marked as benign. Feature paths are also used during the detection. According to our findings we determined fifty features which are frequently used by malware such as CreateService, CreateRemoteThread, FindFirstFile, FindNextFile, MapviewOfFile, CreateFileMapping, QueryDirectoryWriteFile, ReadFile, WriteFile, RegDeleteValue, RegQueryValue, RegSetInfoKey, etc.

The proposed rule-based detection agent detects various forms of unknown and known malware efficiently. It is also quite fast when compared with a learning-based detection agent. After the rule-based detection agent finishes its task, the results are stored in the database and sent back to the behavior-based detection agent. Since the rule-based detection agent is quite fast when compared with a learning-based agent, the detection result is first sent to the client while learning-based detection is still performing. After the learning-based detection process is finished, its results are also sent to the client as well. For future study, we aim to combine learning-based and rule-based detection results. Behavior-based detection agent will compare the results coming from learning-based detection and rule-based detection agents. If there are some differences, the detection process will be repeated for those samples to decrease the misclassification rate. When the same classification results are gathered from both detection agents, the results will be sent to the client.

**V. CASE STUDY**

This section presents case study and experiments. In order to simulate the cloud environments, we used different computers, VMs, switches and routers in the campus network. Different versions of Windows machines are used for test cases including Windows 7, 8, 10, VMs 7, 8 and 10. Proposed dataset creation model is implemented by using Python scripting language. For learning-based detection, Weka and some python libraries are used, and for rule-based

detection the proposed algorithm is implemented in Python language as well. Totally, 7000 malware and 3000 benign portable executables are analyzed. Data collection and

representation, model performance and evaluation are explained in the following subsections.

TABLE 4. List of malware types that are analyzed (The list is shortened).

Malware MD5 Signature	Malware Type	Malware Family/Malware Specific Name
f2c6a6541976bab117d03f7a8c2ccbf7	Trojan	Trojan.Win32.Generic, Trojan.Injector, Trojan.Symmi
f3a6ab31986c928d312699dc7208a211	Ransomware	TR/Crypt.Zpack, Ransom:Win32
28cb0c8083f6a41e7b04137ab166c580	Packed Malware	Gen:Packer.PESpin, Trojan.Win32.Crypt
f448a906cc9906b8f7589b117a079280	Backdoor	Backdoor.BDS, Trojan [Backdoor]/Win32.Hlux
b64f34137982332156e058cd63cf480b	Dropper	Win32/TrojanDropper.VB, Trojan.TR/Drop.VB
996f29ba29a14fc0ebf46ce38675f8cd	Ransomware	Ransom: Win32/Blocker
4a7e35d8c11e213a051462e66e73a3e	Packed Malware	Win32.Packed.VMProtect, HackTool.GameHack!8
f3aa059c23a2080bc0b219eebf5577e0	Virus	Virus:Win32, Win32.Parite.B
9085a7dff20d6a5c287d3056d3ed1cc4	Rootkit	Dropper.Generic_r.AC, Win32:Rootkit-gen
48cd89827939b3a8976d9bb0993bc338	Spyware	Win.Spyware.Zbot, Gen:Variant.Razy
f3b03c25e1a53168a606732fc96707e3	Keylogger	Gen:Application.Keylog.dm0, TrojanSpy.Vwealer
f3c6372f9c95ba38d72a5c2219ce9f8b	Worm	Worm/Win32.Mabezat

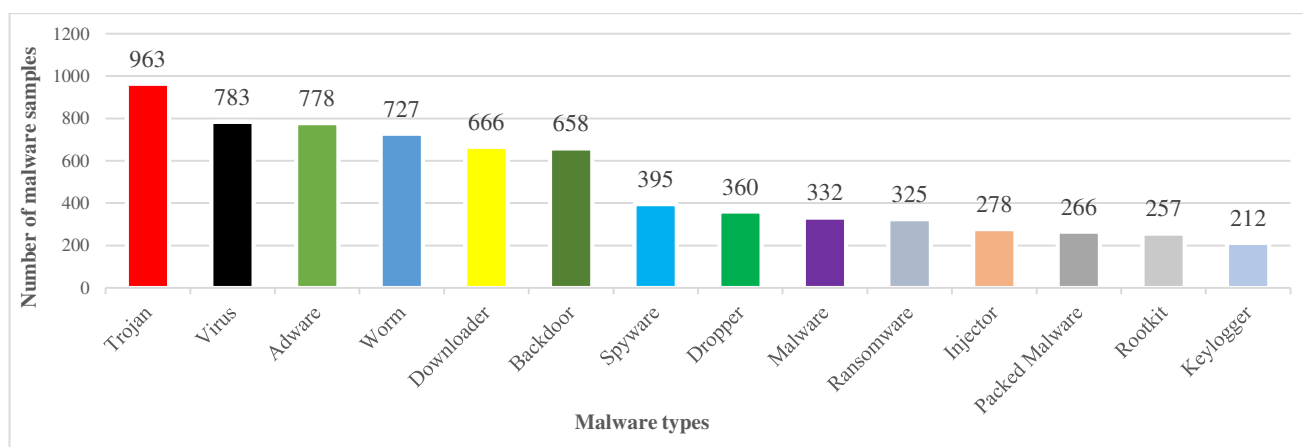


FIGURE 7. Analyzed malware distribution.

