

Mate! Are You Really Aware? An Explainability-Guided Testing Framework for Robustness of Malware Detectors

Ruoxi Sun
CSIRO's Data61
Australia

Minhui Xue
CSIRO's Data61
Cybersecurity CRC
Australia

Gareth Tyson
Hong Kong University of Science and
Technology (GZ)
China

Tian Dong
Shanghai Jiao Tong University
China

Shaofeng Li
Peng Cheng Laboratory
China

Shuo Wang
CSIRO's Data61
Cybersecurity CRC
Australia

Haojin Zhu
Shanghai Jiao Tong University
China

Seyit Camtepe
CSIRO's Data61
Cybersecurity CRC
Australia

Surya Nepal
CSIRO's Data61
Cybersecurity CRC
Australia

ABSTRACT

Numerous open-source and commercial malware detectors are available. However, their efficacy is threatened by new adversarial attacks, whereby malware attempts to evade detection, *e.g.*, by performing feature-space manipulation. In this work, we propose an explainability-guided and model-agnostic testing framework for robustness of malware detectors when confronted with adversarial attacks. The framework introduces the concept of *Accrued Malicious Magnitude (AMM)* to identify which malware features could be manipulated to maximize the likelihood of evading detection. We then use this framework to test several state-of-the-art malware detectors' ability to detect manipulated malware. We find that (i) commercial antivirus engines are vulnerable to AMM-guided test cases; (ii) the ability of a manipulated malware generated using one detector to evade detection by another detector (*i.e.*, transferability) depends on the overlap of features with large AMM values between the different detectors; and (iii) AMM values effectively measure the fragility of features (*i.e.*, capability of feature-space manipulation to flip the prediction results) and explain the robustness of malware detectors facing evasion attacks. Our findings shed light on the limitations of current malware detectors, as well as how they can be improved.

CCS CONCEPTS

• Security and privacy → Software and application security.

KEYWORDS

Malware detectors, Explainability, Robustness

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.
ESEC/FSE '23, December 3–9, 2023, San Francisco, CA, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0327-0/23/12...\$15.00
<https://doi.org/10.1145/3611643.3616309>

ACM Reference Format:

Ruoxi Sun, Minhui Xue, Gareth Tyson, Tian Dong, Shaofeng Li, Shuo Wang, Haojin Zhu, Seyit Camtepe, and Surya Nepal. 2023. Mate! Are You Really Aware? An Explainability-Guided Testing Framework for Robustness of Malware Detectors. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '23)*, December 3–9, 2023, San Francisco, CA, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3611643.3616309>

1 INTRODUCTION

The anti-malware market is at the forefront of cybersecurity innovation, constantly driving anti-malware vendors to update their solutions to protect users against a wide range of malicious software variants. The global antivirus software market is expected to reach more than 4 billion USD in 2025 [18]. Despite this, recent research has shown that the total number of malware infections has continued to rise [4], hitting a new high during the COVID-19 pandemic [13, 66]. There are also numerous high-profile cases of zero-day attacks are being used in offensive campaigns. For example, just a few hours before the Russian-Ukraine conflict, Microsoft's Threat Intelligence Center identified a never-before-seen malware, "FoxBlade", that targeted Ukraine's government ministries and financial institutions [19].

This trend suggests that traditional signature-based and behavior-based methods cannot keep up with the rampant growth of novel malware. Hence, commercial antivirus companies have started using machine learning [1, 68] to enable detection without the need for signatures. However, it has been shown that attackers can evade machine learning-based detectors by manipulating the features that such malware detectors use [36, 37, 56, 58, 74]. Because of this, commercial antivirus systems are susceptible to adversarial attacks [63]. Although there have been numerous works [28, 30, 39, 41, 54] looking at adversarial attacks in computer vision (where adversaries change specific pixels), adversarial attacks on malware are far less understood. We therefore leverage this adversarial approach to build a testing framework to evaluate malware detectors. However, we find that it is difficult to fully identify the root causes that impact the decisions of malware detectors. We therefore reduce this causal

discovery problem to identifying the explainable factors that are measurable to impact the robustness of malware detectors.

