



ANDROID MALWARE BRECH DETECTION FOR IOT DEVICES USING EIGEN SPACE DEEP LEARNING

- 1.Mrs.B.ArunaSri ,Assistant Professor,CSE,MRECW,Hyderabad
- 2.K.Bhavya Goud, Student,CSE,MRECW,,Hyderabad
- 3.P.Jasmi, Student,CSE,MRECW,,Hyderabad
- 4.M.Manasa, Student,CSE,MRECW,Hyderabad

ABSTRACT

Internet of Things (IoT) in military settings generally consists of a diverse range of Internet-connected devices and nodes (e.g. medical devices and wearable combat uniforms). These IoT devices and nodes are a valuable target for cyber criminals, particularly state-sponsored or nation state actors. A common attack vector is the use of malware. In this paper, we present a deep learning based method to detect Internet Of Battlefield Things (IoBT) malware via the device's Operational Code (OpCode) sequence. We transmute OpCodes into a vector space and apply a deep Eigenspace learning approach to classify malicious and benign applications. We also demonstrate the robustness of our proposed approach in malware detection and its sustainability against junk code insertion attacks. Lastly, we make available our malware sample on Github, which hopefully will benefit future research efforts (e.g. to facilitate evaluation of future malware detection approaches)..

INTRODUCTION

Junk Software Injection Attack is a software anti-forensic tactic used against OpCode inspection. As the name suggests, the introduction of junk code that involve the incorporation of innocuous OpCode sequences that do not run in malware, or the inclusion of instructions (e.g. NOP) that do not necessarily make any difference in malware operations. Junk Code Injection Technique is typically intended to obscure the malicious OpCode sequence and that the 'proportion' of malicious OpCodes in malware in our suggested solution, we use affinity-based requirements to minimize junk OpCode injection antiforensics. Specifically, our feature collection approach excludes less detailed OpCodes to minimize the impact of insertion of OpCodes garbage. To show the efficacy of our proposed solution to Code Insertion Attack, In an incremental manner, a specified proportion of all the elements in the graph generated by each sample was chosen randomly and their value increased by one. For example, in the

4th iteration of the evaluations, 20% of the indices in each sample graph were chosen to increase their value by one. In addition, the probability of repeated feature collection for simulate has been included in our assessments and several injections of OpCode. Incrementing $E_{i,j}$ in the sample generated graph is equal to injecting $OpCode_j$ next to $OpCode_i$ in the sample instruction series to deceive the detection algorithm. Algorithm 2 describes the iteration of the junk code insertion during experiments, and this procedure should be repeated for each iteration of the k-fold validation. In order to demonstrate the robustness of our proposed solution and to compare it with existing proposals, two congruent algorithms mentioned in Section 1 are applied to our developed dataset using Adaboost as a classification algorithm.

OBJECTIVE OF THE PROJECT

Robust malware detection for internet of things is a process performed by software and hardware. Input Architecture is the method of translating a user-oriented data definition into a computer-based

program. This architecture is necessary in order to prevent mistakes in the data input process and to display the correct way to the management to get the correct information from the computerized system. This is done by designing userfriendly data entry screens to accommodate huge data volumes. The aim of input design is to make data entry simpler and error-free. The data entry system is designed in such a way that all data processing can be done. It also offers a record screening service. When the data is entered, it must test the authenticity of the results. Data can be entered with the aid of a phone. Reasonable alerts are received as appropriate so that the consumer is not immediately in maize. The goal of the interface design is therefore to create an interface structure that is simple to navigate.