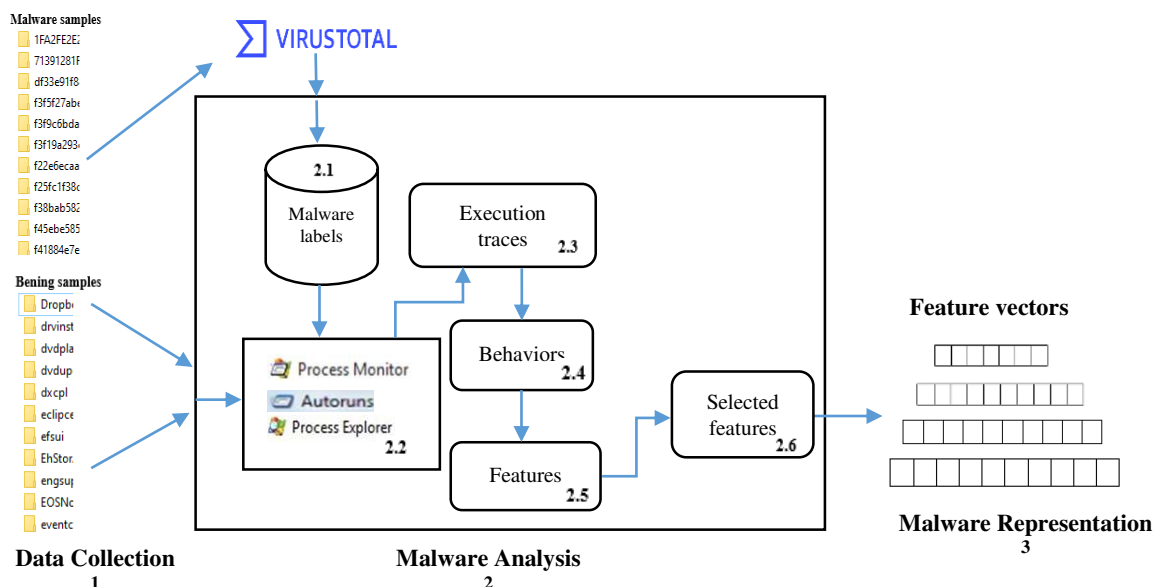


FIGURE 8. Data collection, analysis and representation process.

Process Name, Path, Operation, FileNameDirectory, DirectoryFileName, ...	Process Name, Path, Operation, FileNameDirectory, DirectoryFileName, ...
1 addinutil.exe, C:\Windows\Microsoft.NET\Framework64\v3.5\AddInUtil...	1 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Users\test\Windows7\Desk...
2 addinutil.exe, C:\Windows\Microsoft.NET\Framework64\v3.5\AddInUtil...	2 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Users\test\Windows7\Desk...
3 addinutil.exe, C:\Windows\Microsoft.NET\Framework64\v3.5\AddInUtil...	3 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Users\test\Windows7\Desk...
4 addinutil.exe, C:\Windows\System32\ntdll.dll, LoadImage, sistem, ntdll.c	4 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\ntdll.c
5 addinutil.exe, C:\Windows\Microsoft.NET\Framework64\v3.5\CreateFi...	5 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\SysWOW64\ntdll.c
6 addinutil.exe, C:\Windows\System32\mscoree.dll, CreateFile, sistem, ...	6 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows, CreateFile, siste...
7 addinutil.exe, C:\Windows\System32\mscoree.dll, QueryBasicInforma...	7 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64.c
8 addinutil.exe, C:\Windows\System32\mscoree.dll, CreateFile, sistem, ...	8 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64.c
9 addinutil.exe, C:\Windows\System32\mscoree.dll, LoadImage, sistem, n...	9 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64.c
10 addinutil.exe, C:\Windows\System32\kernel32.dll, LoadImage, sistem, ...	10 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64.c
11 addinutil.exe, C:\Windows\System32\KernelBase.dll, LoadImage, siste...	11 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64wi
12 addinutil.exe, HKLM\System\CurrentControlSet\Control\Terminal Ser...	12 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64wi
13 addinutil.exe, C:\Windows\System32\mscoree.dll, ReadFile, sistem, ms...	13 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64wi
14 addinutil.exe, HKLM\System\CurrentControlSet\Control\Nls\Sorting\...	14 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64wi
15 addinutil.exe, HKLM\System\CurrentControlSet\Control\Nls\Sorting\...	15 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64wi
16 addinutil.exe, HKLM\System\CurrentControlSet\Control\Nls\Sorting\...	16 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64cp
17 addinutil.exe, C:\Windows\Microsoft.NET\Framework64\v3.5\AddInUtil...	17 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64cp
18 addinutil.exe, C:\Windows\System32\msvcrt.dll, LoadImage, sistem, ms...	18 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64cp
19 addinutil.exe, C:\Windows\System32\sechost.dll, CreateFile, sistem, ...	19 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64c
20 addinutil.exe, C:\Windows\System32\sechost.dll, QueryBasicInforma...	20 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64c
21 addinutil.exe, C:\Windows\System32\sechost.dll, QueryBasicInforma...	21 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\wow64c
22 ARP.EXE, C:\Windows\SysWOW64\ARP.EXE, ProcessStart, self, ARP.EXE, Akt...	22 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\SysWOW64\kernel3
23 ARP.EXE, C:\Windows\SysWOW64\ARP.EXE, ThreadCreate, self, ARP.EXE, Akt...	23 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\System32\user32...
24 ARP.EXE, C:\Windows\System32\ntdll.dll, LoadImage, sistem, ntdll.dll, ...	24 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows, CreateFile, siste...
25 ARP.EXE, C:\Windows\System32\ntdll.dll, LoadImage, sistem, ntdll.dll, ...	25 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Users\test\Windows7\Desk...
26 ARP.EXE, C:\Windows>CreateFile, sistem, Windows, Passive, File System, ...	26 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Users\test\Windows7\Desk...
27 ARP.EXE, C:\Windows\System32\wow64.dll, CreateFile, sistem, wow64.dll	27 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\SysWOW64\kernel3
28 ARP.EXE, C:\Windows\System32\wow64.dll, QueryBasicInformationFile, s...	28 f25a13b7c356fa70ae1947d7f4ef0610.exe, C:\Windows\SysWOW64\KernelE
29 ARP.EXE, C:\Windows\System32\wow64.dll, CreateFile, sistem, wow64.dll	29 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\System\CurrentControlS
30 ARP.EXE, C:\Windows\System32\wow64.dll, LoadImage, sistem, wow64.dll, ...	30 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\System\CurrentControlS
31 ARP.EXE, C:\Windows\System32\wow64win.dll, CreateFile, sistem, wow64w...	31 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\Software\Wow6432Node\F
32 ARP.EXE, C:\Windows\System32\wow64cpu.dll, CreateFile, sistem, wow64c...	32 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\SOFTWARE\Policies\Micr
33 ARP.EXE, C:\Windows\System32\wow64cpu.dll, QueryBasicInformationFil...	33 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\SOFTWARE\Policies\Micr
34 ARP.EXE, C:\Windows\System32\wow64cpu.dll, CreateFile, sistem, wow64c...	34 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\SOFTWARE\Policies\Micr
35 ARP.EXE, C:\Windows\System32\wow64log.dll, CreateFile, sistem, wow64l...	35 f25a13b7c356fa70ae1947d7f4ef0610.exe, HKLM\SOFTWARE\Policies\Micr
36 ARP.EXE, C:\Windows\System32\kernel32.dll, LoadImage, sistem, kernel3...	
37 ARP.EXE, C:\Windows\System32\kernel32.dll, LoadImage, sistem, kernel3...	

