# Authorization based Malware Avoidance in Android

**Mr. Vinod R. Thakare[1], Mr. Akash R. Dudhe[2], Mr. Ajay D. Nanure[3], Ms. Kanchan B. Korke [4]**

[1] Assistant Professor, Dept. of E&TC, Jagadambha College of Engg. & Tech., Yavatmal, Maharashtra, INDIA

[2] Assistant Professor, Dept. of E&TC, Jagadambha College of Engg. & Tech., Yavatmal, Maharashtra, INDIA

[3] Assistant Professor, Dept. of E&TC, Jagadambha College of Engg. & Tech., Yavatmal, Maharashtra, INDIA

[4] Assistant Professor, Dept. of E&TC, Jagadambha College of Engg. & Tech., Yavatmal, Maharashtra, INDIA

**Abstract:** A portion of users of feature phones are rapidly switching to smart phones, and many users update and install mobile apps without giving security any thought. As a result, malware finds smart phones to be an appealing target, particularly Android devices. Because of this, it is also crucial to use a framework to identify malicious applications as soon as they are installed on mobile. In this study, we present a powerful framework for detecting malicious programs according to Android authorizations. Authorizations taken from the AndroidManifest.xml record are suggested to be used as highlights when creating machine learning categorization methods in order to extract malware-related highlights. The experiment's findings demonstrate that the suggested approach, when applied to a sample of data of 18,490 applications, achieves an accuracy of above 95%. (11,672 generous; 6,818 malicious).By doing this, it is capable of precisely detecting all malignant as well as generous programs.

**Keywords:** Android, smart phones, machine learning, malware detection, and Android permissions.

**Abbreviations:** APK (Android application package) files, Machine learning (ML), Malware/Benign Group (PMBG), Outlier Malware/Benign Group (OMBG).

## 1. INTRODUCTION

Ironically, Android's ubiquity also piques the interest of cybercriminals, who then produce nefarious apps that can steal sensitive data and undermine mobile platforms. Android does not allow users to install apps from unauthorized sources, such as third-party application stores and file-sharing websites, like certain other rival smart-mobile device networks like iOS do. The problem of malware attack has become so serious that, according to a subsequent assessment, 97% of most of the android malware targets Android smart phones [3].In 2016, there were over 3.25 million new malware apps found. It typically translates to the 10 second introduction of a malicious Mobile app [4].Such malicious apps are designed to carry out many types of attacks using Trojans, worms, abuses, and viruses. Some well-known malicious programmers have more than 50 variants, which makes it incredibly difficult to distinguish from each of them [5].The remainder of the paper is structured as follows: segment 2 surveys prior significant ideas, segment 3 presents the suggested show, segment 4 shows the dataset, segment 5 explores setup, segment 6 discusses experiment outcomes, and segment 7 concludes.

## 2. RELATED WORK

In the Android platform, the application's request for consent is crucial in managing access privileges. Mobile applications lack the necessary license to access client data by default, which compromises security. During setup, the user must provide the app access to all the resources it requests. Developers must note in the AndroidManifest.xml record the permissions requested for the assets. However, not every consent that has been announced has been the required permission for that particular application.

According to Ref. [6], designers frequently claim authorizations that the programme does not truly need, which makes it challenging to identify applications that are acting maliciously. All of the authorizations for the resources needed by the app are listed in the Android Manifest.xml record, which antimalware examines. Stowaway [6] was the first to recognise the licence over benefit issue on Android, in which a software requests extra permissions than it really needs. To ascertain the API calls that the programme generated, Stowaway conducts an inactive assessment. Next, it maps the authorizations demanded by the API calls. In 940 Android application test results, they discovered that one-third of the applications are given

too much privilege. It is unable to interpret API calls generated by programmes using Java reflections.

In [7], authors proposed a simple malware detection technique that, as it were, analyses the exhibit record and extract the information such as authorizations, expectation channels (activity, category, and need), prepares title, and counts the number of permissions that have been redefined in order to detect applications that behave badly. They gather the information, contrast it with the list of keywords supplied by the approach, and then calculate the malignant level. People used Weka [8], a viable tool for data mining, to determine boundary esteem. When all is said and done, they compare the threat score to the edge esteem and label the software as malware if the danger score is higher than the threshold esteem. They tested the effectiveness of the suggested arrangement using 365 samples, and the result provides 90% efficient identification. It can be used in different location designs to reliably discriminate malware, and it just requires the evaluation of show records, making it a sparingly used component. Furthermore, it is unable to detect adware samples.