LITERATURE SURVEY:

EXISTING SYSTEM:

Malware identification approaches may be either static or dynamic. In contextual malware detection methods, the program is executed in a managed environment (e.g. a virtual machine or sandbox) to capture the functional characteristics, such as the necessary resources, the direction of execution, and the desired privilege, in order to identify the program as malware or benign. Static methods (e.g. signature-based detection, byte-based detection, OpCode sequence identification and control flow graph traversal) Statistically check the software code for questionable programs. David et al have proposed Deepsign to automatically detect malware using a signature generation process. The above generates a dataset based on API call activity records, registry entries, site queries, port accesses, etc., in a sandbox and transforms records to a binary matrix. They used the deep-

seated network for classification and allegedly achieved 98.6 percent accuracy. In another study, Pascanu et al. suggested a method for modeling malware execution using natural language processing. They extracted the relevant features using a recurrent neural network to predict future API calls. Both logistic regression and multilayer perceptrons were then used as a classification module. Next API call estimation and use the history of previous events as functionality. It has been recorded that a true positive rate of 98.3 percent and a false positive rate of 0.1 percent is obtained. Demme et al. investigated the feasibility of developing a malware detector on IoT node hardware using output counters as a learning tool and KNearest Neighbor, Decision Tree and Random Forest as classifiers. The reported accuracy rate for specific malware families varies from 25 percent to 100 percent. Alam et al. used Random Forest to identify malware codes on a dataset of Internet-connected mobile apps. They run APKs in an Android emulator and documented different features, such as memory detail, permissions and a network for classification, and tested their approach using different tree sizes. Their results have shown that the ideal classifier includes 40 trees and a mean square root of 0.0171 has been obtained.

LIMITATIONS OF EXISTING SYSTEM

Although dynamic analysis surpasses the static analysis in many aspects, dynamic analysis also has some drawbacks. Firstly, dynamic analysis requires too many resources relative to static analysis, which hinders it from being deploying

on resource constraint smartphone.

On contrast to the above mentioned methods, anomaly detection engine in our proposed detection system performs dynamic analysis through Dalvik Hooking based on Xposed Framework. Therefore, our analysis module is difficult to be detected by avoiding repackaging and injecting monitoring code.

Overall, previous work focuses on detecting malware using machine learning techniques, which are either misuse-based detection or anomaly-based detection. Misuse based detector tries to detect malware based on signatures of known malware.

PROPOSED SYSTEM:

To the best of our knowledge, this is the first OpCodebased deep learning method for IoT and IoBT malware detection. We then demonstrate the robustness of our proposed approach, against existing OpCode based malware detection systems. We also demonstrate the effectiveness of our proposed approach against junk-code insertion attacks. Specifically, our proposed approach employs a class-wise feature selection technique to overrule less important OpCodes in order to resist junk-code insertion attacks. Furthermore, we leverage all elements of Eigenspace to increase detection rate and sustainability. Finally, as a secondary contribution, we share a normalized dataset of IoT malware and benign applications², which may be used by fellow researchers to evaluate and

benchmark future malware detection approaches. On the other hand, since the proposed method belongs to OpCode based detection category, it could be adaptable for non-IoT platforms. IoT and IoBT application are likely to consist of a long sequence of OpCodes, which are instructions to be performed on device processing unit. In order to disassemble samples, we utilized Objdump (GNU binutils version 2.27.90) as a disassembler to extract the OpCodes. Creating n-gram OpCode sequence is a common approach to classify malware based on their disassembled codes. The number of rudimentary features for length N is CN , where C is the size of instruction set. It is clear that a significant increase in N will result in feature explosion. In addition, decreasing the size of feature increases robustness and effectiveness of detection because ineffective features will affect performance of the machine learning approach.

ADVANTAGES OVER EXISTING SYSTEM

The choices made in choosing the detection technique can determined the reliability and effectiveness of the Android malware detection system.

By using this approach the malicious application can be quickly detected and able to prevent the malicious application from being installed in the device.

Hence, by taking advantages of low false-positive rate of misuse detector and the ability of anomaly detector to detect zero-day malware, a hybrid malware detection method is proposed in this paper, which is the novelty in this paper.

