

# Evaluation of Android Malware Detection Based on System Calls

Marko Dimjašević, Simone Atzeni,  
Zvonimir Rakamarić  
University of Utah, USA  
{marko, simone, zvonimir}@cs.utah.edu

Ivo Ugrina  
University of Zagreb, Croatia  
ivo@iugrina.com

## ABSTRACT

With Android being the most widespread mobile platform, protecting it against malicious applications is essential. Android users typically install applications from large remote repositories, which provides ample opportunities for malicious newcomers. In this paper, we evaluate a few techniques for detecting malicious Android applications on a repository level. The techniques perform automatic classification based on tracking system calls while applications are executed in a sandbox environment. We implemented the techniques in the MALINE tool, and performed extensive empirical evaluation on a suite of around 12,000 applications. The evaluation considers the size and type of inputs used in analyses. We show that simple and relatively small inputs result in an overall detection accuracy of 93% with a 5% benign application classification error, while results are improved to a 96% detection accuracy with up-sampling. Finally, we show that even simplistic feature choices are effective, suggesting that more heavyweight approaches should be thoroughly (re)evaluated.

## 1. INTRODUCTION

The global market for mobile devices has exploded in the past several years, and according to some estimates the number of smartphone users alone reached 1.7 billion worldwide in 2014. Android is the most popular mobile platform, holding nearly 85% of the global smartphone market share. One of the main advantages of mobile devices such as smartphones is that they allow for numerous customizations and extensions through installing applications from public application markets. The largest of such markets (e.g., Google Play, Apple App Store) have more than one million applications available for download each, and there are more than 100 billion mobile device applications installed worldwide.

This clearly provides a fertile environment for malicious activities, including development and distribution of malware. A recent study [23] estimates that the total amount of malware across all mobile platforms grew exponentially at the rate of 600% between 03/2012 and 03/2013. Around 92% of malware applications found in this study target Android. In a related study [34], similar statistics are reported — the number of malicious applications in the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*IWSPA'16, March 11 2016, New Orleans, LA, USA*

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ISBN 978-1-4503-4077-9/16/03...\$15.00

DOI: <http://dx.doi.org/10.1145/2875475.2875487>

Google Play store grew around 400% from 2011 to 2013, while at the same time the number of malicious applications removed annually by Google has dropped from 60% in 2011 to 23% in 2013. Due to the sharp increase in the total amount of malware, the percentage of removed malware dropped significantly despite the fact that the absolute number actually increased from roughly 7,000 in 2011 to nearly 10,000 in 2013. While companies such as Google regularly scan their application markets using proprietary tools, this process is often ineffective as the above numbers illustrate. There are also unofficial, open markets where often no scanning is being performed, partially because there is a lack of solid freely available solutions and tools. As a consequence, Android malware detection has been an active area of research in the past several years, both in industry and academia.

In this paper, we evaluate existing and propose novel dynamic Android malware detection techniques based on tracking system calls, all of which we implemented as a free and open-source tool called MALINE. Our work was initially inspired by a similar approach proposed for desktop malware detection [30], albeit we provide simpler feature encodings and an Android-specific tool flow. We provide several encodings of behavior fingerprints of applications into features for subsequent classification. We performed an extensive empirical evaluation on a set of more than 12,000 Android applications. We analyze how the quality of malware classifiers is affected across several dimensions, including the choice of an encoding of system calls into features, the relative sizes of benign and malicious data sets used in experiments, the choice of a classification algorithm, and the size and type of inputs that drive a dynamic analysis. Furthermore, we show that the structure of system call sequences observed during application executions conveys in itself a lot of information about application behaviors. Our evaluation sheds light on several such aspects, and shows that the proposed combinations can be effective: our technique yields an overall detection accuracy of 93% with a 5% benign application classification error. Finally, we provide guidelines for domain experts when making choices on malware detection tools for Android, such as MALINE.

Our approach provides several key benefits. By guarding the users at the repository level, a malicious application is detected early and before it is made publicly available for installation. This saves scarce energy resources on the devices by delegating the detection task to a trusted remote party, while at the same time protecting users' data, privacy, and payment accounts. System call monitoring is out of reach of malicious applications, i.e., they cannot affect the monitoring process. Hence, our analysis that relies on monitoring system calls happens with higher privileges than those of malicious applications. In addition, tracking system calls entering the kernel (and not calls at the Java library level) enables us to

capture malicious behavior potentially hidden in native code. Since our approach is based on coupling of a dynamic analysis with classification based on machine learning, it is completely automatic. We require no source code, and we capture dynamic behavior of applications as opposed to their code properties such as call graphs; hence, our approach is mostly immune to common, simple obfuscation techniques. The advantages of our approach make it complementary to many existing approaches, such as the ones based on static analysis.

Our contributions are summarized as follows:

- We propose a completely automatic approach to Android malware detection on the application repository level using system calls tracking and classification based on machine learning, including a novel heuristics-based encoding of sequences of system calls into features.
- We implement the approach in a tool called MALINE, and perform extensive empirical evaluation on more than 12,000 applications. We show that MALINE effectively discovers malware with a very low rate of false positives.
- We compare several feature extraction strategies and classifiers. In particular, we show that the effectiveness of even very simplistic feature choices (e.g., frequency of system calls) is comparable to much more heavyweight approaches. Hence, our results provide a solid baseline and guidance for future research in this area.

## 2. OUR APPROACH

Our approach is a three-phase analysis, as illustrated in Fig. 1. The first phase is a dynamic analysis where we track system calls<sup>1</sup> during execution of an application in a sandbox environment and record them into a log file. In the second phase, we encode the generated log files into feature vectors according to several representations we define. The last phase takes the feature vectors and applies machine learning [20] to learn to discriminate benign from malicious applications.