To implement the testing framework for malware detectors, it is necessary to build techniques that can (i) generate adversarial malware variants; and (ii) measure the explainable features that drive the malware detector’s ultimate decision (benign or malicious). One existing approach for generating adversarial malware test samples is *obfuscation*. This focuses on changing the semantic meanings of code snippets in the problem-space (i.e., source code), and further obfuscating the malicious signatures or patterns, including hiding the control flow, inserting dummy code, and manipulating variable names [8, 11, 14, 17, 23, 47], thereby fooling rule-based malware detectors. However, defenses against obfuscation are well-researched. For example, recent research [25, 26, 32] looks into malware behavior distillation or program behavioral variability analysis towards training robust malware classifiers. Furthermore, Pierazzi *et al.* [58] argue that the use of mass obfuscation may be counterproductive, rendering antivirus companies to be on the alert [70, 71].

In this research, we aim to test the robustness of malware detectors against a new type of adversarial generation: *feature-space manipulation*. This involves manipulating the malware to trigger changes in the feature space used by detectors. Such threats are becoming more prominent because machine learning-based detectors have reduced the efficacy of problem-space attacks. Feature-space manipulation aims to introduce the intended changes to the feature space (i.e., extracted feature vectors) by precisely modifying the problem-space (i.e., code). These manipulating actions could be, for example, adding a redundant section (e.g., adding a new code section without linking its address in the section table) or injecting dead code that is unreachable (e.g., adding a file I/O request under an always-false condition, so that the dummy code will never be executed). Similar techniques have been implemented by Demetrio *et al.* [34] in a black-box optimization of adversarial Windows malware. Although they still focus on problem-space, it is possible to apply such techniques in a feature-space manipulation. Due to its urgency, we argue that building a comprehensive testing framework that can automate the identification of limitations in malware detectors is vital.

Several challenges need to be solved during the establishment of a testing framework for robustness of malware detectors. The *first challenge* is how we can mimic a random attacker in the real world, who may have limited capacity and knowledge. This means that we are limited to off-the-shelf and easy-to-obtain techniques in the testing. The *second challenge* is, since most commercial malware detectors are not open-source, this must be done in a detector-agnostic manner. The *third challenge* is how to interpretively understand the testing results (e.g., why test cases work across different detectors). With the above challenges in mind, to ensure the reproducibility and coverage of our framework, we have several criteria on the development of our testing framework: (i) *Easy-to-obtain*, we only utilize open-source and off-the-shelf tools or techniques; (ii) *Model-agnostic*, we decouple the testing strategy from the specifics of the detector; and (iii) *Explainable*, an explainable approach is proposed, which will help us to *explain* the root cause of test cases’ transferability and identify potential weaknesses in malware detectors.

In this paper, we first propose *Accrued Malicious Magnitude (AMM)* to guide feature space manipulation, finding the correct

action(s) on the problem-space that will influence the feature values but without changing run-time functionality. We then evaluate the robustness of state-of-the-art malware detectors with generated test cases. To establish such a strategy, we generate test cases by perturbing the feature space and converting the manipulations back into the problem space. To achieve our goal, the problem is split into two sub-problems: (i) generating test cases; and (ii) testing malware detectors. Our research helps security researchers as well as vendors who develop anti-malware solutions to better understand the robustness of malware detectors and provide insights into how to improve malware defense strategies in an interpretive manner. The main contributions of this paper are four-fold:

- We propose an explainability-guided and model-agnostic testing framework for malware detectors (§4). Our framework generates test cases while preserving the malicious functions of the malware. We introduce the concept of *Accrued Malicious Magnitude (AMM)* to guide the feature selection approach for feature-space manipulation. We further project the manipulated features back into the problem space with a binary builder.
- We use AMM to test the robustness of state-of-the-art malware detectors (§5). We show that commercial antivirus engines are vulnerable to AMM-based test cases (§6.2). Experimental results indicate that feature-space manipulated test cases have significant evasion capability, which decreases the detection rates of 8 state-of-the-art Android malware detectors by 91.75% on average, and bypasses an average of 37.35 (62.25%) antivirus engines in VirusTotal [16]. We also highlight the generalizability of our AMM approach by applying it to WinPE malware detectors (§6.4).
- We explain how manipulations trained on one detector can work on another detector (i.e., transferability) through our explainability-guided approach (§6.3), indicating that this transferability relies on the overlap features that have large AMM values between different machine learning models.
- We further investigate an approach to improve machine learning-based detectors through excluding *high sensitive but less important* features during training (§6.5). Results show that AMM values can effectively measure the capability of features of flipping classification results. We suggest that machine learning-based anti-virus products should consider using the AMM values to improve their robustness.