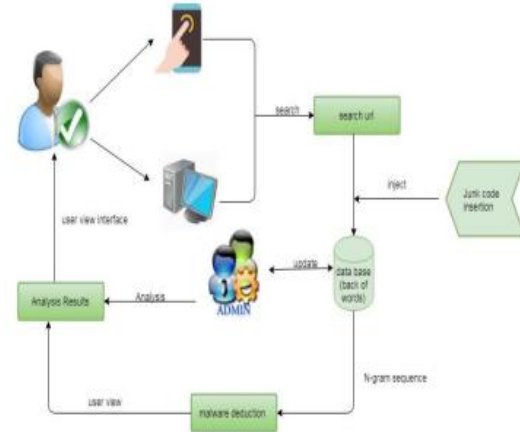
LITERATURE SURVEY:

YEAR	TITLE	METHODOLOGY	RESEARCH PROPOSAL	ALGORITHM
2019	Robust Malware Detection for Internet Of Things Devices Using Deep Eigenspace Learning	Using opcodes via natural processing model	Accuracy-93.1%	Random forest as a classifier
2017	Malware Detection for IOT using Deep EIGEN learning	The data from benign and malignant folders are read, the data are converted to integer value. The Integer value are applied counter function; in which it converts all data to number of counts.	Accuracy-99.8%	Deep neural network to classify the dataset - KNN, SMO, Naive bayes.
2020	Robust malware detection using Space Learning	Solution for mobile malware detection using network traffic flows, which assumes that each HTTP flow is a document and analyzes HTTP flow requests using NLP string analysis.	Accuracy-99.15%	Junk code insertion procedure, N-gram SVM.
2022	Malware Detection in Internet of Things (IOT)	Features are selected and one hot encoding is applied followed by the ensemble classifier based on CNN and LSTM outputs for detection of malware.	Accuracy-99.5%	ANN, CNN, RNN, LSTM LR and Fuzzy C-Mean, GA and PSO are used for better feature selection.
2019	Robust Intelligent Malware Detection using Deep Learning	We present a deep learning based method to detect Internet of Battlefield Things malware via the devices Operational Code sequence.	Accuracy-90.03%	Advanced MLAs such as deep learning

REQUIREMENTS SPECIFICATIONS:

A software requirements specification (SRS) is a document that describes what the software will do and how it will be expected to perform. A software requirements specification is the basis for entire project. It lays the framework that every team involved in development will follow. It's used to provide critical information to multiple teams - development, quality assurance, operation and maintenance. This keeps everyone on the same page. Using the SRS helps to ensure requirements are fulfilled. And it can also help you make decisions about your product's lifecycle - for instance, when to retire a feature. Writing an SRS can also minimize overall development time

and costs. Embedded development teams especially benefits from using an SRS.



SYSTEM ARCHITECTURE

Software Requirements

For developing the application the following are the Software Requirements:

1. Python
2. Django

Operating Systems supported

1. Windows7
2. Windows XP
3. Windows8

Technologies and Languages used to Develop

1. Python

Debugger and Emulator

- Any Browser (Particularly Chrome)

Functional Requirements

- Graphical User interface with the User.

SYSTEM SPECIFICATION:

SOFTWARE REQUIREMENTS:

- **Operating system :** Windows 7 Ultimate.
- **Coding Language :** Python
- **Front-End:** Python
- **Designing:** Html,css,javascript.
- **Data Base:** MySQL

HARDWARE REQUIREMENTS:

For developing the application the following are the Hardware Requirements:

- Processor: Pentium IV or higher
- RAM: 256 MB
- Space on Hard Disk: minimum 512MB
- **System** : Pentium IV 2.4 GHz.
- **Hard Disk** : 40 GB.
- **Floppy Drive** : 1.44 Mb.
- **Monitor** : 14' Colour Monitor.
- **Mouse** : Optical Mouse.
- **Ram** : 512 Mb

MODULE DESCRIPTION

There are three modules can be divided here for this project they are listed as below

- User Activity
- Malware Deduction
- Junk Code Insertion Attacks
- From the above three modules, project is implemented. Bag of discriminative words are achieved.