### 2.1 Dynamic Analysis Phase

As our approach is based on concrete executions of applications, the first phase tracks and logs events at the operating system level that an application causes while being executed in a sandbox environment. The generated event logs serve as a basis for the subsequent phases of our analysis. Unlike numerous static analysis techniques, this approach reasons only about events pertaining to the application that are actually observed in the operating system.

A user’s interaction with Android through an application results in events being generated at the operating system level, which are rendered as system calls. In our work, we automatically emulate this interaction as explained in detail in § 3. For that reason, we execute every application in a sandbox environment and observe resulting system calls in a chronological order, from the very beginning until the end of its usage. The output of this phase is a log file containing chronologically ordered system calls: every line consists of a time stamp, the name of the system call, its input values, and the return value, if any. Having the system calls recorded chronologically enables us to construct various feature vectors that

<sup>1</sup>A system call is a mechanism for a program to request a service from the underlying operating system’s kernel. In Android, system calls are created by information flowing through its multi-layered architecture, starting from an application on top; hence, they capture its behavior.

characterize the application’s behavior with different levels of precision, as explained in the next section.

More formally, let  $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$  be a set of system call names containing all the system calls available in the Android operating system for a given processor architecture. Then a system call sequence  $\sigma$  of length  $m$ , representing the chronological sequence of recorded system calls in a log file, is a sequence of instances of system calls  $\sigma = (q_1, q_2, \dots, q_m)$ , where  $q_i \in \mathcal{S}$  is the  $i$ th observed system call in the log file. Such call sequences are passed to the feature extraction phase.

### 2.2 Feature Extraction Phase

As explained earlier, how features are picked for the feature vector is important for the machine learning classification task. Therefore, we consider two representations for generating a feature vector from a system call sequence  $\sigma$ . Our simpler representation is concerned with how often a system call happens, while our richer representation encodes information about dependencies between system calls. Both representations ignore system call information other than their names and sequence numbers (e.g., invocation time, input and output values), as it can be seen from the definition of  $\sigma$ . Once we compute a feature vector  $\mathbf{x}$  for every application under analysis according to a chosen representation, we form a feature matrix by joining the feature vectors such that every row of the matrix corresponds to one feature vector.

**System Call Frequency Representation.** How often a system call occurs during an execution of an application carries information about its behavior [6]. A class of applications might be using a particular system call more frequently than another class. For example, some applications might be making considerably more I/O operation system calls than known goodware, indicating that the increased use of I/O system calls might be a sign of malicious behavior. Our simple system call frequency representation tries to capture such features. In this representation, every feature in a feature vector represents the number of occurrences of a system call during an execution of an application. Given a sequence  $\sigma$ , we define a feature vector  $\mathbf{x} = [x_1 x_2 \dots x_{|\mathcal{S}|}]$ , where  $x_i$  is equal to the frequency (i.e., the number of occurrences) of system call  $s_i$  in  $\sigma$ . In experiments in § 4, we use the system call frequency representation as a baseline comparison against the richer representation described next.

**System Call Dependency Representation.** Our system call dependency representation was inspired by previous work that has shown that a program’s behavior can be characterized by dependencies formed through information flow between system calls [16]. However, we have not been able to find a tool for Android that would provide us with this information and also scale to analyzing thousands of applications. Hence, we propose a novel scalable representation that attempts to capture such dependencies by employing heuristics. As we show in § 4, even though our representation is simpler than the one based on graph mining and concept analysis from the original work [16], it still produces feature vectors that result in highly accurate malware detection classifiers.

For a pair of system calls  $q_i$  and  $q_j$  in a sequence  $\sigma$ , where  $i < j$ , we define the distance between the calls as  $d(q_i, q_j) = j - i$ . We then approximate a potential data flow relationship between a pair of system calls using the distance between the calls in a sequence (i.e., log file). For example, if two system calls are adjacent in  $\sigma$ , their distance will be 1. Furthermore, let  $w_{g,h}$  denote the weight of a directed edge  $(s_g, s_h)$  in a *system call dependency graph* we generate. The system call dependency graph is a complete digraph with the set of vertices being the set of all the system call names  $\mathcal{S}$ , and hence having  $|\mathcal{S}|^2$  edges. Then,  $w_{g,h}$  for a sequence  $\sigma$  is

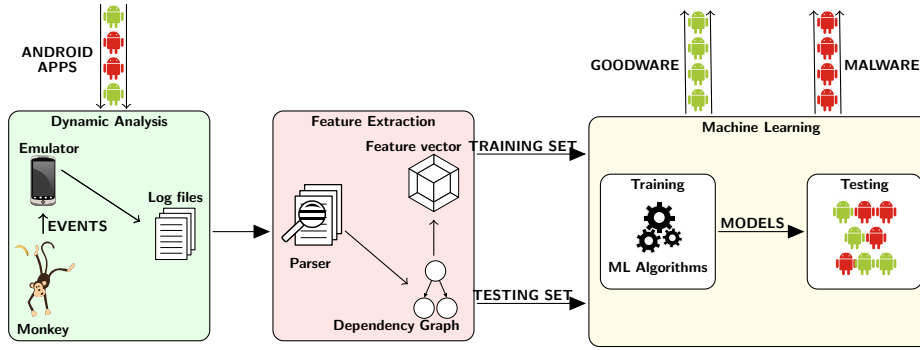


Figure 1: MALINE tool flow divided into three phases.

computed as:

$$w_{g,h} = \begin{cases} 0, & \text{if } g = h \\ \sum_{\substack{i < j < k \\ q_i = s_g, q_j = s_h}} \frac{1}{d(q_i, q_j)}, & \text{otherwise} \end{cases}$$

where  $k$  is the minimal index such that  $q_i = q_k$  and  $i < k \leq |\sigma|$ . Informally, the closer the pair is in a sequence, the more it contributes to its edge weight in the graph. Hence, instead of explicitly observing a data flow between system calls, our representation captures it implicitly: it is based on a simple observation that the closer a pair of system calls is in a sequence, the more likely it is that there is a data flow between the pair.