Given that malware authors frequently announce more authorizations than they need for the app inside this abstract record,Y. Haung and colleagues [9] developed an approach for significantly improved location of consent-based spyware discovery that takes into account the analysis of both asked and required authorizations. Additionally, it examines the aspects that are easy to recover before labeling the application as malicious or generous. This calls for the use of three distinct labeling types: location-based labeling, scanner-based labeling, and mixed labeling. If an app is downloaded via the official Google app store, it is classified as benign; however, the program is labeled as hazardous if it was downloaded from a malicious source. Currently, if the antivirus scanner is used for labeling, deems a programme to be generous, the app is named as helpful, and the same is true for malware cases. The app is tagged using both location-based as well as scanner-based labels as part of blended labeling. Following labeling, all tests are divided into three datasets, and access to these datasets is requested. Machine learning algorithms like Gullible Bayes, Ada Boost, Support Vector Machine, and Choice Tree are then used to evaluate the datasets

[9]. Upon the basis of the data produced by these classifiers, it is possible to assess how well the consent-based discovery strategy performs. In [9], authors tested a data collection that included 124,769 different kinds of apps and 480 malicious ones. They investigated how authorization malware detection actually worked and found that it can successfully identify more than 81 percent of dangerous software tests. The suggested method provides a rapid filter for locating malware, however, fully depend on the outputs of the classifiers because the performance metrics they give are not optimal.

Jaguar was demonstrated by Sanz Borja et al. [10] for the tracking of fraudulent apps by looking at the sought consents in the application. People connected several classifier calculations using a dataset of 357 good applications and 249 bad apps to examine the malicious behavior of apps using approval labels such as and display in the AndroidManifest.xml record. The system has a high localization rate, but the results it produces have a high false positive rate, and it is insufficient for effectively discovering malware because it still needs information about other aspects and dynamic behavior.

KIRIN is a tool designed by Enck et al. [11] that provides simple certification at the moment of establishment. It describes the security requirements and, in essence, checks the authorizations requested by the app with its security requirements. If the programme does not pass all of the security requirements, it is certified as malware. If an app is identified as malware, the establishment of the app is terminated. After trying 311 applications obtained from the official Android display, creators discovered that 5 failed to meet the requirements. The suggested method is compact because it essentially scans the Menifest.xml file. Because the information provided for application certification is insufficient for malware location, KIRIN's limitations include the possibility that it may potentially label certain legitimate apps as spyware.

It's possible that DroidMat [12] is a tool that pulls data from manifest records, such as access rights, message conveying through motives, including System calls tracking, in order to assess application activity. It uses K-means clustering to increase the ability to find

malware and to categorize programmes as benign or malicious using KNN calculations [13].It is speedier than Androgaurd [14] because it requires less time to classify the 1,738 apps as harmful or benign.As it doesn't demand manual labour or vigorous reenactment, it is also less taxing. However, as an inactive based location technique, it is unable to identify malware like DroidKngFu and Base Bridge that powerfully stack the evil substance.

## 3. PROPOSED MODEL

### 3.1 APK files

The subsection depicts the overall structure and details of Android applications, as seen in Figure 1.The fundamental tools for implementation, operation and application, including source code, materials, resources, certifications, and manifest records,are all contained in zip files, which are essentially what Android applications are when viewed at a high level. The APK file's content that is most important to our work is explained in the following subsections. Users familiar with Android applications and their strategy may want to skip this area.
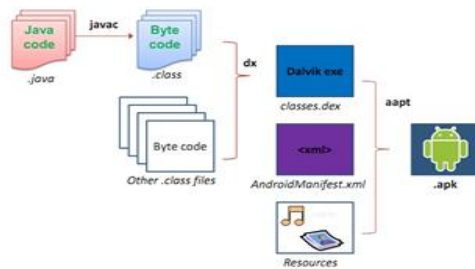


Figure 1: APK Architecture

### 3.2 AndroidManifest.xml

Each Mobile application has a display record that contains essential information that the Android framework needs to obtain in order to run the app's code.More specifically, the system won't permit an application to use or access any assets or components that aren't specified in the display record [15].