1. USER ACTIVITY:

User handling for some various times of IOT(internet of things example for Nest Smart Home, Kisi Smart Lock, Canary Smart Security System, DHL's IoT Tracking and Monitoring System, Cisco's Connected Factory, ProGlove's Smart Glove, Kohler Verdera Smart Mirror. If any kind of devices attacks for some unauthorized malware softwares. In this malware on threats for user personal dates includes for personal contact, bank account numbers and any kind of personal documents are hacking in possible.

2. MALWARE DEDUCTION:

1.

Users search the any link notably, not all network traffic data generated by malicious apps correspond to malicious traffic. Many malware take the form of repackaged benign apps; thus, malware can also contain the basic functions of a benign app. Subsequently, the network traffic they generate can be characterized by mixed benign and malicious network traffic. We examine the traffic flow header using N-gram method from the natural language processing (NLP).

2. JUNK CODE INSERTION ATTACKS:

Junk code injection attack is a malware anti-forensic technique against OpCode inspection. As the name suggests, junk code insertion may include addition of benign OpCode sequences, which do not run in a malware or inclusion of instructions (e.g. NOP) that do not actually make any difference in malware activities. Junk code insertion technique is generally designed to obfuscate malicious OpCode sequences and reduce the 'proportion' of malicious OpCodes in a malware.

2.2 ALGORITHM:

Algorithm: Junk Code Insertion Procedure
Input: Trained Classifier D, Test Samples S, Junk Code Percentage k
Output: Predicted Class for Test Samples P

1: $P = fg$

2: for each sample in S do

3: W = Compute the CFG of sample based on Section 4.1

4: R = fselect k% of W's index randomly (Allowd uplicate indices)g

5: for each index in R do
6: $W_{index} = W_{index} + 1$
7: end for
8: Normalize W
9: $e_1; e_2 =$ 1st and 2nd eigenvectors of W
10: $l_1; l_2 =$ 1st and 2nd eigenvalues of W
11: $P = P S D(e_1; e_2; l_1; l_2)$
12: end for 13: return P

2.3 INPUT AND OUTPUT DESIGN

2.3.1 Input Design

The configuration of the input is the relation between the information system and the customer. It involves the creation of requirements and procedures for data preparation and these measures are required to position transaction data in a functional form for analysis and can be accomplished by checking a device for reading data from a written or printed record or by making people lock the data directly into the database. The input architecture focuses on managing the amount of input required, reducing errors, preventing delays, avoiding unnecessary steps and making the process quick. The feedback is built in such a way as to maintain protection and ease of use while maintaining.

2.3.2 Output Design

A standard performance is one that meets the requirements of the end user and communicates the details clearly. In any system, the effects of the processing are transmitted by outputs to users and to other systems. In the production process, it is decided how the material is to be transferred for immediate use, as well as the output of the hard copy. This is the most critical and clear source information to the customer. Effective and insightful performance architecture strengthens the interaction of the device and help users make decisions. □ Designing the output of the machine should continue

in an coordinated, well thought-out manner; the correct output should be produced thus ensuring that every output feature is configured so that the program can be used conveniently and efficiently. When evaluating the program output configuration, they will define the unique performance required to satisfy the requirements. □ Pick the methods to display the details. □ Build a text, study or other format containing information generated by the device. **2.4 DEEP EIGENSPACE LEARNING AND DEEP LEARNING**

Deep Eigenspace Learning

A prevalent data type in machine learning is graphs as a complex data structure for representing relationships between vertices. There are very few algorithms for data mining and deep learning that consider a graph as an input. A logical alternative is therefore to integrate a graph into a vector space. Graph embedding is essentially a bridge between recognizing statistical patterns and graph mining. Eigenvectors and individual values are two characteristic elements in the continuum of the graph, which could turn the adjacency matrix of a graph linearly into a vector space. It denotes ownvectors, uniqueness values and the adjacency or affinity matrix of a line. In this article, for the learning process, we employ a sub-set of v and \dot{y} . $Av = \lambda v$ To obtain substantive knowledge of the structure of CGFs generated, a graph is produced which illustrates the cumulative of all samples in our dataset. The figure below consists of two major diagonal building blocks (marked with red boundaries), indicating that the samples contain two main data distributions. Based on the spectrum theory of the graph, there should be an explicit own gap in the proper values

of the matrix in this case, and it depicts the presence of a gap between π_2 and $\pi_k(k>2)$.