From a sequence  $\sigma$ , we compute weights  $w_{g,h}$  for every system call pair  $(s_g, s_h) \in \mathcal{S}^2$ . For  $g$  and  $h$  such that  $w_{g,h} = 0$ , we still consider edge  $(s_g, s_h)$  to exist, but with the weight of 0. Since each application is executed only once during our dynamic analysis phase, we generate one system call dependency graph per application.

We generate a feature vector  $\mathbf{x}$  of an application by taking edge weights from its system call dependency graph. For every directed edge  $(s_g, s_h)$  there is a corresponding feature in  $\mathbf{x}$ , and hence the dimensionality of  $\mathbf{x}$  is  $|\mathcal{S}|^2$ . Given a sequence  $\sigma$ , we define a feature vector  $\mathbf{x} = [x_1 x_2 \dots x_{|\mathcal{S}|^2}]$ , where  $x_i$  is equal to  $w_{g,h}$  such that  $i = (g-1) \cdot |\mathcal{S}| + h$ .

### 2.3 Machine Learning Phase

We use the generated feature vectors for our applications (i.e., feature matrices) together with provided malware/goodware labels to build classifiers. We experimented with several of the most popular and effective classifiers: support vector machines (SVMs), random forest (RF), LASSO and ridge regularization [20] and used the double cross-validation approach to tune parameters of classifiers.

When a probabilistic classifier is used, a threshold that appropriately tunes the trade-off between sensitivity and specificity can be studied using *receiver operating characteristic* (ROC) curves [20]. Generating ROC curves is especially valuable to the users of malware detectors such as ours, since they can use them to fine-tune sensitivity vs. specificity depending on the intended usage. Hence, we generate ROC curves for the most interesting classifiers.

We have chosen to use around 33% samples as malware and the rest as goodware. Although this approach does not generate a perfectly balanced design, it tries to represent the goodware population as best as possible while still keeping the high percentage of malware samples and computational costs at a practical level. In addition, we explored what can be achieved by balancing the design through resampling strategies of *up-sampling* (or over-sampling) the minority class and *down-sampling* (or under-sampling) the majority class [25] implemented through bootstrapping.

## 3. IMPLEMENTATION

We implemented our approach in a tool called MALINE, and Fig. 1 shows its tool flow. The implementation comes as a free and open reproducible research environment to foster further evaluation, development, and research in this area.<sup>2</sup> MALINE heavily utilizes our own build of the Android Software Development Kit (SDK). The SDK includes the Android Emulator, which runs a virtual machine (VM) with the Android operating system. Every application MALINE analyzes is installed, executed, and monitored in the VM. The tool primarily resides on the host machine and relies on the Android Debug Bridge (*adb*) to communicate with the VM.

### 3.1 Host and Emulator

MALINE consists of a number of smaller components. We implemented multiple interfaces on the host side, ranging from starting and monitoring an experiment with multiple emulator instances running in parallel to machine-learning differences between applications based on the processed data obtained from emulator instances. It is the host side that coordinates and controls all such activities. For example, it creates and starts a pristine installation of Android in an emulator instance, then installs an application in it, starts the application, and waits for it to finish so it can analyze system calls the application has made during its execution.

We use the emulator, which is built on top of QEMU [5], in the dynamic analysis phase of our approach (see Fig. 1). For every application we create a pristine sandboxed environment since the emulator enables us to easily create a clean installation of Android. It is important that each application is executed in a clean and controlled environment to make sure nothing is left behind from previous executions and to be able to monitor the execution. Hence, every application’s execution is completely independent of executions of all the other analyzed applications.

### 3.2 Automatic Execution of Applications

In order to scale to thousands of applications, our dynamic analysis phase implements an automatic application execution process. The process starts with making a clean copy of our default VM. The copy contains only what is installed by default in a fresh installation of the Android operating system from the Android Open Source Project. Once the installation boots, we use *adb* to send an application from the host machine to the VM for installation. Next, we start the application and immediately begin tracing system calls related to the operating system process of the application with the *strace* tool. The system calls are recorded into a log file.

<sup>2</sup>The MALINE tool is available from <https://github.com/soarlab/maline> under the GNU Affero GPLv3 license.

We simulate a user interaction with an Android device by injecting both *internal* and *external* events into the emulator. Internal events are sent to the application itself, such as screen clicks, touches, and gestures. We use the Monkey tool [2] as our internal event generator (see Fig. 1). It sends a parameterized number of the events to the application, with a 100 ms pause period between consecutive events if applicable.<sup>3</sup> Unlike internal events, which are delivered to the application, external events are delivered to the emulator and include events that come from interacting with an external environment. In our experiments, for external events we focus on generating text messages and location updates only since those are sometimes related to malicious behaviors.

Even though a system call log file contains rich information (e.g., time stamp, input and output values), thereby forming a chronological sequence of low-level operations, we preserve only system call names and their order. We stop an application execution when all internal events generated by Monkey are delivered and executed, and then we pull the log file from the VM to the host machine for parsing. Next, we apply a feature vector representation, either the system call frequency or dependency representation as explained in § 2. The output is a textual feature vector file per log file, i.e. per application, listing all the features. Finally, we combine all the feature vectors into a single matrix where each matrix row corresponds to one feature vector, i.e. one application.

### 3.3 Classification

Using the feature matrix generated from logs and previously obtained labels denoting malware/goodware for applications, we proceed with performing a classification. We experimented with several classification algorithms: random forest, SVMs, LASSO, and ridge regression. An implementation of SVMs is based on libSVM [7], while all the other algorithms are implemented in R [32] using the language’s libraries. The scripts are heavily parallelized and adjusted to be run on large machines or clusters. For example, running a random forest model on a feature matrix from a system call dependency graph sample requires 32 GB of RAM in one instance of 5-fold cross-validation.