### 3.2.1 Application components

The application components are the initial crucial part of the display file and each programme. A small amount of an application's overall usefulness is incorporated into each component. When a client taps an application's icon to open the client interface, the framework can then carry out client interaction. Thus, it is crucial to be aware that one framework may contain multiple components, each of which aids in carrying out a different aspect of its unique capabilities. The below will illustrate them:

Client interfaces can be seen via the exercises component. Every activity component provides a specific screen on which clients can carry out a few operations or make requests. Activities include, for instance, screens with calendars to explore in a hotel booking application or displays with an active camera and a button to request a picture in a photo application. The function of this component is to give the application a means of interacting with clients.

- The administrations unit typically works in the background on specific duties continuously without interaction with clients. A music player app, for instance, might be activated as a bonus when a client exits the most webpages. The practical element may be used by a malicious application to send data to inaccessible sites, receive commands from them, or monitor its environment.

- By utilizing the recipient's component, often referred to as broadcast recipients, apps are able to receive notifications or messages from the system or from other applications. Time-zone changes and low-battery alerts are frequent system notifications. As not all communications on an Android device are useful to an application, a channel can be used to get rid of messages that are not interesting to a particular application.

- A vendor's component, also known as a substance supplier, connects a database of one application to the database of another application that needs to exchange information. A content provider component must be named in the display file of an application that offers shared databases. Additionally, a uses-permission for

restricting access to a supplier must be specified in the display record.

### 3.2.2 Uses-Permissions

When an application permits other applications to access its information, the other is given permission. At the time of installation, the client grants the uses-permissions, which are explicit authorizations that the application needs from the Android framework.All preferred usage permissions for an Android application must be specified inside the AndroidManifest.xml file.READ CONTACTS is an illustration of a uses-permission; if granted, it permits reading of all contacts in the contact list.
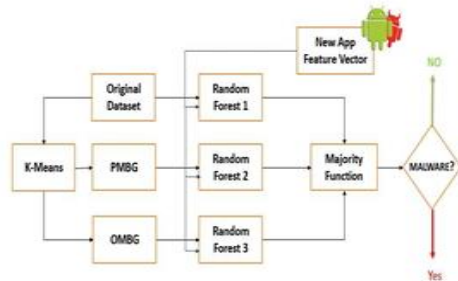


Figure 2: Proposed Model for Android Malware Detection

Over time, machine learning (ML) classifiers have contributed to the creation of smart frameworks in a number of fields. Spyware on PC and mobile platforms is increasingly being found using machine learning (ML) techniques .The supervised machine learning foundation of our research which is taken from a named dataset and used to create and prepare a demonstration, and unsupervised machine learning, which uses the highlights displayed in the past area.ML calculations used in our analysis include K-Means and Irregular Woodland.

In this research, we offer a machine learning-based approach for the localization of zero-day Android malware that uses a composite model of the homogeneous classifiers employed in numerous simultaneous combination plans. In order to offer a single categorization option for various uses, the system seeks to combine the advantages of several machine learning approaches. The original dataset discussed in area V was divided into two groups using K-Means Analysis in order to identify pure malware and benign groups (PMBG) and outlier malware and

benign groups (OMBG), respectively. Random Forest Algorithm was then used to forecast any Malicious programs. A location approach's building squares are shown in Figure 2.The parallel detector's ML computations for its component parts include:

### 3.2.3 K-Means Clustering Algorithm

In this concept, the clustering analysis is carried out using K-Means calculations. Here, characterizing k centroids, each one associated with a cluster, requires the most attention. Dealing out each application to the closest centroid is the next stage. In order to reduce the cost function's sum of squares, the K-means clustering divides the dataset into sections. [16]There aren't many doubts about the transmission of underlying information when using the information-driven K-means clustering technique. When expanding, it makes sure that there is a local minimum of model effort, which shortens the time it takes for clusters to converge on big datasets. [14]K-means is a straightforward method for dividing a data set into k clusters. [16]These are the specifications of the K-means clustering calculation:

1. Choose k irregular occurrences from the subset of information being prepared to serve as the centroids of the clusters C1 through Ck.