To the best of our knowledge, this is the *first* paper to systematically test the robustness of malware detectors in a way that combines feature-space manipulations with semantic explainability.

2 PRELIMINARIES

In this section, we introduce a motivating example and related research.

2.1 Motivating Example

To motivate the need for our testing framework, we analyze the source code from an example Android malware, which is tagged as malicious by 39/60 detectors from VirusTotal (VT) [16]. From the source, we find a snippet of malicious code shown in the top part of Figure 1. As shown in lines 3 and 4, the malware executes a native scripts via root permissions by `su -c ./script1` command.

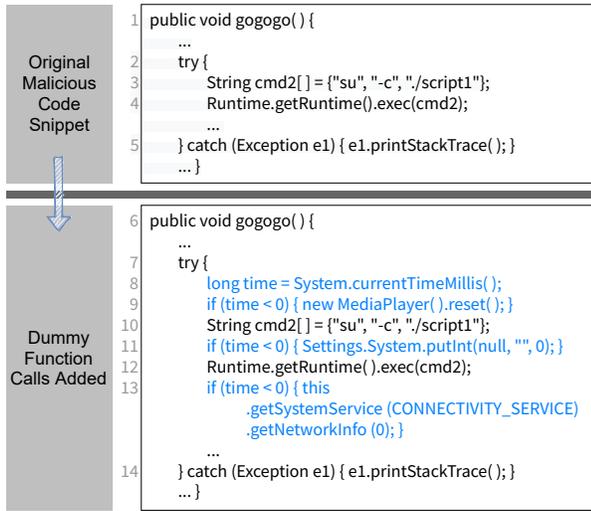


Figure 1: An intuitive example of injecting dummy ‘benign’ function calls into malicious malware, which bypasses 6 antivirus engines in VirusTotal and evades Drebin.

Static features are commonly used in malware detectors where the results is determined through pattern recognition, weighted algorithms, or signature matching. Thus, it is possible to mislead detectors by introducing more ‘benign-oriented’ elements into malware. Hence, we insert several ‘benign’ function calls with always-false condition closure (e.g., $time < 0$) to ensure they are unreachable during run-time, preserving the original (malicious) functionality. After rebuilding the source code, the modified binary is identified by 33 scanners – 6 fewer than originally, and it bypasses the machine-learning detector provided by Drebin [24]. The remainder of this paper develops an explainability-guided testing framework and problem-space rebuilding tool that can automate this intuitive idea for the testing of detector robustness.

2.2 Related Work

Malware detectors. Many modern antivirus engines utilize rule-based analysis, such as signature matching, static unpacking, heuristics matching, and emulation techniques [6, 50]. However, rule-based antivirus engines rely heavily on expert knowledge. With the advantage of feature extraction derived from machine learning techniques, there has been a flurry of work that integrates machine learning models into malware detectors [22, 24, 40, 49, 50, 57, 72, 73]. We focus our evaluation on detectors that use static features due to their prevalence in providing pre-execution detection and prevention for many commercial endpoint protection solutions, such as Kaspersky [7], Avast [9], and ESET [5].

Evaluation of malware detectors. A few studies [56, 59] have explored the effect of obfuscations on anti-malware products, utilizing off-the-shelf tools. Hammad *et al.* [44] conducted a large-scale empirical study that evaluates the effectiveness of the top anti-malware products, including 7 open-source, academic, and commercial obfuscation tools. Several studies [31, 62, 65] have evaluated machine learning-based malware classifier models with the adversarial samples generated by generative adversarial networks

(GANs) or automated poisoning attacks. Recent research [26, 52] proposed methods to cope with concept drift or dataset shift, which may lead to performance degradation of malware detectors. Compared to our research, the scope of these studies only covers either the rule-based products or the machine learning-based models in isolation (rather than both).

Adversarial samples against malware detectors. The goal of the adversarial attacks is to generate a small perturbation for a given malware sample that results in it being misclassified. This type of attack has been extensively explored in computer vision, and previous research efforts have also investigated the applicability of such techniques to malware classification. Xu *et al.* [74] proposed a genetic programming-based approach to perform a directed search for evasive variants for PDF malware. Demetrio *et al.* [35] demonstrated that genetic programming based adversarial attacks are applicable to portable executable (PE) malware classifiers. Two recent works [21, 64] also apply deep reinforcement learning to generate adversarial samples for PE malware to bypass machine learning models. Compared to our research, these studies focus on proposing adversarial attacks rather than testing the robustness of malware detectors and further explaining it.