## 4. EVALUATION

We evaluated MALINE using a set of 32-core machines with 128 GB of RAM running Ubuntu 12.04. The machines are part of the Emulab infrastructure [37]. We wrote scripts to automatize and parallelize our experiments, without which our extensive experimental evaluation would not be possible. In our experiments, we use only the x86 Android emulator; the resulting x86 system call set  $\mathcal{S}$  has 360 system calls. Note that our accompanying technical report presents more detailed experimental results [9].

### 4.1 Input Data Set

We obtained applications from Google Play as goodware and from the Drebin dataset [3] as malware. Before we could start using the collected applications in MALINE, we performed a filtering step. First, we removed applications that we failed to consistently install in the Android emulator. For example, even though every Android application is supposed to be self-contained, some applications had dependencies that were not installed at the time; we do not include such applications in our final data set. Second, we removed all applications that we could not consistently start or that would crash immediately. For example, unlike typical Android applications, application widgets are miniature application views that

<sup>3</sup>The pause between two consecutive events may not be applicable for actions that are time-dependent, such as screen tapping.

do not have an Android Activity, and hence they cannot be started from a launch menu.

Applications in the Drebin dataset were collected between Aug 2010 and Oct 2012, and filtered by their collectors to contain only malicious applications. To the best of our knowledge, this is the latest verified malware application collection of its size used by researchers. The malicious applications come from more than 20 malware families, and are classified based on how an application is installed and activated, or based on its malicious payloads [41]. The aim of our work is not to explore specifics of the families; many other researchers have done that. Therefore, in our experiments, we make no distinction between malicious applications coming from different families. The Drebin dataset contains 5560 malware applications; after filtering, our malicious data set contains 4289 of those applications.

We obtained the benign data set in Feb. 2014 by utilizing a crawler tool. The tool searched Google Play for free-of-charge applications in all usage categories (e.g., communication, music and audio, business), and randomly collected applications with at least 50,000 downloads. We stopped our crawler at 12789 collected Google Play applications; after filtering, our benign data set contains 8371 of those applications. Note that we make a reasonable assumption that applications with more than 50,000 downloads are benign.

### 4.2 Configurations

We explore effects of several parameters in our experiments. The first parameter is the number of events we inject with Monkey into the emulator during an application execution. The number of events is directly related to the length of the execution. We insert 1, 500, 1000, 2000, and 5000 events. It takes 229 seconds on average (with a standard deviation of 106 seconds) for an application execution with 500 events and 823 ( $\pm 816$ ) seconds with 5000 events.<sup>4</sup> That includes the time needed to make a copy of a clean virtual machine, boot it, install the application, run it, and download log files from the virtual machine to the host machine.

The second parameter is a flag indicating if a benign background activity should be present while executing the applications in the emulator. The activity comprises of inserting SMS text messages and location updates into the emulator. We experiment with the activity only in the 500-Monkey-event experiments, while for all the other experiments we include no background activity.

It is important to ensure that consistent sequences of events are generated across executions of all applications. As Monkey generates pseudo-random events, we use the same pseudo-random seed value in all experiments.

### 4.3 Experimental Results

The total number of system calls an application makes during its execution directly impacts its feature vector, and potentially the amount of information it carries. Hence, we identified the number of injected events, which directly influences the number of system calls made, as an important metric to track. The number of system calls observed per application in the dynamic analysis phase of an experiment varies greatly. For example, in an experiment with 500 Monkey events it ranges from 0 (for applications that failed to install and are filtered out) to over a million. Most of the applications in this experiment had less than 100,000 system calls in total.

**Feature Matrices.** After the dynamic analysis and feature extraction phases (see § 2) on our filtered input set, MALINE generated 12

<sup>4</sup>The standard deviations are relatively large compared to the averages because some applications crash in the middle of their execution. We take recorded system call traces up to that point as their final execution traces.