2.for each preparation event X:

    a. Determine the Euclidean separate D $(Ci,X)$, where i=1...k. Locate the cluster Cq that is most near X.

b. Give X to Cq. change the Cq centroid. (The mathematical mean of the instances within a cluster serves as the centroid.)

3. Recur with Step 2 until the cluster centroid for C1, C2... Ck stabilizes according to the mean-squared-error criterion.

4. For every test Z instance:

    a. When i=1...k, calculate the distinct D $(Ci,Z)$. Find Z's closest Cr cluster by doing some analysis.

    b. Using the Edge run the show, categorise Z as a novelty or a frequent occurrence.

### 3.2.4  Random Forest Classification Algorithm

Calculating an arbitrary woods is intended to create a forest of random trees. It is a classification learning approach  functions by creating a decision tree during training and producing the lesson via individual trees. The learners are given the generic bootstrap conglomerating, or sacking, approach by the training algorithm for irregular woodlands. The benefits of arbitrary woodland calculations include their effectiveness in handling enormous databases with hundreds of input variables without wiping any variables .Furthermore, it retains correctness even when a significant portion of the information is lost thanks to an efficient technique for assessing missing data. [17]In the examination of permission-based malware finding, the third decision tree technique is determined to be random timberland calculation. Here are the steps for calculating irregular woodland:

1. Assume that M factors make up the classifier and that N cases are being prepared.

2. The choice at a hub of the tree is determined by m input factors; m should be significantly smaller than M.

3. Choose a preparation set from all N available preparation scenarios N times using substitution. By predicting their classes, use the remaining cases to evaluate the tree's error.

4. Randomly choose m variables to use as the basis for the decision at each hub of the tree. Based on these m variables in the training set, choose the best portion.

5. No trees are clipped and are fully established.

## 4. DATASET

The Android Dataset has been used in this paper in [18].This dataset, named a Unique Dataset and extracted by [18], contains 18,490 Records and 151 authorization highlights, as shown in Table 1.



| Features (Type) | Features (keywords) |
|---|---|
| Permissions | ACCESS_CHECKIN_PROPERTIES<br>ACCESS_COARSE_LOCATION<br>ACCESS_FINE_LOCATION<br>ACCESS_LOCATION_EXTRA_COMMANDS<br>ACCESS_MOCK_LOCATION<br>BIND_DEVICE_ADMIN<br>CALL_PHONE<br>CALL_PRIVILEGED<br>CAPTURE_AUDIO_OUTPUT<br>CAMERA[1] |

[1] A Total of 151 Permissions are used.

Table 1: Permissions Features

## 5. EXPERIMENT SETUP

In this data analysis, the collection of data was split into two groups based on label agreement (11672 Kind and 6818 Malware).The exclusions from each group were extracted using K-Means computations, as shown in Table 2.This data is split into four groups: 348 Outlier Kind, 11324 Unadulterated Kind, 1096 Exceptional Malware, and 5722 Pure Malware. The Unadulterated generous group and the Unadulterated malware bunch are combined to form the Pure malware/benign group (PMBG).Exception Malware and Exception Generous are both members of the Outlier malware/benign group (OMBG).For both these categories and a special dataset, the Irregular Timberland Classifier was used. For each classification, a Decision Tree Show is generated to determine Android malware. For the latest Android app, first extract its 151 authorization Highlights. Secondly, employ the yield methods to ascertain whether this Android software is harmful or not by employing the subsequent strategy: if two classifiers predict the same name (Malware), this malicious Android software is malicious; if they predict something other, it's safe to use the Android app.

| | PMBG | OMBG | Total |
|---|---|---|---|
| Benign | 11324 | 348 | 11672 |
| Malware | 5722 | 1096 | 6818 |
| Total | 17046 | 1444 | 18490 |

**Table2:**Dataset Groups

The 10-fold cross approval method is coupled to the matrix in order to evaluate how the classifier display is being used.As a result, the data is separated into 10 equal, non-overlapping sections, numbered k1, k2, k3, to k10.Each step of the evaluation uses a prepared display from the previous 9 parts and one partition as test information. The results are discovered to be the middle value, which provides the classifier with its final performance data. The k-fold cross approval method, a popular ML assessment technique, is appropriate for our objective of determining the relative effectiveness of the ML classifiers in identifying obscure harmful programmes that are simulated by the non-overlapping testing subsets. To investigate the parallel classifier technique to Android malware specific area, the following execution metrics are used.