3 THREAT MODEL & PROBLEM DEFINITION

In this section, we define the threat model we use in testings, and present our problem definition.

3.1 Threat Model

Our testing framework generates adversarial samples to test the robustness of malware detectors. For this, we must define our assumed threat model. We follow the methodology by Carlini *et al.* [29] and describe this threat model with the adversary’s goals, capabilities, and knowledge.

Adversary’s goal. The adversary’s goal is to manipulate malware samples to evade the detection of malware detectors, including both white-box and black-box detectors. In the testing, we only use binary detectors which determine if the software under test is benign or malicious. Thus, the goal of attackers is to cause the malicious samples to be misclassified as benign.

Adversary’s capability and knowledge. In this work, we assume that an attacker has full knowledge of *one* machine learning model that has been trained for malware detection, including its architecture and training dataset. This is reasonable as many machine learning models, such as LGBM [48], are open-source. Such a white-box model will be used as the source of adversarial sample generation. With this white-box model, the attacker is capable of ranking the contribution of features and manipulating a malicious sample accordingly to influence its representation in the feature space. For other black-box detectors under-test, including machine learning classifiers and antivirus engines, the adversary has no knowledge about the detectors’ training dataset, inner structure, or detection mechanism. For instance, the adversary cannot inject poisoned data into the training dataset or manipulate any code of detectors. However, they will still have some basic knowledge about machine learning-based detectors, including access to open-source datasets which could be part of the actual training

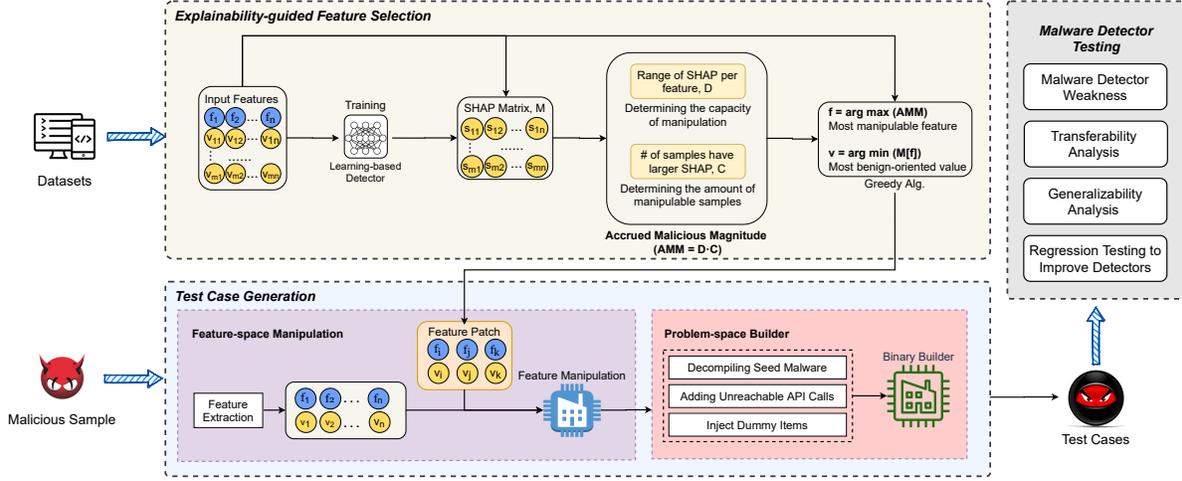


Figure 2: The overview of our testing framework.

dataset, generic or popular feature extraction methods [61], and off-the-shelf machine-learning detectors.

3.2 Problem Definition

Our goal is to test the robustness of malware detectors using the generated test cases. Consider a malware detector mapping a piece of software $x \in X$ to a classification label $l \in \{0, 1\}$ (where 0 represents benign and 1 represents malicious). The adversary is trying to mislead this prediction. Our key test metric is the detection rate on test cases:

$$d(x_m) = 1, x_t = \text{Gen}(x_m), d(x_t) = 0, \quad (1)$$

$$\text{DetectionRate}(d) = \frac{1}{n} \sum_{i=1}^n d(x_t^i)$$

where d is the malware detector which could be either a machine learning detector (a feature extraction method plus a trained model) or an antivirus engine. x_m is the original malware sample, and Gen is the generator of test case, x_t , while keeping its malware functionality the same as x_m . $\text{DetectionRate}(d)$ is defined as a detector's detection rate on n test cases, indicating the robustness of d .