FIGURE 9. Extracted behaviors for 4 samples (Figure 8. 2.4) (The list is shortened).

result, Benign	result, Malware
1 name, addinutil.exe	1 name, f3f46b1f53a62b77236950cbdad3f3a0.exe
2 ProcessStartSF , 1	2 ThreadCreateSF , 1
3 ProcessStartThreadCreateSF , 1	3 RegOpenKeySY , 272
4 ProcessStartLoadImageSF , 1	4 RegOpenKeyRegQueryKeySY , 220
5 ThreadCreateSF , 1	5 RegQueryKeySY , 703
6 ThreadCreateLoadImageSF , 1	6 RegQueryValueSY , 575
7 LoadImageSY , 198	7 QueryNameInformationFileSY , 121
8 CreateFileSY , 2251	8 QueryNameInformationFileSF , 5
9 CreateFileQueryBasicInformationFileSY , 1488	9 ProcessStartTP , 6
10 CreateFileLoadImageSY , 276	10 ProcessStartThreadCreateTP , 6
11 CreateFileReadFileSY , 211	11 ThreadCreateTP , 16
12 QueryBasicInformationFileSY , 789	12 LoadImageSF , 5
13 QueryBasicInformationFileCreateFileSY , 1366	13 LoadImageSY , 202
14 QueryBasicInformationFileLoadImageSY , 132	14 CreateFileSY , 328
15 QueryBasicInformationFileReadFileSY , 62	15 CreateFileQueryBasicInformationFileSY , 159
16 LoadImageReadFileSY , 21	16 CreateFileLoadImageSY , 174
17 RegQueryValueSY , 46	17 QueryBasicInformationFileSY , 120
18 ReadFileSY , 303	18 QueryBasicInformationFileCreateFileSY , 168
19 RegOpenKeySY , 18	19 QueryBasicInformationFileLoadImageSY , 86
20 RegOpenKeyRegQueryKeySY , 6	20 CreateFileQueryNameInformationFileSY , 5
21 RegQueryKeySY , 1182	21 CreateFileTP , 319
22 RegOpenKeyTP , 1518	22 RegSetInfoKeySY , 169
23 RegQueryKeyTP , 721	23 RegOpenKeyTP , 1110
24 RegQueryKeyRegEnumKeyTP , 24	24 RegSetInfoKeyTP , 523
25 RegEnumKeyTP , 120	25 CreateFileReadFileSY , 57
26 RegEnumKeyRegQueryKeyTP , 24	26 CreateFileCreateFileMappingSY , 129
27 RegOpenKeyRegQueryKeyTP , 60	27 QueryBasicInformationFileReadFileSY , 13
28 RegOpenKeyRegEnumValueTP , 12	28 QueryBasicInformationFileCreateFileMappings
29 RegQueryKeyRegEnumValueTP , 6	29 ReadFileSY .62
30 RegEnumValueTP , 19	
31 RegQueryValueTP , 2076	
32 QueryDirectorySY , 66	
33 QueryDirectoryCreateFileSY , 42	
1 result, Benign	1 result, Malware
2 name, ARP.EXE	2 name, f25a13b7c356fa70ae1947d7f4ef0610.exe
3 ProcessStartSF 1, 1	3 ProcessStartSF , 1
4 ProcessStartThreadCreateSF 1, 1	4 ProcessStartThreadCreateSF , 1
5 ProcessStartLoadImageSF 1, 1	5 ThreadCreateSF , 1
6 ThreadCreateSF 1, 1	6 ThreadCreateLoadImageSF , 1
7 ThreadCreateLoadImageSF 1, 1	7 LoadImageSF , 2
8 LoadImageSF 1, 1	8 LoadImageSY , 84
9 LoadImageSY 1, 27	9 CreateFileSY , 76
10 CreateFileSY 1, 26	10 CreateFileQueryBasicInformationFileSY , 38
11 CreateFileQueryBasicInformationFileSY 1, 14	11 CreateFileLoadImageSY , 58
	12 QueryBasicInformationFileSY , 36
	13 QueryBasicInformationFileCreateFileSY , 45
	14 QueryBasicInformationFileLoadImageSY .29
	15