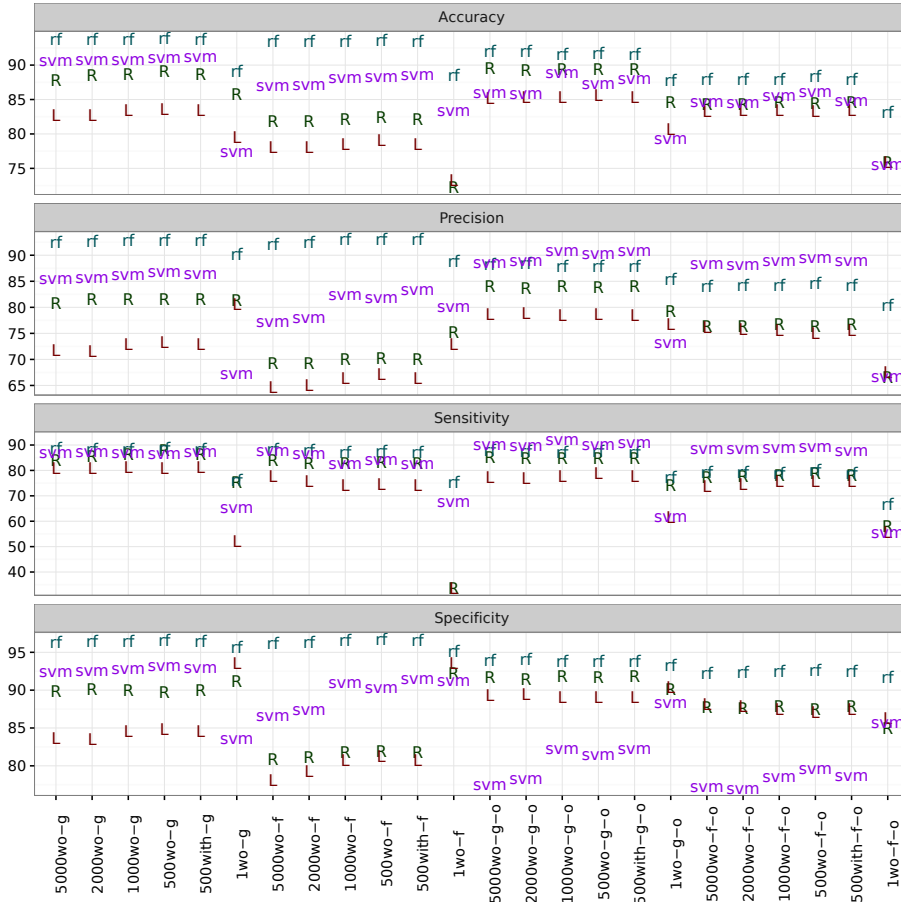


Figure 2: Comparison of the quality of our classifiers through different quality measures (averaged on cross-validation folds). Labels on the  $x$  axis are written in the short form where  $wo$  stands for *without background*,  $with$  stands for *with background*,  $f$  stands for *freq*,  $g$  stands for *graph*,  $o$  at the end denotes that 0-1 matrices were used, and the numbers at the beginning represent numbers of generated events. 1-event experiments have a 2.3% smaller set of applications.

different feature matrices. The matrices are based on varying experiment configurations including: 5 event counts (1, 500, 1000, 2000, 5000), 2 system call representations (frequency- and dependency-graph-based), and the inclusion of an optional benign activity (SMS messages and location updates) for experiments with 500 events. We refer to these matrices with  $X_{rep}^{size}$ , where  $rep \in \{freq, graph\}$  is the used representation of system calls and  $size$  is the number of generated events. In addition, we denote an experiment with the benign background activity using an asterisk.

Obtained feature matrices generated according to the system call dependency representation exhibited high sparsity. This is not surprising since the number of possible system call pairs is 129600. Hence, all columns without a nonzero element were removed from our matrices. Both the frequency and dependency feature vector representations resulted in different nonzero elements in the feature matrices. However, those differences could have only a small or no impact on the quality of classification, i.e., it might be enough only to observe if something happened encoded as zero/one values. Therefore, we have created additional feature matrices by replacing all nonzero elements with ones to measure the effect of feature matrix structure on the classification.

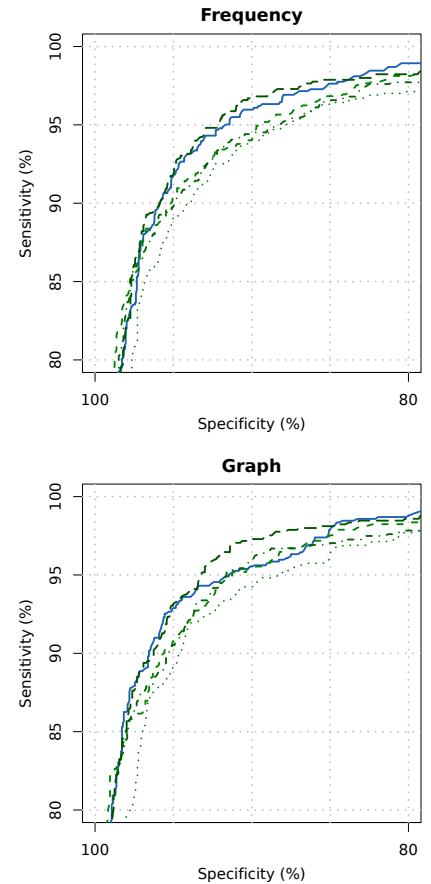


Figure 3: ROC curves for 5-fold cross-validation with RF model on  $X_{freq}^{5000}$  and  $X_{graph}^{5000}$ .

**Cross-validated Comparison of Classifiers.** Reduced feature matrices (just feature matrices from now on) and goodwill/malware labels are input to the classification algorithms we used: support vector machines (SVMs), random forest (RF), LASSO, and ridge regression. To avoid possible overfitting, we employed double 5-fold cross-validation on the set of applications to tune parameters and test models. To enable comparison between different classifiers for different feature matrices, the same folds were used in the model building among different classification models. Prior to building the final model on the whole training set, all classifiers were first tuned by appropriate model selection techniques to derive the best parameters.

The built classifiers were then validated on the appropriate test sets. Fig. 2 shows measures of the quality of prediction (i.e., accuracy, sensitivity, specificity, and precision) averaged between cross-validation folds for different classifiers. If a positive example (i.e., malware in our case) is classified into the positive (resp., negative) group, we obtained a *true positive* (resp., *false negative*). Analogously, we define *true negative* and *false positive*. The threshold for probabilistic classifiers was set at the usual level of 0.5. Since changes to this threshold can have an effect on the sensitivity and

the specificity of classifiers, a usual representation of the effect of these changes is given by ROC curves (see Fig. 3 for an example). In the paper we give ROC curves only for the random forest models (as the best classifiers judging from the cross-validated comparison) with the largest number of events (5000).

As it can be seen from Fig. 2, 1-event quality measures are consistently the worst in each category, often with a large margin. This indicates the importance of leveraging the information gathered while driving an application using random events. Moreover, the random forest algorithm consistently outperforms all other algorithms across the four quality measures. In the case where feature matrices have weights instead of zeros and ones, it shows only small variations across all the input parameters, i.e., the number of events inserted by Monkey, whether there was any benign background activity, and the chosen feature vector representation. Other classification algorithms perform better on the dependency than on the frequency representation. Of the other algorithms, SVM is most affected by the presence of the background activity, giving worse sensitivity with the presence, but on the other hand giving better specificity.