4 TESTING FRAMEWORK

Our testing framework consists of three key components (see Figure 2): (i) explainability-guided feature selection, to select the features for manipulation; (ii) a test case generator that relies on the previously selected features; and (iii) malware detector testing to identify which detectors are robust against the adversarial malware samples. Before diving into the testing framework, we would like to introduce preliminary knowledge about SHAP [53], the model explanation technique we used in our testing.

4.1 A Primer on SHAP

Research into explainable machine learning has proposed multiple systems to interpret the predictions of complex models. We rely on SHAP [53] (based on the coalitional game theory concept of Shapley values) as prerequisite knowledge to bootstrap the testing framework. The SHAP framework subsumes several earlier

model explanation techniques together, including LIME [60] and Integrated Gradients [67]. SHAP has the objective of explaining the final value of a prediction by attributing a value to each feature based on its contribution to the final result. To accomplish this task, the SHAP frameworks train a surrogate linear explanation model g of the form:

$$f(x) = g(x'),$$

$$g(x') = \phi_0 + \sum_{j=1}^M \phi_j x'_j, \quad (2)$$

where f is the original model, x is the input sample to be attributed, x' is the coalition vector of x . $\phi_0 = \mathbb{E}_X(f(X))$ is the average prediction of the original model on sampled dataset X . The Shapley value $\phi_j \in \mathbb{R}$ is the feature attribution for the j^{th} feature x'_j to the model's decision. Summing the effects of all feature attributions approximates the difference of prediction for x and the average of the original model, aiming to explain any machine learning-based model without internal knowledge.

LIME uses a linear explanation model $g(x')$ to locally approximate the original model, where locality is measured in the simplified binary input space, i.e., $x' \in \{0, 1\}^M$. To find ϕ , LIME minimizes the following objective function:

$$\xi = \arg \min_{g \in G} L(f, g, \pi_x) + \Omega(g), \quad (3)$$

where L is the squared loss over a set of samples in the simplified input space weighted by the kernel function π_x , and Ω penalizes the complexity of $g \in G$ where G is hypothesis space. Therefore, based on the input feature vectors and the output predictions of the model, in this research, we use SHAP to approximate the fragility of each feature, selecting features that can be manipulated to flip the prediction.

4.2 Step 1: Feature Selection

In the first step of our methodology, we utilize SHAP to create explainability-guided adversarial test cases. In this step we use a single detector model to generate SHAP values for the input dataset,

Table 1: Feature extraction methods used in different Android malware detectors and in our measurement.

Features	Drebin [24]	MaMaDroid [57]	RevealDroid [40]	DroidSpan [27]	Our Measurement
Permissions	●				●
Hardware Components	●				
App Components	●				
API Calls	●	○	●	●	●
Strings (e.g., Network Addresses)	●				●
Call Graphs		●		●	
Native Call			●		

●: the feature is involved; ○: the feature is indirectly involved.

calculating how much one feature contributes to an individual prediction. It is assumed that an adversary would have access to this single model. The workflow of the explainability-guided feature selection is illustrated in Algorithm 1. For a set of seed malware, X_s , we generate a corresponding test case set, X_t .

Pre-processing. We extract features from the training samples X of a trained machine learning model m (line 2). Here we will use a generic feature extraction approach that is adopted from open-sourced detectors to mimic the adversary’s capability and knowledge. Using Android as an example, Table 1 summarizes the feature extraction methods of several state-of-the-art for Android malware detection. We adopt features that are representative and easy to extract. It is possible to generate stronger evasive adversarial test cases with more features involved. However, our goal is to evaluate malware detectors with an easy-to-obtain approach, which should provide a lower bound of robustness. Note, we use a similar strategy to determine the feature extraction methods for other operating systems, as described in §6.4.

Then the vectorized samples X' and the model are input to `shap()` to calculate the SHAP value matrix M (line 3). The matrix is then used to select the most evasive features and the most benign-oriented values (*i.e.*, the most negative values that exist in the selected features, as 0 represents benign).