FIGURE 10. Generated features for 4 samples (Figure 8. 2.5) (The list is shortened).



feature set is increasing over time. For 10,000 samples we gathered 751 features. On the other hand, for rule-based detection feature vectors, only features that are performed by related samples are taken into consideration. The repetition of each feature is written. For 10,000 samples, we got an average of 60 features, which is quite fewer when compared to the other datasets.

**B. MODEL PERFORMANCE AND EVALUATION**

After the feature selection process is completed, learning-based and rule-based detection processes start. To measure the performance of proposed model holdout and cross-validation methods, as well as detection rate (*DR*), false positive rate (*FPR*), f-measure, and accuracy metrics were used. At first, when the generated dataset was small, the cross-validation method returned more feasible results than holdout. However, when the dataset has grown timely, the holdout method returned favorable results as well. Best performances are obtained when *k* is chosen 10 for cross validation, and the data is divided into 80% training and 20% test for holdout method. TP represents the number of malware samples correctly classified as malware, TN the number of benign samples correctly classified as benign, FP the number of benign samples being mistakenly classified as malware, and FN the number of malware samples being mistakenly classified as benign. *DR*, *FPR*, f-measure, and accuracy metrics are calculated by using confusion matrix (Table 5) as follows:

$$DR = Recall = TP / (TP + FN) \tag{7}$$

$$FPR = FP / (FP + TN) \tag{8}$$

$$Precision = TP / (TP + FP) \tag{9}$$

$$F\text{-Measure} = (2 * precision * recall) / (precision + recall) \tag{10}$$

$$Accuracy = TP + TN / (TP + TN + FP + FN) \tag{11}$$

TABLE 5. Confusion matrix.

		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

**VI. RESULTS AND DISCUSSION**

This section summarizes the test results and discusses the proposed system performance. When performance is evaluated, various learning algorithms are used. During the training and testing, best performances are gathered by using cross validation *k* = 10 and holdout method which is using 80% training and 20% testing sets. *DR*, *FPR*, f-measure, and accuracies are used as metrics to compare test results. For rule-based detection, we did not use any training and testing phases. The performance is obtained in real time by using predefined features. The cloud environment has provided a fast and scalable environment for our learning- and rule-based system. The test results can be seen in Figure 13, Table 6, Table 7, Table 8, Table 9, and Table 10.

Table 6 demonstrates the proposed model performance when 10,000 program samples are analyzed. For learning-based detection, both cross validation and holdout methods are performed fairly well. Cross validation results are slightly higher than holdout results. The best results are obtained when ML classifiers such as decision trees (J48, RF, LMT)

and KNN are used. For instance, in J48 *DR*, *FPR*, f-measure and accuracy are measured as 99.8%, 0.4%, 99.8% and 99.75%, respectively. In the same way, RF algorithm achieved 100% for *DR*, 0.6% for *FPR*, 99.6% for f-measure, and 99.83% for accuracy; LMT achieved 99.3% for *DR*, 0.5% for *FPR*, 99.6% for f-measure, and 99.38% for accuracy; and KNN achieved 100% for *DR*, 1.2% for *FPR*, 99.7% for f-measure, and 99.64% for accuracy. The obtained results on SLR and SMO classifiers are slightly lower than J48, RF, LMT and KNN. Using appropriate kernels for SMO and reducing bias for SLR can increase the performance. On the other hand, the obtained test results are satisfactory for rule-based detection which do not use any learning algorithm and do not require any training phase. rule-based detection achieved 97.8%, 6.6%, 97.4%, and 96.5% for *DR*, *FPR*, f-measure and accuracy, respectively.

Table 7 shows performance on *n*-gram, ClaMP and our dataset based on selected classifiers. It can be clearly seen that J48, RF, and KNN classifiers perform better on our dataset. For example, the J48 algorithm performance on *n*-gram is measured as 98.1% for *DR*, 2% for *FPR*, and 98.05% for accuracy; on ClaMP dataset [37] it is measured as 98.2% for *DR*, 2.6% for *FPR*, and 97.8% for accuracy; and on our dataset it is measured as 99.8% for *DR*, 0.4% for *FPR*, and 99.75% for accuracy. Similar results are obtained by using different ML classifiers as well.

The number of features is reasonable on our dataset when it is compared with the *n*-gram dataset (Table 7). Figure 13 and Table 8 indicate the number of properties in the feature vectors for learning-based detection and rule-based detection, respectively. Until a certain number, the numbers of properties are increased while analyzed program samples are increasing (Figure 13). This is because some program samples exhibit different features. Furthermore, when properties are combined, each sample’s features are added to the feature vector. For repeated frequency, frequency number is written for feature value. If property is not presented in the feature vector, 0 is written for that property. This conversion also raises the number of features in our dataset. However, for rule-based detection we have not created a feature vector. Thus, the number of features is reasonable when compared to other feature extraction methods (Table 8). For rule-based detection, our dataset consists of an average of 60 features with a minimum of 5 and a maximum of 150 (Table 8). Most of the time, the extracted malware features are more than the benign features for each sample for our dataset.