When the weights in the feature matrices are replaced with zeros and ones, thereby focusing on the structure of the features and not their values, all the algorithms consistently perform better on the dependency than on the frequency feature vector representation. However, a comparison within an algorithm based on the weights or zeros and ones in the feature matrices is not straightforward. Random forest clearly performs worse when zeros and ones are used in the feature matrices. LASSO and ridge typically perform better in all the quality measures apart from sensitivity for the zeros and ones compared to the weights.

If a domain expert in Android malware detection is considering to apply MALINE in practice, there are several practical lessons to be learned from Fig. 2. The expert can choose to use only the random forest algorithm as it consistently provides the best outcomes across all the quality measures. To reduce the time needed to dynamically analyze an application, it suffices to provide 500 Monkey events as an application execution driver. Furthermore, the presence of the benign background activity does not make much of a difference. On the other hand, to provide few execution-driving events to an application does not suffice. Finally, if the time needed to learn a classifier is crucial and more important than sensitivity, the expert can choose the frequency feature vector representation since it yields almost as good results as the dependency one, but with far smaller feature vectors.

Fig. 3 shows that there is not much variability between 5 different folds from the cross-validation of the best-performing algorithm, namely random forest. This indicates a high stability of the random forest model on the input data set regardless of the choice of training and test sets. It is up to the domain expert to make the trade-off choice in tuning a classifier towards either high sensitivity or specificity. The choice is directly related to the cost of having false positives, the benefits of having more true positives, etc. For example, the domain expert may choose the dependency graph feature vector representation and fix the desired specificity level to 95%; from the graph ROC curve in Fig. 3 it follows that the sensitivity level would be around 93%.

**Exploring the Effect of Matrix Sparsity.** Sparsity of feature matrices can sometimes lead to overfitting. Although we significantly reduce the sparsity with the removal of columns with all zeros, this just removes non-informative features and sparsity is still relatively high (25% for graph representations). To be sure that the effect seen in the cross-validation comparison is real, we performed additional exploration by adopting the idea of permutation tests [29].

Due to prohibitively high computational costs, we used only one classification model to explore the effect of sparsity. We chose the random forest classifier, since it gave the best results on the cross-validation comparison, and the 5000-event matrices. Prior to building a classifier, we permute application labels. As before, we applied 5-fold cross-validation on permuted labels, thus obtaining quality of prediction on the permuted sample. This procedure is repeated 1000 times. Average accuracies of the obtained classifiers were compared to the accuracy of the RF model from Fig. 2 and they were all significantly lower — the best is at 83% for the system call dependency representation. Although 1000 simulations is not much in permutation models, it still reduces the probability of accidentally obtaining high quality results just because of sparsity. **Exploring the Effect of Unbalanced Design.** Since the number of malware applications in our input set is half the number of goodware, we have an *unbalanced design*. Hence, we employed down/up-sampling through bootstrapping to explore if we could get better results using balanced designs (the same number of malware and goodware). Here, we used only the RF classifier to keep computational costs feasible.

Up- and down-sampling exhibited the same effect on the quality of prediction for all feature matrices: increasing sensitivity at the cost of decreasing specificity. This does not come as a surprise since we have equated the number of malware and goodware applications, thereby giving larger weights to malware applications in the model build. However, the overall accuracy for models with down-sampling was lower than for the unbalanced model, while for models with up-sampling it was higher (up to 96.5% accuracy with a 98% sensitivity and 95% specificity). To explore the stability of results under down- and up-sampling, these methods were repeated 10 times; the standard deviation of accuracies between repeats (on the percentage scale) was 0.302.

## 5. RELATED WORK

There is a large body of research on malware detection in contexts other than Android (e.g., [8, 16, 21, 24, 26, 30, 31, 33, 35, 43]). While our work was originally inspired by some of these approaches, we primarily focus in this section on more closely related work on Android malware detection. Ever since Android as a mobile computing platform has become popular, there is an increasing body of research on detecting malicious Android applications, and we split it into static and dynamic analysis techniques.

**Static Techniques.** Static techniques are typically based on source code or binary analyses that search for malicious patterns (e.g., [15, 36]). For example, static approaches include analyzing permission requests for application installation [1, 14, 18], control flow [27, 28], signature-based detection [15, 19], and static taint-analysis [4].

Stowaway [13] is a tool that detects over-privilege requests during the application install time. Enck et al. [38] study popular applications by decompiling them back into their source code and then searching for unsafe coding security issues. Yang et al. [40] propose AppContext, a static program analysis approach to classify benign and malicious applications. AppContext classifies applications using machine learning based on the contexts that trigger security-sensitive behaviors. It builds a call graph from an application binary and after different transformations it extracts the context factors via information flow analysis. It is then able to obtain the features for the machine learning algorithms from the extracted context. In the paper, 202 malicious and 633 benign applications from the Google Play store are analyzed. AppContext correctly identifies 192 malicious applications with an 87.7% accuracy. Gascon et al. [17] also use call graphs to detect malware. Once they extract function call graphs from Android applications, they apply

a linear-time graph kernel in order to map call graphs to features. These features are given as input to SVMs to distinguish between benign and malicious applications. They conducted experiments on 135,792 benign and 12,158 malware applications, detecting 89% of the malware with 1% of false positives.

**Dynamic Techniques.** Dynamic analysis techniques consist of running applications in a sandbox environment or on real devices in order to gather information about the application behavior. Dynamic taint analysis [11, 39] and behavior-based detection [6, 10] are examples of dynamic approaches. Our approach analyzes Android applications dynamically and captures their behavior based on the execution pattern of system calls. Some existing works follow similar approaches.