Feature selection. To select the feature that has largest malicious magnitude, we propose the concept of *Accrued Malicious Magnitude (AMM)*. The AMM is defined as the product of the magnitude of SHAP values in each feature and the number of samples that have malicious-oriented values (*i.e.*, values towards the positive side, as 1 represents malicious) in the corresponding feature. By calculating AMM values, we select the feature that has the largest modifiable capability and has the most samples to be modified as the test cases, *i.e.*, samples that have SHAP values towards the positive (malicious) side, which have the potential to be manipulated to benign. Specifically, starting from the `getRange(M)` in line 5, we first calculate the range of SHAP values in each feature and store the results in a one-dimension vector D . D indicates the potential magnitude we can modify on each feature, *i.e.*, each $d_i \in D$ presents the difference between the maximum SHAP value and the minimum SHAP value of feature f_i . Next, for each feature, we count how many samples have a SHAP value larger than the mean SHAP value of that feature (the `countLarge(M)` in line 6) and collect the results in a one-dimensional vector C . Therefore, a larger $c_i \in C$ means that, for feature f_i , there are more samples that have a SHAP value towards malicious, such that more samples can be manipulated towards benign. Therefore, we select the most evasive feature according to the AMM values, denoting the dot product of the range of SHAP

Algorithm 1: AMM-based Feature-Space Selection

Input: Machine learning model m , dataset X , and the number of features to be selected N .

Output: Feature patch P .

```

1  $P = \text{map}(\text{Feature}, \text{Value});$ 
2  $X' \leftarrow \text{vectorize}(X);$ 
3  $M \leftarrow \text{shap}(X', m);$ 
4 while  $\text{size}(P) < N$  do
5    $D \leftarrow \text{getRange}(M);$ 
6    $C \leftarrow \text{countLarge}(M);$ 
7    $AMM \leftarrow D \cdot C;$ 
8    $f \leftarrow \arg \max(AMM);$ 
9    $v \leftarrow \arg \min(M[f]);$ 
10  if isManipulatable( $f$ ) then
11     $P \leftarrow P \cup (f, v);$ 
12    for each  $x' \in X'$  do
13      if  $x'[f] \neq v$  then
14         $idx \leftarrow \text{getIndex}(X', x');$ 
15         $M \leftarrow M \setminus M[idx];$ 
16         $X' \leftarrow X' \setminus x';$ 
17 return  $P;$ 

```

values (D) and the number of SHAP values greater than mean (C) (line 7).

Value selection. Once we have identified the feature f to compromise, the next step is to choose the value for the selected feature to guide the manipulation. We select the most benign-oriented value, v , in the feature space. This corresponds to the most **negative** value in $M[f]$, the SHAP value of the feature f (line 9).

Updating the feature. After obtaining (f, v) , if the selected feature f is manipulable, we add the pair into a map, P , as the *Feature Patch* to be used in the feature-space manipulation (line 11). Note that, due to the strong semantic restrictions of the binaries, we cannot simply choose any arbitrary pairs of feature and values for the test manipulation. Instead, we restrict the feature-space manipulation to only features and values that are independent (IID) and can be modified with original functionalities preserved. For example, consider the feature that counts the size of a binary: When we modify the value of another feature, the former will be modified indirectly. Therefore, the features and values we select to be manipulated follow two principles employed by the previous literature [42, 43, 61]. These principles are: (i) features are manipulable in the original problem space; and (ii) selected features have no dependencies or cannot be affected by other features. We described manipulable features for Android and WinPE datasets later in §5.2 and §6.4.

Greedy strategy. After obtaining feature-value pairs, we conduct a greedy strategy, removing samples that have the same value, v , for feature f from the dataset (lines 13 to 17). We do this to make sure that the same feature-value pair will not be selected again. The procedure repeats until we find N feature-value pairs. These N pairs are then used as feature patch in the next stage to generate test cases.

4.3 Step 2: Test Case Generator

In the next step, the test case generator conducts feature manipulation according to the feature patch obtained from prior steps.

Feature-space manipulation. Equation 4 summarizes the feature-space manipulation:

$$\begin{aligned} x'_m &= \text{vectorize}(x_m), \\ x'_t &= \text{manipulateFeature}(x'_m, P), \\ x_t &= \text{Gen}(x_m, x'_t), \end{aligned} \quad (4)$$

where x'_m is the result of applying feature extraction on a malicious sample x_m using `vectorize()`. `manipulateFeature()` manipulates the sample in feature-space guided by the selected feature and value pairs, P . Note that, the sample generator `Gen()` will take the manipulated feature-space sample x'_t and the original seed sample x_m as input, and implement the changes in feature-space back to problem-space to generate the test case, while keeping its malware functionality (as detailed in Binary builder below).