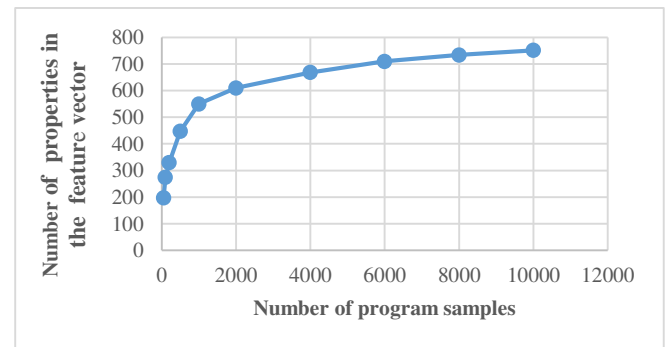


FIGURE 13. Number of analyzed program samples versus properties in the feature vectors for learning-based detection.



TABLE 6. Proposed model results on various classifiers and rule-based detection.

Method	Classifier	DR (%)	FPR (%)	F-measure (%)	Accuracy
Learning-based detection using Cross-validation	J48	99.8	0.4	99.8	99.75
	RF	100	0.6	99.6	99.83
	LMT	99.3	0.5	99.6	99.38
	KNN	100	1.2	99.7	99.64
	SLR	95.4	1.6	97.3	96.29
	SMO	93.4	6.7	95.2	93.37
Learning-based detection using Holdout	J48	99.3	0.3	99.6	99.45
	RF	99.9	1	99.7	99.6
	LMT	99.4	0	99.7	99.6
	KNN	100	1.5	99.7	99.55
	SLR	93.4	3.6	95.8	94.35
	SMO	92.5	7.6	94.4	92.45
Rule-based detection	No training	97.8	6.6	97.4	96.5

TABLE 7. Performance comparison of *n*-gram, ClaMP and our dataset based on selected classifiers.

Classifier	<i>n</i> -gram Dataset(1386 samples, 14,520 features)			ClaMP Dataset (5210 samples, 69 features)			Our Dataset (10,000 samples, 751 features)		
	DR	FPR	Accuracy	DR	FPR	Accuracy	DR	FPR	Accuracy
J48	98.1	2	98.05	98.2	2.6	97.8	99.8	0.4	99.75
RF	99	0.4	99.2	99.7	2.1	98.8	100	0.6	99.83
KNN	94.9	0	97.47	98	2	97.9	100	1.2	99.64

TABLE 8. Number of features for rule-based detection.

	Min	Max	Average
Each sample	5	150	60

To evaluate the efficiency of the proposed system more accurately, model performance for different feature extraction methods is compared in Table 9. In addition, *DR*, *FPR* and accuracies are compared on the same classifiers for various cloud-based and other studies in the literature (Table 10). The proposed feature extraction method has generated considerably better results than other methods (Table 9) including studies in [38], [39], [40], [41], [42]. The 99.8% performance is measured for the proposed feature extraction method versus 96.4%, 99.6%, 97.6%, 89.92%, and 99.28% in [38], [39], [40], [41], [42] studies, respectively. Even though the performance of some feature extraction methods is quite good in the literature, it cannot be confirmed that they are as successful as the proposed method due to the other deficiencies such as the low number of analyzed samples and higher number of extracted features. The performance of various ML classifiers such as J48, RF, KNN, SVM, SLR, and neural networks on different studies are measured (Table 10). It can be clearly seen that combining the proposed feature extraction method with an appropriate ML algorithm produces more satisfactory results in terms of *DR*, *FPR* and accuracies when compared with other studies in the literature.

In this section, the performance of the proposed cloud-based malware detection system and the leading methods in the literature are compared. The proposed system successfully detects both different types and families of malware, as well as the new generation and previously unknown malware. The measured performance values increase as the number of programs analyzed increase. In addition, the number of features does not increase after a certain number of programs being analyzed in the proposed system. The results obtained were higher than the pioneering

method results in the literature. In addition, testing the proposed system in the cloud environment and using 2 different detection mechanisms provide a distinct advantage. On the other hand, current studies in the literature face some insufficiency for malware detection:

1. Behaviors are not clearly determined.
2. The number of extracted features is high.
3. Perform well for only certain types and families of malicious software.
4. Inadequate for detection of new generation malware.
5. Not resistant to evasion and stealth techniques which leads to decreasing performance.

In the proposed system, these deficiencies were identified and necessary contributions were made to increase the performance.

In addition, a number of key findings were obtained during the malware analysis. These findings should be taken into account while creating a fast and effective detection method. The main findings identified can be listed as follows:

1. Malware creates random files with meaningless file names.
2. Several malware types use newly created processes and existing processes for malicious purposes.
3. Some of the malware injects itself into operating system exe files including svchost.exe, conhost.exe, winlogon.exe, etc. and system DLLs on Windows operating systems.
4. Some malware variants hide themselves by creating similar systems' and third-parties' file names.
5. Some malware variants disable the existing security software such as firewall, IDS, antivirus scanner, etc. whenever they are performed.
6. Some malware variants perform malicious activities in the temporary files.