**True positive ratio (TPR):**It is frequently the proportion of dangerous apps inside the dataset that were accurately identified as malicious apps.

**TNR (True Negative Dataset Ratio):** The ratio of correctly categorized good apps to all the generous apps in the dataset.

**False positive ratio (FPR):** The percentage of falsely labeled good apps to all the good apps in the dataset.

**False negative proportion (FNR):** The ratio of malicious applications that were incorrectly identified as being in the dataset's total number of noxious apps.

**Accuracy (ACC):** This is the degree to which the classifier provided by (TPR + TNR)/(TPR + TNR + FPR + FNR) is accurate.

Error ratio (ERR): Typically, AUC (computed from: 1 - Area under ROC): ROC is the receiver operation characteristics bent.

AUC is a measurement of the area under the ROC that illustrates the classifier's foresight. Higher AUC

classifiers provide far better predictive capability and can provide for better affectability adjustment.

## 6. Results and Discussions

The initial round of tests was carried out utilizing each of the several possible classification computations for the dimensional model in order to acquire consistent results for researching the simultaneous classifiers technique to new malware discovery. The outcomes from all of these classifiers are compiled in Table 3 and opposed with those of other classifiers, like ID3 and optimistic Bayes classifiers, in Table. **4,5**

| Dataset | Performances Metrics | | | | | | |
|---------|------|------|------|------|------|------|------|
|         | TPR  | TNR  | FPR  | FNR  | ACC  | ERR  | AUC  |
| Original | 0.836 | 0.941 | 0.059 | 0.164 | 0.889 | 0.111 | 0.959 |
| PMBG    | 0.813 | 0.950 | 0.045 | 0.186 | 0.881 | 0.119 | 0.957 |
| OMBG    | 0.913 | 0.858 | 0.141 | 0.087 | 0.885 | 0.115 | 0.951 |

Table 3: Performance Results from the Three Individual Classifiers (Random Forest)

| Dataset | Performances Metrics | | | | | | |
|---------|------|------|------|------|------|------|------|
|         | TPR  | TNR  | FPR  | FNR  | ACC  | ERR  | AUC  |
| Original | 0.873 | 0.963 | 0.038 | 0.127 | 0.918 | 0.082 | 0.971 |
| PMBG    | 0.848 | 0.966 | 0.034 | 0.152 | 0.907 | 0.093 | 0.974 |
| OMBG    | 0.984 | 0.925 | 0.075 | 0.016 | 0.954 | 0.046 | 0.999 |

Table 4: Performance Results from the Three Individual Classifiers (ID3)

| Dataset | Performances Metrics | | | | | | |
|---------|------|------|------|------|------|------|------|
|         | TPR  | TNR  | FPR  | FNR  | ACC  | ERR  | AUC  |
| Original | 0.695 | 0.944 | 0.056 | 0.305 | 0.819 | 0.181 | 0.896 |
| PMBG    | 0.732 | 0.952 | 0.046 | 0.225 | 0.842 | 0.158 | 0.915 |
| OMBG    | 0.941 | 0.873 | 0.127 | 0.059 | 0.907 | 0.093 | 0.963 |

Table 5: Performance Results from the Three Individual Classifiers (Naïve Bayes).

According to Table 3, among the classification techniques, the Moment Irregular Timberland

classifier has the lowest detection ratio (TPR) and overall exactness.The Primary Arbitrary Forest classifier exhibited the best malware location performance with around 91.7% location proportion, while the Third Arbitrary Forest classifier displayed higher location with 95.4%.Overall accuracy/error rates were best achieved by the Third Irregular Forest Classifier.The following criteria, which are listed in Table 6, were used to decide whether or not a new Android app was malware.This harsh Android app is categorized as malware in the uncommon circumstance that two classifiers predict the same label (Malware), whilst another is categorized as kind.

| Predication Sequence of Random Forest Classifiers (RF1,RF2,RF3) | The Final Output of Majority Function Malware(M) or Benign(B)? |
|---|---|
| (B,B,B) | B |
| (B,B,M) | B |
| (B,M,B) | B |
| (M,B,B) | B |
| (M,M,M) | M |
| (M,M,B) | M |
| (M,B,M) | M |
| (B,M,M) | M |

Table 6: Majority Function Outputs

We see the recommended classification method as a genuinely workable solution to locate Android malware and improve existing systems.In specifically, for finding Android malware that is a zero-day flaw for which no genuine results have been produced.

