# AUTOMATED ANDROID MALWARE DETECTION USING DEEP LEARNING AND ML MODELS FOR CYBERSECURITY

[1]**D. Srivalli**, [2]**A. ADITHI REDDY**, [3]**A. JAHNAVI REDDY**, [4]**AMRUTH DEEPTHI**, [5]**B.SAANVI REDDY**

[1]Assistant Professor, Department of CSE(DS), Malla Reddy Engineering College for Women (Autonomous Institution – UGC, Govt. of India), Hyderabad, INDIA.

[2345]UG,Department of CSE(DS), Malla Reddy Engineering College for Women (Autonomous Institution – UGC, Govt. of India), Hyderabad , INDIA.

**Abstract:** Current technological advancement in computer systems has transformed the lives of humans from real to virtual environments. Malware is unnecessary software that is often utilized to launch cyber-attacks. Malware variants are still evolving by using advanced packing and obfuscation methods. These approaches make malware classification and detection more challenging. New techniques that are different from conventional systems should be utilized for effectively combating new malware variants. Machine learning (ML) methods are ineffective in identifying all complex and new malware variants. The deep learning (DL) method can be a promising solution to detect all malware variants. This paper presents an Automated Android Malware Detection using Optimal Ensemble Learning Approach for Cybersecurity (AAMD-OELAC) technique. The major aim of the AAMD-OELAC technique lies in the automated classification and identification of Android malware. To achieve this, the AAMD-OELAC technique performs data preprocessing at the preliminary stage. For the Android malware detection process, the AAMD-OELAC technique follows an ensemble learning process using three ML models, namely Least Square Support Vector Machine (LS-SVM), kernel extreme learning machine (KELM), and Regularized random vector functional link neural network (RRVFLN). Finally, the hunter-prey optimization (HPO) approach is exploited for the optimal parameter tuning of the three DL models, and it helps accomplish improved malware detection results. To denote the supremacy of the AAMD-OELAC method, a comprehensive experimental analysis is conducted. The simulation results portrayed the supremacy of the

AAMD-OELAC technique over other existing approaches.

## 1. INTRODUCTION

Malware is a software which can cause potential threats to a computer, server, client, or computer network. Malware causes damage after it is implanted or introduced into a target's computer and it is in the form of an executable codes, script files and other softwares. The codes are known as viruses, ransomware, spyware, adware, worms , Trojans ,scareware and many other forms. The commonly used methods for protecting against malware is to prevent the software from gaining access to the target computer includes antivirus software, firewalls and many other techniques. The main uses of them include preventing their access to target computers, checking the presence of suspicious activities, recover from malware attacks. Another strategy to differentiate malware apps from genuine Android apps includes sophisticated dynamic and static analysis tools to detect and classify malicious apps automatically. There are encryption techniques which will decrease the chances of malwares from being

detected. To avoid this problem, we can study Android apps to extract permissions which are sensitive that are widely used in Android malwares. An automated malware detection system is used to fight against malwares and assist Android app marketplaces to detect and remove unknown malicious apps. Static analysis tools are used to extract source codes or byte codes, often traversing the paths of programs to check for some unique and hidden resources. Static analysis approaches are used for different tasks which includes the behaviour assessment of Android apps, detection of application clones, automatic test case generations, or for uncovering non functional issues related to performance.

The important point which is to be noted is that the code is not executed or run but the tool itself is executed. The source code is the input to the tool and the mined features are the output.eg:- Drebin Dynamic program analysis is the analysis of Android applications by executing the programs on a virtual environment like Android Studio. The target programs must be executed with test inputs to produce the behavior.

System calls are analyzed to monitor the behaviour of Android applications.eg:- TaintDroid Malware classification is an open problem commonly rectified by employing machine learning techniques. Permissions and API calls are extracted w Man is able to detect behaviors which are sensitive from Android applications. Most of the detections are based on the difference of permissions detected by benign apps and malware apps. By analysing the permissions requested and api call usages, benign app and malware app samples can effectively expose abnormal behaviors and finally distinguish malware from many genuine applications.

So considering the drawbacks of the above techniques we propose a new model which is based on the extracted permissions from the apks and uses deep learning techniques to formulate the model.

## 2. SYSTEM DESIGN

Most of the malware detection tools uses the manual of lists of features based on permissions, api calls, sensitive resources, intents, etc., which are difficult to come by. To address this problem, we study the different real Android applications to mine hidden patterns of malware and are able to extract highly sensitive permissions that are widely used in Android malware. Benign apps are downloaded from apkpure.com which is a free site of benign apks from google playstore. Malicious apps are downloaded and are extracted from virusShare.com and Contagio Mini Dump. Features like Api related Permissions are considered to develop the system. Permission Distribution Permissions[1] from malwares and benign apps are identified. By analysis, Access_wifi_state,SendSms etc are commonly used by malwares. The requested permissions of the android applications are declared in a file called Android manifest of the respective apks. From the manifest files, permissions are extracted and are converted to a csv file. A large number of permissions are identified in the previous step. Out of which a few must be selected for further processing. For that Mann Whitney test[2] is employed. For each permission, if a particular app uses that permission,the corresponding permission is set as 1 or else it is set to 0. These values are indicated by p values.[2] Therefore two sets of samples

are generated. One to represent one specific permission usage of malicious apps and the other to represent specific permission usage of benign apps. In the previously created input file, a comparison test is applied. For each permission, the average values are computed for each of the feature vector. So from two sets of samples, we compute the average values. And those permissions with higher average values will be selected as the feature vector for training.