Dini et al. [10] propose a framework (MADAM) for Android malware detection which monitors applications at the kernel and user level. MADAM detects system calls at the kernel level and user activity/idleness at the user level to capture the application behavior. Their extremely preliminary and limited results, considering only 50 goodwill and 2 malware applications, show 100% of an overall detection accuracy. Crowdroid [6] is another behavior-based Android malware detector that uses system calls and machine learning. As opposed to our approach, Crowdroid collects information about system calls through a community of users. A lightweight application, installed in the users' devices, monitors system calls (frequency) of running applications and sends them to a centralized server, which performs a classification. Crowdroid was evaluated on a limited number of goodwill and only 2 real malware applications, obtaining detection accuracies of 100% for one and 85% for the other.

## 6. THREATS TO VALIDITY

**Application Crashes.** Given that Monkey generates sequences of pseudo-random input events, it is to be expected that it can drive an application into a state that does not handle certain kinds of events, causing a crash. Depending on an experiment, we observed from 29% to 49% applications crash, which could bias our empirical results. However, it is important to note that the crash rate of goodwill and malware applications is roughly the same. Therefore, application crashes do not bring in a classification bias.

**Age of Applications.** Our goodwill data set was downloaded in 2014, while our malware applications are from 2010 – 2012. Because the Android operating system's API evolved from 2010 to 2014, it could mean our approach learns differences between APIs, and not between benign and malicious behaviors. Unfortunately, we could not obtain older versions of applications from Google Play as it hosts only the most recent versions. In addition, to the best of our knowledge, a more recent malware data set does not exist. Hence, we manually downloaded 2010 – 2012 releases of 92 applications from F-Droid [12], an Android application repository offering multiple releases of free and open-source applications; we assumed the applications to be benign. We classified them using MALINE, and we got specificity of around 88%. Compared to the specificities from Fig. 2, which are around 96%, this might indicate that MALINE performs API difference learning to some extent. A comparison using a much larger set of applications across different releases would need to be performed to draw strong conclusions.

**Hidden Malicious Behavior.** Malicious behavior may occasionally be hidden and triggered only under very specific circumstances. As our approach is based on random testing, we might miss such hard-to-reach behaviors, which could affect our ability to detect such application as malicious. Such malware is not common though, and ultimately we consistently get sensitivity of 87% and more using MALINE.

**Detecting Emulation.** As noted in previous work [8, 22, 31], malware could potentially detect it is running in an emulator, and alter its behavior accordingly. MALINE does not address this issue directly. However, an application trying to detect it is being executed in an emulator triggers numerous system calls, which likely leaves a specific signature that can be detected by MALINE. We consistently get sensitivity of 87% and more using MALINE, and hence at most 13% of malware in our experiments successfully disguised as goodwill. Finally, Chen et al. [8] show that only less than 4% of malware in their experiments changes its behavior in a virtualized environment.

**System Architecture and Native Code.** While the majority of Android-powered devices are ARM-based, MALINE uses an x86-based emulator for performance reasons. Few Android applications — less than 5% according to Zhou, et al. [42] — contain native libraries typically compiled for multiple platforms, including x86, and hence executable by MALINE. Nonetheless, the ARM and x86 architectures have different system calls: with the x86- and ARM-based emulators we observed applications utilizing 360 and 209 different system calls, respectively. Our initial MALINE implementation was ARM-based, and switching to x86 yielded roughly the same classification results in preliminary experiments, while it greatly improved performance.

**Randomness in MALINE.** We used only one seed value for Monkey's pseudo-random number generator; it is possible the outcome of our experiments would have been different if another seed value was used. However, as the seed value has to be used consistently within an experiment consisting of thousands of applications, it is highly unlikely the difference would be significant.

## 7. CONCLUSIONS AND FUTURE WORK

We performed a preliminary feature selection exploration, but were not successful in obtaining consistent results. The reason could be a high dimensionality of the classification problem (15,000 features for our dependency representation) or a strong correlation between features. We left a more extensive feature selection exploration for future work. We also want to explore combining several learning techniques to investigate ensemble learning. We already do use a form of ensemble learning in the random forest algorithm, but we are planning to look at combinations of various algorithms too.

In this paper, we proposed a free and open-source reproducible research environment MALINE for dynamic-analysis-based detection of malware in Android. We performed an extensive empirical evaluation of our novel system call encoding into a feature vector representation against a well-known frequency representation across several dimensions. The novel encoding showed better quality than the frequency representation. Our evaluation provides numerous insights into the structure of application executions, the impact of different machine learning techniques, and the type and size of inputs to dynamic analyses, serving as a guidance for future research.

## Acknowledgments

We thank Raimondas Sasnauskas for helping us with the Android platform, and Geof Sawaya for providing feedback on an early draft of this paper. Thanks to Andreas Zeller and his group for providing us with a tool for bulk application downloading from Google Play, and the Flux Research Group of the University of Utah for the Emulab infrastructure. This work was supported in part by NSF CCF 1421678.