## 7. CONCLUSIONS

We conducted a permission-based study of malware categorization for Android applications in three main steps as part of the display consider. The Android dataset was split into two groups in the first stage, according to their names (11672 Generous and 6818 Malware).K-Means calculations have been used to separate the exceptions from each bunch inside the Moment step. The dataset is divided into four groups: 348 Exception Generous, 11324 Pure Generous, 1096 Exception Malware, and 5722 Pure Malware. One bunch called Pure Malware and Benign Bunch combines Immaculate Malware Bunch and Pure Kind Bunch (PMBG).The Outlier Malware and Generous Bunch is a collection of this Exception Malware and Outlier Benign (OMBG).In the final phase, a decision tree show was built by the Random Timberland

Classifier employed for both these classes (PMBG and OMBG) and for each unique dataset, which is used to base malicious Android applications. We anticipate that the results of the current study may serve as the basis for next malware detection research.

## References

1.Available: https://www.idc.com/promo/smartphonemarket-share/os , IDC, "Smartphone os market share, 2017 q1." [Online]. last accessed 2023/07/23.

2. Available: https://www.statista.com/statistics/281106/number-ofandroid-app-downloads-from-google-play/ , Statista, "Cumulative number of apps downloaded from the google play as of may 2016." [Online]. last accessed 2023/07/23.

3. G. Kelly. "Report: 97 percent of mobile malware is on Android,""This is the simple way to stay safe," according to Forbes Tech in 2014.

4.G.DATA, Available: https://www.gdatasoftware.com/blog/2017/04/29712-8-400-new-android-malware-samples-every-day "8,400 new android malware samples every day." [Online].

5. Symantec, "Latest intelligence for march 2016," in Symantec Official Blog, 2016.

6. Available: https://www.truststc.org/pubs/848.html. [Accessed: 06- Nov-2015]. Android Permissions Demystified." [Online].

7. R. Sato, D. Chiba, and S. Goto, "**Detecting Android Malware by Analyzing Manifest Files**," pp. 23–31, 2013.

8. [Online]. http://www.cs.waikato.ac.nz/ml/weka/ is available. "Weka 3 - Data Mining in Java Using Open Source Machine Learning Software. [Accessed on December 16, 2012]

9. C.-Y. Huang, Y.-T. Tsai, and C.-H. Hsu, "Performance evaluation on permission-based detection for android malware," Adv. Intell. Syst. Appl. - Vol. 2, vol. 21, pp. 111–120, 2013.

10. B. Sanz, I. Santos, C. Laorden, X. Ugarte-Pedrero, P. G. Bringas, and G. Álvarez, "PUMA: Permission usage to detect malware in android," Adv. Intell. Syst. Comput., vol. 189 AISC, pp. 289–298, 2013.

11. W. Enck, M. Ongtang, and P. McDaniel, "On lightweight mobile phone application certification,"

Proc. 16th ACM Conf. Comput. Commun. Secur. - CCS "09, pp. 235–245, 2009.

12D.-J. Wu, C.-H. Mao, T.-E. Wei, H.-M. Lee, and K.-P. Wu, . "DroidMat: Android Malware Detection through Manifest and API Calls Tracing," 2012 Seventh Asia Jt. Conf. Inf. Secur., pp. 62–69, 2012.

13. L. Kozma, "k Nearest Neighbors algorithm ( kNN )," 2008.

14. Zaw, W. \& Aung, Z., Permission-Based Android Malware Detection. IJSTR. 2013.

15. Enck, W., Ongtang, M., \& McDaniel,). Understanding android security. IEEE Security and Privacy, 7, 50-57. P. (2009b

16. Abu Samra, A. A., Yim, K. \& Ghanem, O. A., 2013. Analysis of Clustering Technique in Android Malware Detection. IEEE.

17. Breiman, L., 2001. Random Forests. Machine Learning... pp. 45 (1): 5-32.

18. Mike Yang: CAPIL: Component-API Linkage for Android Malware Detection (2015, master project).