2.4.2 Deep Learning

Deep Learning or Deep structural learning is an evolved version of Neural Network. There are few or several basic, interconnected nodes called neurons in a standard NN. In a few layers, NN's neurons are arranged, namely: an input layer, several unseen layers and an output layer. DL as a "upgraded" NN phenomenon, focuses on deeper understanding of the data structure by focusing on the strengths and functionalities of the hidden layer. Recently, deep learning has been successfully applied to tackle problems across a range of applications, including speech recognition and machine vision. DL types, such as Convolutional Networks, Limited Boltzmann Machines and Sparse Coding.

RESULTS:

SCREENSHOTS:

HOME PAGE: This is screen that is opened once the project is run as shown in screen 1.

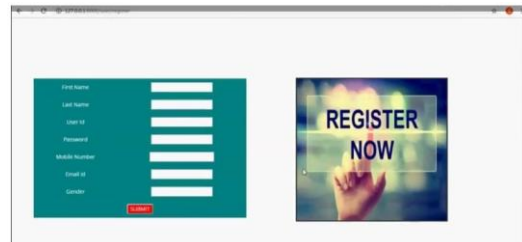
The screen has the title name followed by the tabs which are follows-



HOME PAGE

Registration: A registered user is a user of a website, program, or other system who has previously registered. Registered users normally provide some sort of credentials (such as a username or e-

mail address, and a password) to the system in order to prove their identity: this is known as logging in. Systems intended for use by the general public often allow any user to register simply by selecting a register or sign up function and providing these credentials for the first time.

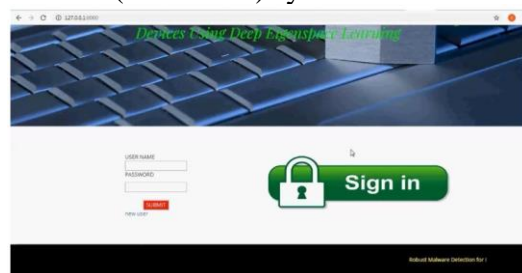


REGISTRATION PAGE



MY DETAILS

Sign in: After registration process user must be sign in with his username and password .A sign in is the period of activity between a user sing in and sign out of a (multi-user) system.



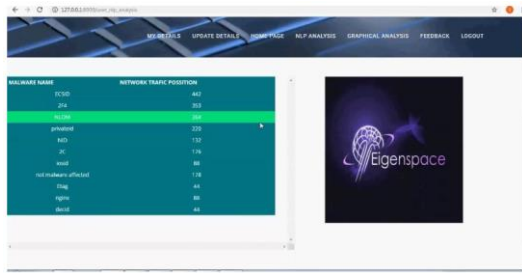
SIGNIN

UPDATE DETAILS: In update tab the user can update his details



UPDATE DETAILS

NLP ANALYSIS:



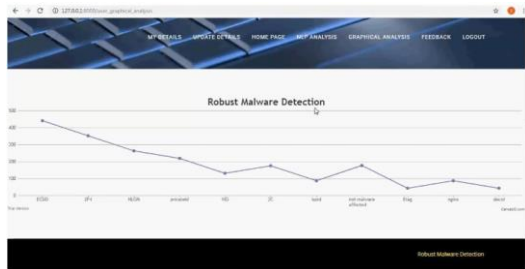
GRAPHICAL ANALYSES:



FEEDBACK:



NLP ANALYSIS GRAPHICAL ANALYSIS:

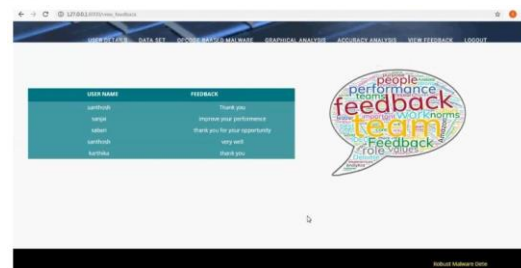


GRAPHICAL ANALYSIS

ADMIN PAGE:



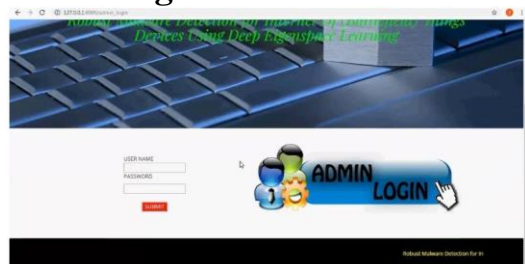
ADMIN PAGE



TESTCASES SCENARIOS:

AND

Admin login:



ADMIN LOGIN

OPCODE BASED MALWARE:



Opcode sequence analyses

TEST CASES	DESCRIPTION	PREDICTIONS	RESULTS
1	USER Sign In	If user sign in should contain username with alphabets and password with only digits	Pass
2	Updating information	Related information must be enter in specified box	Pass
3	Information invalid	If information is not founded	Fail
4	Searching URL	All the category of URL's are listed to search	Pass
5	Verify	If the given URL data is found in database	Pass
6	NLP analysis	Viewing results and graphical representation	Pass
7	Feedback	If the feedback is successfully taken	Pass

CONCLUSION and FUTURE WORK:

IoT, particularly IoBT, will be increasingly important in the foreseeable future. No malware detection solution will be foolproof but we can be certain of the constant race between cyber attackers and cyber defenders. Thus, it is important that we maintain persistent pressure on threat actors. In this paper, we presented an IoT and IoBT malware detection approach based on class-wise selection of Op- Codes sequence as a feature for classification task. A graph of selected features was created for each sample and a deep Eigenspace learning approach was used for malware classification. Our evaluations demonstrated the robustness of our approach in malware detection with an accuracy rate of 98.37% and a precision rate of 98.59%, as well as the capability to mitigate junk code insertion attacks.

FUTURE ENHANCEMENT:

Android is a new and fastest growing threat to malware. Currently, many research methods and antivirus scanners are not hazardous to the growing size and diversity of mobile malware. As a solution, we introduce a solution for mobile malware detection using network traffic flows, which assumes that each HTTP flow is a document and analyzes HTTP flow requests using NLP string analysis. The N-Gram line generation, feature selection algorithm, and SVM algorithm are used to create a useful malware detection model. Our evaluation demonstrates the efficiency of this solution, and our trained model greatly improves existing approaches and identifies malicious leaks with some false warnings. The harmful detection rate is 99.15%, but the wrong rate for harmful traffic is 0.45%. Using the newly discovered malware further verifies the performance of the proposed system. When used in real environments, the sample can detect 54.81% of harmful applications, which is better than other popular anti-virus scanners. As a result of the test, we show that malware models can detect our model, which does not prevent detecting other virus scanners. Obtaining basically new malicious models Virus Total detection reports are also possible. Added, Once new tablets are added to training.

REFERENCES:

- <https://ieeexplore.ieee.org/document/8302863/>
- https://www.google.com/search?rlz=1C1CHBF_enIN854IN854&sxsrf=ALeKk03MZoDsC7Y3dMmotBRglFMni5-FUw%3A1590205665898&ei=4ZzIXtG4NqOY4-EP4PK3GA&q=robust+malware+detection+for+iot+devices+using+deep+eigens



- pace+lear
ning+github&oq=robust+m
alware+detect&gs_lcp
- https://www.researchgate.net/publication/323405239_Robust_Malware_Detection_for_Internet_Of_Battlefield_Things_Devices_Using_Deep_Eigenspace_Learning
 - <https://github.com/nsslabcuus/Malware>
 - <https://towardsdatascience.com/malware-detection-using-deep-learning-6c95dd235432>
 - <https://www.jetbrains.com/help/pycharm/configuring-project-and-ide-settings.html>
 - <https://sourceforge.net/projects/staruml/>