Binary builder. To ensure that no loss of functionality is inadvertently introduced as a side effect of feature manipulation, we only apply these changes to unreachable areas of binaries, so that these changes will never be executed during run-time. Therefore, we guarantee that test cases are executable and can be applied in the testing of malware detectors in the wild. Then, we apply these changes on seed binaries with the help of open-source binary builders.

For simplicity, we present our tooling for Android malware, yet we emphasize that our framework works with other operating systems (see §6.4). We adopt a similar feature extraction method of Drebin on the Android APK. Since features are a vector of boolean values representing the existence of a feature, the feature value could only be modified from 0 (absence) to 1 (presence) to preserve original functionalities. We first leverage Apktool [3] to decompile an APK file into Smali [12] code, a structured assembly language. API calls and network URLs are transformed to smali instruction code, which is wrapped by an unreachable disclosure, e.g., an always-false condition closure such as `if(time < 0)`. The Smali code is then inserted into the Smali file of the main activity. Features representing Android manifest components are inserted into `AndroidManifest.xml` file directly. Finally, we utilize Apktool to assemble all decompiled and manipulated files into an adversarial APK sample. If a feature manipulation cannot be implemented in this way, we skip it and continue with the next feature in the selected patch.

4.4 Step 3: Malware Detector Testing

After the test cases are generated, our framework programmatically executes a series of tests on the malware detectors. This involves (i) calculating the per-detector robustness; and (ii) testing the transferability of attacks across detectors. We conclude this subsection by proposing techniques that can improve detector performance.

Testing detector robustness. The methodology to test detector robustness is straightforward. We first input the seed malware into each detector and collect the detection rate on the unaltered malware dataset. Next, the test cases generated from the white-box model will be input to the detectors (including the white-box model itself). We then compare the difference between the detection rate of the seed vs. test cases. This allows us to test which detectors are vulnerable to test cases, i.e., whether the manipulation on AMM features changes the detector results.

Transferability analysis. Considering that machine learning-based detectors may use similar feature extraction methods, it is possible that multiple detectors are susceptible to the same feature manipulations. Specifically, we posit the test cases generated from one machine learning model are likely to be effective on other models that are trained on the same data distribution, due to the similarity of decision boundaries. Therefore, the test framework also calculates the ability to transfer an evasion trained on one detector to another, i.e., transferability. This is important as, the more powerful the transferability a malware test case has, the more effective it is against other detectors. To evaluate the transferability of test cases generated by the AMM-based approach, we generate cases from white-box models and apply them to black-box models. If transferability exists, the test framework calculates the feature-space overlaps among models to explore the root cause of transferability. This can be used by developers to understand the weaknesses in their feature engineering.

Regression testing for detector improvement. Inspired by recent research [46], a detection model can be improved by removing important features from the training phase, where the important features refer to the ones that are highly contributing to the model prediction accuracy. Shapley Additive Global importance (SAGE) [33] is a framework that measures how much a feature contributes to the prediction accuracy of a model. We apply the improvement on the detection models with SAGE and test their robustness against adversarial samples. Specifically, we first calculate SAGE values of each feature. Since the most important features are the ones with largest SAGE values, we sort the features by SAGE values in a descending order and select the top features. Then we remove the top features from samples and generate a new training set. Finally, an improved model m_i is trained with the new training set. To establish a thorough comparison, we train another model m_a that excludes top AMM-based features to compare and explain the improvement of malware detectors.

5 TESTING FRAMEWORK SETUP

In this section, we describe the setup of our tests, including the detectors under test, the datasets, and how we trained the models.

5.1 Detectors under Test

To showcase our testing framework, we experiment with a number of malware detectors. For Android malware detectors, we test the robustness of 8 state-of-the-art machine learning-based detectors (a combination of 2 feature extraction methods and 4 machine learning models); for WinPE detectors, we test 4 detectors (1 feature extraction method accompany 4 models). We also test 60 antivirus engines available using VirusTotal [16]. The detectors under test are listed in Table 2. Specifically, we involve Drebin [24], MaMaDroid [57], and Ember [22] in our testing, because (i) they proposed unique features extraction methods; (ii) high accuracy rates are reported on large datasets; and (iii) their source code and datasets have been made publicly available.