7. Malware becomes persistent in the computer system by locating itself in the automatic startup file and registry locations.

TABLE 9. Comparison of different malware dataset creation methods.

Paper	Feature Extraction Method	Performance (%)	Year
Anderson et al. [38]	Markov chain weighted directed graph	96.40	2011
Chandramohan et al. [39]	Bounded feature space behavior modeling	99.60	2013
Das et al. [40]	Semantics of malicious behaviors	97.60	2016
Narayanan et al. [41]	Context sensitive, adaptable and scalable rules	89.92	2017
Jeon et al. [42]	Dynamic analysis with CNN	99.28	2020
<b>Proposed Method</b>	<b>Cloud-based behavior centric model</b>	<b>99.80</b>	<b>2021</b>

TABLE 10. Performance of ML classifiers on different studies.

Study	Classifier	DR (%)	FPR (%)	Accuracy (%)	Year
Santos et al. [43]	DT: J48	92	9	91.25	2013
	KNN K = 1	93	7	92.8	
	KNN K = 3	91	8	91.7	
	SVM: Polynomial	88	9	89.65	
Yousefi-Azar et al. [44]	SVM	-	5.07	93.44	2018
	RF	-	6.82	90.05	
	KNN	-	10	91.28	
Yadav [22]	Neural network	86.4	-	-	2019
Kumar et al. [45]	Decision tree	-	12.5	95.7	2020
	RF	-	6.7	97.9	
	LR	-	4.2	94.3	
	NB	-	-	32.52	
Azeez et al. [46]	Decision tree	95	5.36	98.29	2021
	RF	98	2.13	99.24	
	J48	<b>99.8</b>	<b>0.4</b>	<b>99.75</b>	
<b>Proposed Method</b>	RF	<b>100</b>	<b>0.6</b>	<b>99.83</b>	2021
	KNN	<b>100</b>	<b>1.2</b>	<b>99.64</b>	
	SVM	<b>93.4</b>	<b>6.7</b>	<b>93.37</b>	
	SMO	<b>93.4</b>	<b>6.7</b>	<b>93.37</b>	

## VII. LIMITATIONS AND FUTURE WORKS

Although CBCM is quite effective in detecting different kinds of malware, there are some limitations that need to be addressed. Malware samples were selected randomly among several malicious software variants, but malware types were not equally distributed. For example, most of the malware samples analyzed were Trojan, virus, adware, worm and downloader. More ransomware, spyware, rootkit, and packed malware need to be analyzed. In total, 10,000 program samples were analyzed, in the future the number of program samples will be increased. Even though our feature extraction and selection algorithms work quite well, there are still some benign samples misclassified as malware. This is because there are some features which are frequently seen in malware but rarely seen in benign. The numbers of these features increased when more program samples are analyzed. In the future, we will improve our feature selection algorithm and also use well-known algorithms mentioned in the literature to decrease the features that lead to misclassification. Even though the proposed model can detect some portion of the obfuscated and packed malware samples, it cannot detect all of them. Thus, the proposed model will be improved more in order to detect those malware samples. We analyze malware only on various versions of Windows machines. We will extend our system to other operating systems including different Linux distributions and macOS.

In this study, we classify the analyzed program samples into two categories including malware and benign classes. In the future, we will also classify the malware types into different

classes such as virus, worm, Trojan, rootkit, ransomware, etc. In the cloud environment, a limited number of servers and VMs are used during the analysis, these numbers can be increased in the feature. Rule-based detection agent works in real-time, but learning-based detection agent does not work in real time. In the future, we are planning to combine learning-based and rule-based detection agents to work together in real time. In addition, we aim to build a deep learning-based detection agent on different servers on the cloud as well.

## VIII. CONCLUSION

In this paper, a malware detection system which works in the cloud computing environment is presented. There are two parts including client and cloud environment. A client sends suspicious file samples to the cloud, and receives the analyzed results which show whether the suspicious samples are malware or benign. In the cloud, our system consists of three phases. In the first phase, file samples are analyzed by using relevant tools to gather execution traces on different VMs and sent to the behavior-detection agent. In behavior-based detection agent, behaviors and features are generated by using proposed CBCM. In this phase, system calls, system call types, system call paths, system resources and different file types are considered. By this way, malicious features patterns are segregated from benign ones. In the third phase, selected features are sent to the learning-based and rule-based agents to classify file samples as malware or benign. The

results are sent back to the behavior-based detection agent, evaluated and sent back to the client.

Our test results confirm that combining proposed feature extraction and selection phases with appropriate learning- and rule-based detection agents increase the performance. The proposed system can effectively detect both known and unknown malware for different data samples. When the proposed system is compared to other systems in the literature, the obtained *DR* and accuracies are quite higher while *FPR* and *FNR* are lower. On the other hand, some portion of the malware samples are remaining undetected due to the use of advanced code obfuscation techniques. Increasing the analysis time, as well as determining more specific features may increase the *DR*. We also aim to extend our system to work on different cloud provider premises such as AWS, IBM Cloud Foundry, Salesforce Platform as well. We hope that our proposed system and its algorithms will assist those who would like to develop an influential detection system on the cloud for daily evolving malware.