## 8. REFERENCES

- [1] Y. Aafer, W. Du, and H. Yin. DroidAPIMiner: Mining API-level features for robust malware detection in Android. In *SecureComm*, 2013.
- [2] UI/application exerciser monkey.
- [3] D. Arp, M. Spreitzenbarth, M. Huebner, H. Gascon, and K. Rieck. Drebin: Effective and explainable detection of android malware in your pocket. In *NDSS*, 2014.
- [4] S. Arzt, S. Rasthofer, C. Fritz, E. Bodden, A. Bartel, J. Klein, Y. Le Traon, D. Octeau, and P. McDaniel. FlowDroid: Precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for Android apps. In *PLDI*, 2014.
- [5] F. Bellard. QEMU, a fast and portable dynamic translator. In *USENIX Annual Technical Conference*, 2005.
- [6] I. Burguera, U. Zurutuza, and S. Nadjm-Tehrani. Crowdroid: behavior-based malware detection system for Android. In *SPSM*, 2011.
- [7] C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3), 2011.
- [8] X. Chen, J. Andersen, Z. Mao, M. Bailey, and J. Nazario. Towards an understanding of anti-virtualization and anti-debugging behavior in modern malware. In *DSN*, 2008.
- [9] M. Dimjašević, S. Atzeni, I. Ugrina, and Z. Rakamarić. Android malware detection based on system calls. Technical Report UUCS-15-003, University of Utah, 2015.
- [10] G. Dini, F. Martinelli, A. Saracino, and D. Sgandurra. MADAM: A multi-level anomaly detector for Android malware. In *CNS*, 2012.
- [11] W. Enck, P. Gilbert, B.-G. Chun, L. P. Cox, J. Jung, P. McDaniel, and A. N. Sheth. Taintdroid: An information-flow tracking system for realtime privacy monitoring on smartphones. In *OSDI*, 2010.
- [12] F-droid, free and open source android app repository.
- [13] A. P. Felt, E. Chin, S. Hanna, D. Song, and D. Wagner. Android permissions demystified. In *CCS*, 2011.
- [14] A. P. Felt, E. Ha, S. Egelman, A. Haney, E. Chin, and D. Wagner. Android permissions: User attention, comprehension, and behavior. In *SOUPS*, 2012.
- [15] Y. Feng, S. Anand, I. Dillig, and A. Aiken. Apposcopy: Semantics-based detection of android malware through static analysis. In *FSE*, 2014.
- [16] M. Fredrikson, S. Jha, M. Christodorescu, R. Sailer, and X. Yan. Synthesizing near-optimal malware specifications from suspicious behaviors. In *SP*, 2010.
- [17] H. Gascon, F. Yamaguchi, D. Arp, and K. Rieck. Structural detection of Android malware using embedded call graphs. In *AISec*, 2013.
- [18] A. Gorla, I. Tavecchia, F. Gross, and A. Zeller. Checking app behavior against app descriptions. In *ICSE*, 2014.
- [19] M. Grace, Y. Zhou, Q. Zhang, S. Zou, and X. Jiang. RiskRanker: scalable and accurate zero-day Android malware detection. In *MobiSys*, 2012.
- [20] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, 2nd edition, 2009.
- [21] R. K. Jidigam, T. H. Austin, and M. Stamp. Singular value decomposition and metamorphic detection. *Journal of Computer Virology and Hacking Techniques*, 11(4), 2014.
- [22] Y. Jing, Z. Zhao, G.-J. Ahn, and H. Hu. Morpheus: automatically generating heuristics to detect Android emulators. In *ACSAC*, 2014.
- [23] Third annual mobile threats report: March 2012 through March 2013. Juniper Networks Mobile Threat Center.
- [24] C. Kolbitsch, P. M. Comparetti, C. Kruegel, E. Kirda, X.-y. Zhou, and X. Wang. Effective and efficient malware detection at the end host. In *USENIX Security*, 2009.
- [25] M. Kuhn and K. Johnson. *Applied Predictive Modeling*. Springer, 2013.
- [26] A. Lanzi, D. Balzarotti, C. Kruegel, M. Christodorescu, and E. Kirda. AccessMiner: Using System-centric Models for Malware Protection. In *CCS*, 2010.
- [27] S. Liang, A. W. Keep, M. Might, S. Lyde, T. Gilray, P. Aldous, and D. Van Horn. Sound and precise malware analysis for Android via pushdown reachability and entry-point saturation. In *SPSM*, 2013.
- [28] S. Liang, W. Sun, and M. Might. Fast flow analysis with Gödel hashes. In *SCAM*, 2014.
- [29] M. Ojala and G. C. Garriga. Permutation tests for studying classifier performance. *Journal of Machine Learning Research*, 11, 2010.
- [30] S. Palahan, D. Babić, S. Chaudhuri, and D. Kifer. Extraction of statistically significant malware behaviors. In *ACSAC*, 2013.
- [31] R. Paleari, L. Martignoni, G. F. Roglia, and D. Bruschi. A fistful of red-pills: How to automatically generate procedures to detect CPU emulators. In *USENIX WOOT*, 2009.
- [32] The R project for statistical computing.
- [33] K. Rieck, P. Trinius, C. Willems, and T. Holz. Automatic analysis of malware behavior using machine learning. *Journal of Computer Security*, 19(4), 2011.
- [34] RiskIQ's report on malicious mobile apps, 2014.
- [35] T. Singh, F. D. Troia, V. A. Corrado, T. H. Austin, and M. Stamp. Support vector machines and malware detection. *Journal of Computer Virology and Hacking Techniques*, 2015.
- [36] F. Wei, S. Roy, X. Ou, and Robby. Amandroid: A precise and general inter-component data flow analysis framework for security vetting of Android apps. In *CCS*, 2014.
- [37] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. An integrated experimental environment for distributed systems and networks. *SIGOPS Oper. Syst. Rev.*, 36(SI), 2002.
- [38] E. William, D. Octeau, P. McDaniel, and S. Chaudhuri. A study of Android application security. In *USENIX Security*, 2011.
- [39] L. K. Yan and H. Yin. Droidscope: Seamlessly reconstructing the OS and Dalvik semantic views for dynamic Android malware analysis. In *USENIX Security*, 2012.
- [40] W. Yang, X. Xiao, B. Andow, S. Li, T. Xie, and E. William. Appcontext: Differentiating malicious and benign mobile app behavior under contexts. In *ICSE*, 2015.
- [41] Y. Zhou and X. Jiang. Dissecting Android malware: Characterization and evolution. In *SP*, 2012.
- [42] Y. Zhou, Z. Wang, W. Zhou, and X. Jiang. Hey, you, get off of my market: Detecting malicious apps in official and alternative Android markets. In *NDSS*, 2012.
- [43] D. Y. Zhu, J. Jung, D. Song, T. Kohno, and D. Wetherall. TaintEraser: Protecting sensitive data leaks using application-level taint tracking. *SIGOPS Oper. Syst. Rev.*, 45(1), 2011.