## 2.1 Malware Detection

This feature vector is divided into two. One can be used for training the model and the other can be used for determining the model parameters. The first feature vector is fed to the classifier. The classifier employed here is the Neural Networks and K-Means Clustering Algorithm. Two trained models will be created. The second feature vector is given as input to the model to determine the model parameters like accuracy, precision, recall etc. Unknown apks are then given as input to the model so that the model will predict these apks as benign or malicious.

## 3. IMPLEMENTATION

Here, we take a closer look at how the system was implemented. The whole system was developed using python language.

Benign apps are downloaded from apkpure.com. Malicious apps are downloaded and extracted from virusshare.com and Contagio Minidump. A total of 135 benign apps were collected. A total of 327 malicious apps were collected. The features namely permissions are extracted using Python 3.7 in Spyder. A package called Androguard[5] is used to extract manifest files from apks. The extracted permissions are correctly displayed on the screen. Feature Selection is done using Extra Tree Classifier which is included as a built in package in python. Feature selection is performed successfully using the dataset. Feature Vectors are generated by using Mann-Whitney test[3]. It is implemented using the inbuilt package called scipy.stats in Python 3.7. The weights and their corresponding feature names are written to a csv file.

Training phase receives a training dataset which is a csv file. The model is trained using Neural networks and k-means clustering algorithm. Output of this phase is a confusion matrix and graphs showing the dataset which are classified correctly and incorrectly. The model is tested using different samples of both malware and benign apps. The output of the feature map as well as the prediction will be printed on the screen. The feature map generated for the training data sample is given in the figure below:-



Figure 1: Feature Map of the training dataset

The extracted permissions from the testing data sample is shown in the figure below:-



Figure 2: Extracted Features from the test dataset

The feature map of the test data sample and the prediction is shown in the figure below:-



Figure 3: Feature Map of test data

4 CONCLUSIONS

The implemented system collects datas in the form of Android apks from various Internet sources. The apks are extracted to collect features which are basically permissions. A feature vector is created based on the permissions and the given apks. This is the input to the ML algorithms to build a trained model.

Unknown applications are used as input. An overall accuracy of 88 percent is achieved. The main limitations of the model include:-    A large dataset must be collected to avoid overfitting problem.

The extracted permissions are limited because the number of malicious applications on the internet are very less.

This system considers the differences of malware and benign apps but it does not consider the categories of benign apps which can be useful for malware detection.    This system is open to Mimicry and Pollution attacks.    This system bypass malwares using Java Reection and bytecode encryption.

## 5 REFERENCES

[1] H. Rathore, A. Nandanwar, S. K. Sahay, and M. Sewak, ''Adversarial superiority in Android malware detection: Lessons from reinforcement learning based evasion attacks and defenses,'' Forensic Sci. Int., Digit. Invest., vol. 44, Mar. 2023, Art. no. 301511.

[2] H. Wang, W. Zhang, and H. He, ''You are what the permissions told me! Android malware detection based on hybrid tactics,'' J. Inf. Secur. Appl., vol. 66, May 2022, Art. no. 103159.

[3] A. Albakri, F. Alhayan, N. Alturki, S. Ahamed, and S. Shamsudheen, ''Metaheuristics with deep learning model for cybersecurity and Android malware detection and classification,'' Appl. Sci., vol. 13, no. 4, p. 2172, Feb. 2023.

[4] M. Ibrahim, B. Issa, and M. B. Jasser, ''A method for automatic Android malware detection based on static analysis and deep learning,'' IEEE Access, vol. 10, pp. 117334–117352, 2022.

[5] L. Hammood, İ. A. Doğru, and K. Kılıç, ''Machine learning-based adaptive genetic algorithm for Android malware detection in auto-driving vehicles,'' Appl. Sci., vol. 13, no. 9, p. 5403, Apr. 2023.

[6] P. Bhat and K. Dutta, ''A multi-tiered feature selection model for Android malware detection based on feature iscrimination and information gain,'' J. King Saud Univ.-Comput. Inf. Sci., vol. 34, no. 10, pp. 9464–9477, Nov. 2022.

[7] D. Wang, T. Chen, Z. Zhang, and N. Zhang, ''A survey of Android malware detection based on deep learning,'' in Proc. Int. Conf. Mach. Learn. Cyber

Secur. Cham, Switzerland: Springer, 2023, pp. 228–242.

[8] Y. Zhao, L. Li, H. Wang, H. Cai, T. F. Bissyandé, J. Klein, and J. Grundy, ''On the impact of sample duplication in machine-learning-based Android malware detection,'' ACM Trans. Softw. Eng. Methodol., vol. 30, no. 3,

pp. 1–38, Jul. 2021.

[9] E. C. Bayazit, O. K. Sahingoz, and B. Dogan, ''Deep learning based malware detection for Android systems: A comparative analysis,'' Tehnički vjesnik, vol. 30, no. 3, pp. 787–796, 2023.

[10] H.-J. Zhu, W. Gu, L.-M. Wang, Z.-C. Xu, and V. S. Sheng, ''Android malware detection based on multi-head squeeze-and-excitation residual network,'' Expert Syst. Appl., vol. 212, Feb. 2023, Art. no. 118705.

[11] K. Shaukat, S. Luo, and V. Varadharajan, ''A novel deep learning-based approach for malware detection,'' Eng. Appl. Artif. Intell., vol. 122, Jun. 2023, Art. no. 106030.

[12] J. Geremias, E. K. Viegas, A. O. Santin, A. Britto, and P. Horchulhack, ''Towards multi-view Android malware detection through image-baseddeep learning,'' in Proc. Int. Wireless Commun. Mobile Comput. (IWCMC), May 2022, pp. 572–577.