## REFERENCES

- [1] O. Aslan and R. Samet, "Investigation of possibilities to detect malware using existing tools," in Proc. IEEE/ACS 14th Int. Conf. Comput. Syst. Appl. (AICCSA), Oct. 2017.
- [2] O. Aslan and R. Samet, "A comprehensive review on malware detection approaches," IEEE Access, vol. 8, pp. 6249-6271, Jan. 2020.
- [3] M.R. Watson, A. K. Marnerides, A. Mauthe, and D. Hutchison, "Malware detection in cloud computing infrastructures", IEEE Transactions on Dependable and Secure Computing, vol. 13 no. 2, pp. 192-205, 2015.
- [4] Website explaining cloud computing. Accessed: Jan. 12, 2021. [Online]. Available: <https://www.investopedia.com/terms/c/cloud-computing.asp>
- [5] Website explaining cloud computing. Accessed: Jan. 14, 2021. [Online]. Available: <https://azure.microsoft.com/en-us/overview/what-is-cloud-computing/#benefits>
- [6] O. Aslan, M. Ozkan-Okay, and D. Gupta, "A Review of Cloud-Based Malware Detection System: Opportunities, Advances and Challenges", European Journal of Engineering and Technology Research, vol. 6 no. 3, pp. 1-8, March. 2021.
- [7] D. Gupta, S. Bhatt, M. Gupta, and A. S. Tosun, "Future smart connected communities to fight covid-19 outbreak," Internet of Things, vol. 13, 100342, March. 2021.
- [8] D. Gupta, S. Bhatt, M. Gupta, O. Kayode, and A. S. Tosun, "Access control model for google cloud iot," in Proc. IEEE 6th Intl Conf. on Big Data Security, pp. 198-208, May. 2020.
- [9] L. Martignoni, R. Paleari, and D. Bruschi, "A framework for behavior-based malware analysis in the cloud," in Proc. Int. Conf. Inf. Syst. Secur. Berlin, Germany: Springer, 2009.
- [10] S. K. Cha, I. Moraru, J. Jang, J. Truelove, D. Brumley, and D. G. Andersen, "SplitScreen: Enabling efficient, distributed malware detection," J. Commun. Netw., vol. 13, no. 2, pp. 187-200, Apr. 2011.
- [11] T. Y. Win, H. Tianfield, and Q. Mair, "Detection of malware and kernel-level rootkits in cloud computing environments," in IEEE 2nd Int. Conf. on Cyber Security and Cloud Computing, pp. 295-300, Nov. 2015.
- [12] M. K. Gupta, S. Shaw, and S. Chakraborty, "Pattern Based Malware Detection Technique in Cloud Architecture," Sep. 2017.
- [13] N. Rakotoniravony, B. Taubmann, W. Mandarawi, E. Weishäupl, P. Xu, B. Kolosnjaji, and H.P. Reiser, "Classifying malware attacks in IaaS cloud environments," Journal of Cloud Computing, vol. 6 no. 1, pp. 1-12, Dec. 2017.
- [14] H. Sun, X. Wang, R. Buyya, and J. Su, "CloudEyes: Cloud-based malware detection with reversible sketch for resource-constrained Internet of Things (IoT) devices," Softw. Pract. Exper., vol. 47, no. 3, pp. 421-441, Mar. 2017.
- [15] N. M. Babu, and G. Murali, "Malware detection for multi cloud servers using intermediate monitoring server," in Int. Conf. on Energy, Communication, Data Analytics and Soft Comput. (ICECDS), pp. 3609-3612, Agu2017.
- [16] L. Xiao, Y. Li, X. Huang, and X. Du, "Cloud-based malware detection game for mobile devices with offloading," IEEE Trans. Mobile Comput, vol. 16, no. 10, pp. 2742-2750, Oct. 2017.
- [17] M. Abdelsalam, R. Krishnan, Y. Huang and R. Sandhu, "Malware detection in cloud infrastructures using convolutional neural networks," in IEEE 11th Int. Conf. on Cloud Computing (CLOUD), pp. 162-169, July.2018.
- [18] Q. K. Ali Mirza, I. Awan, and M. Younas, "CloudIntell: An intelligent malware detection system," Future Gener. Comput. Syst., vol. 86, pp. 1042-1053, Sep. 2018.
- [19] Q. K. Ali Mirza, I. Awan, and M. Younas, "A Cloud-Based Energy Efficient Hosting Model for Malware Detection Framework," in IEEE Global Communications Conference (GLOBECOM), pp. 1-6, Dec. 2018.
- [20] S. Shen, L. Huang, H. Zhou, S. Yu, F. Fan, and Q. Cao, "Multistage signaling game-based optimal detection strategies for suppressing malware diffusion in fog-cloud-based IoT networks," IEEE Internet of Things Journal, vol. 5 no. 2, pp. 1043-1054, Apr. 2018.
- [21] W. Zhou and B. Yu, "A cloud-assisted malware detection and suppression framework for wireless multimedia system in IoT based on dynamic differential game," China Communications, vol. 15 no. 2, pp. 209-223, Feb. 2018.
- [22] R. M. Yadav, "Effective analysis of malware detection in cloud computing," Comput. Secur., vol. 83, pp. 14-21, Jun. 2019.
- [23] P. Indirapriyadarsini, M. U. Mohiuddin, M. Taqeeuddin, C. S. Reddy, and T. Koushik, "Malware Detection using Machine Learning and Cloud Computing," international Journal for Research in Applied Science & Engineering Technology (IJRASET), vol. 8, pp. 101-104, Jun. 2020.
- [24] D. Deyannis, E. Papadogiannaki, G. Kalivianakis, G. Vasiliadis, and S. Ioannidis, "Trustav: Practical and privacy preserving malware analysis in the cloud," in Proc. of the Tenth ACM Conference on Data and Application Security and Privacy, pp. 39-48, March, 2020.
- [25] O. Aslan, R. Samet, and O. O. Tanrıöver, "Using a Subtractive Center Behavioral Model to Detect Malware," Security and Communication Networks, vol. 2020, Feb. 2020.
- [26] J. Friedman, T. Hastie, and R. Tibshirani, "Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors)," Annals of statistics, vol. 28 no. 2, pp. 337-407, 2000.
- [27] Das malwerk, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <https://dasmalwerk.eu/>
- [28] MalwareBazaar, malware downloading website. Accessed: Jan. 15 2021. [Online]. Available: <https://bazaar.abuse.ch/>
- [29] TheZoo aka, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <https://thezoo.morirt.com/>
- [30] Malwarebenchmark, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <http://malwarebenchmark.org/>
- [31] Malshare, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <https://malshare.com/>
- [32] Tekdefense, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <http://www.tekdefense.com/downloads/>
- [33] Virussign, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <https://virussign.com/>
- [34] Virusshare, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <https://virusshare.com/>
- [35] Kernelmode, malware downloading website. Accessed: Jan. 15, 2021. [Online]. Available: <https://www.kernelmode.info/forum/>
- [36] VirusTotal, online malware detection scanner. Accessed: Jan. 19, 2021. [Online]. Available: <https://www.virustotal.com/>
- [37] Classification of Malware with PE headers (ClAMP), malware dataset. Accessed: Feb. 23, 2021. [Online]. Available: <https://github.com/urwithajit9/ClAMP>
- [38] B. Anderson, D. Quist, J. Neil, C. Storlie, and T. Lane, "Graph-based malware detection using dynamic analysis," J. Comput. Virol., vol. 7, no. 4, pp. 247-258, Nov. 2011.
- [39] M. Chandramohan, H. B. K. Tan, L. C. Briand, L. K. Shar, and B. M. Padmanabhuni, "A scalable approach for malware detection through bounded feature space behavior modeling," in Proc. 28th IEEE/ACM Int. Conf. Autom. Softw. Eng. (ASE), pp. 312-322, Nov. 2013.
- [40] S. Das, Y. Liu, W. Zhang, and M. Chandramohan, "Semantics-based online malware detection: Towards efficient real-time protection against malware," IEEE Trans. Inf. Forensics Security, vol. 11, no. 2, pp. 289-302, Feb. 2016.

- [41] A. Narayanan, M. Chandramohan, L. Chen, and Y. Liu, "Context-aware, adaptive, and scalable Android malware detection through online learning," IEEE Trans. Emerg. Topics Comput., vol. 1, no. 3, pp. 157-175, Jun. 2017.
- [42] J. Jeon, J. H. Park, & Y. S. Jeong, "Dynamic analysis for IoT malware detection with convolution neural network model," IEEE Access, vol. 8, pp. 96899-96911, Jun. 2020.
- [43] I. Santos, F. Brezo, X. Ugarte-Pedrero, & P. G. Bringas, "Opcode sequences as representation of executables for data-mining-based unknown malware detection," Information Sciences, vol. 231, pp. 64-82, May. 2013.
- [44] M. Yousefi-Azar, L. G. Hamey, V. Varadharajan, and S. Chen, "Malytics: a malware detection scheme," IEEE Access, vol. 6, pp. 49418-49431, 2018.
- [45] R. Kumar, K. Sethi, N. Prajapati, R. R. Rout, and P. Bera, "Machine Learning based Malware Detection in Cloud Environment using Clustering Approach," In 2020 11th International Conf. on Computing, Communication and Networking Technologies (ICCCNT) pp. 1-7, July. 2020.
- [46] N. A. Azeez, O. E. Odufuwa, S. Misra, J. Oluranti, and R. Damaševičius, "Windows PE Malware Detection Using Ensemble Learning," in Informatics, Multidisciplinary Digital Publishing Institute, vol. 8, no. 1, p.p. 1-10, Feb. 2021.



**ÖMER ASLAN** is a Dr. researcher in the computer engineering department at the University of Siirt, Turkey. He received his PhD in cyber security field in 2020 from University of Ankara, Turkey, MSc in information security field in 2014 from University of Texas at San Antonio, United States of America (USA), and BSc in computer engineering department in 2009 at University of Trakya, Turkey. He is working on computer systems, information security, cyber security, malware analysis, cloud computing,

and IoT device security. He has published several papers on international journals and conferences. He has been also serving as a reviewer in some prestigious journals.



**Merve Ozkan-Okay** received B.S. and M.S. degree in computer engineering from Ankara University, in 2014 and 2016 respectively. She is a Research Assistant and doing Ph.D. in Department of Computer Engineering, Ankara University. Her current research interests include cyber security, cloud-based systems, machine learning and image processing. She has published several papers on international journals and conferences.



**Deepti Gupta** is pursuing Ph.D. in computer science at University of Texas at San Antonio (UTSA). She received B.S. and M.S. degree in mathematics from Chaudhary Charan Singh University, India, the M.Tech. degree in computer engineering from Shobhit University, India and the M.S. degree in computer science from The University of Texas at San Antonio (UTSA). She has worked as an Adjunct Faculty in the Department of Computer Science at St. Edward

University, Austin. Her primary area of research includes security and privacy in cloud computing and Internet of Things, security models, and deep learning. She has also served as reviewer and committee member in conferences.