To follow the conventions of prior studies [24, 27, 40, 57, 61], we select 4 off-the-shelf machine learning-based malware models that are commonly used in malware detectors. In our threat model, the adversary has full knowledge of one machine learning-based

Table 2: Detectors under test.

Type	Name	Description
Feature Extraction Methods	Drebin [24]	A lightweight method for Android malware detection, extracting 8 sets of features from an application’s code and manifest.
	MaMa-Droid [57]	A static-analysis-based system that abstracts app’s API calls and builds a model from the call graph of an app as Markov chains.
	Ember [22]	A static WinPE malware classifier extracting eight groups of raw features that include both parsed features and format-agnostic histograms and counts of strings.
Machine Learning Models	LGBM [48]	LightGBM, an open-sourced gradient boosting framework, based on the decision tree algorithm.
	SVM	Support Vector Machine, a supervised learning method based on statistical learning frameworks.
	RF	Random Forests, an ensemble learning method that combines decision trees to provide classification.
	DNN	A feed-forward neural network with with one input layer and three fully-connected hidden layers (the last one ends with a Softmax function).
Antivirus Engines	VT [16]	VirusTotal, a website that aggregates more than 70 antivirus products and online scan engines, allowing a user to check for viruses that the user’s own antivirus software may have missed.

model. Without loss of generality, we set LGBM as this *white-box* model. Note, any machine learning model could serve this role as our approach is model-agnostic. For the antivirus engines, we use VirusTotal, an online service that provides over 70 antivirus scanners to detect malicious files and URLs. We find that 60 scanners are always available while the others are not stable. Therefore, we include these 60 scanners in our evaluation. All antivirus engines fall into the black-box category as the attacker has no specific knowledge about them.

5.2 Datasets and Model Training

In our experiments, we conduct testings using an Android dataset and a WinPE dataset.

The Android Application Package (APK) is the package file format used by the Android operating system for distribution of mobile apps. We use the well-studied Drebin [24] dataset, which contains 123,453 benign samples and 5,560 malicious. We also include 8,351 malware samples that have been uploaded on VirusShare [15]. Since the ratio of malicious and benign apps is unbalanced, we randomly select 5,560 benign and 5,560 malicious samples, making up 11,120 samples. Further, we create a random 50:20:30 split of samples for training, validation, and testing, to train a LightGBM model. 3,000 test cases are generated from 3,000 randomly selected malicious samples. The details about the WinPE dataset establishing and model training are provided in §6.4.

To train the ML models mentioned above, we employ Androguard [2] to extract raw features from APK samples. Androguard is a python tool to analyze and manipulate Android files. It disassembles an APK file and converts its byte code and resource files into a readable and structured format. We further extract features from the manifest file and from the Dalvik Executable (dex) file. These features are then used to train and evaluate the machine learning based detectors. All Android features are manipulable as they are independent to each other.

6 TESTING FRAMEWORK RESULTS

We next employ our explainability-guided test framework to evaluate the robustness of the detectors. Specifically, we compare the detection rates between *original* samples and the *test cases* generated by our proposed strategy. To choose the number of features to manipulate, we conduct pilot experiments on Android and WinPE datasets.

6.1 Pilot Experiment

Number of Features. To choose the number of features to manipulate in the explainability-guided test case generation, we conduct pilot experiments on the Android and WinPE datasets. Here we only report the experiment on the Android dataset as they follow the same strategy. The pilot experiment compares the detection rate of the LGBM model while manipulating 10, 25, 50, 75, 100, 125, and 150 features on 20 sets of 100 seed samples (2,000 samples in total). As shown in Figure 3, with more features manipulated, the detection rate decreases. The trend turns at 75 features, with 5.1% test cases detected. We therefore choose $N = 75$ for the APK binaries.

6.2 Testing Results on Android Malware Detectors

In our Android malware detector testing, we generate APK test cases from the LGBM model, and test on 8 machine learning-based detectors using 2 feature extraction methods and 4 models (1 in white-box and 3 in black-box), and 60 antivirus engines (black-box) from VT.

Testing machine learning detectors. All machine learning-based detectors have reasonable detection rate (above 85% on average) on the original malware samples, as shown in Figure 4. Note, the y-axis refers to the detection rates of samples and each detector is presented on the x-axis. On the Drebin-based detectors, the test cases averagely reduce the detection rate by 74.11% (ranging from 50.37% on Drebin-RF to 98.97% on Drebin-SVM). Noticeably, all test cases bypass the detection of Drebin-SVM. This may be attributable to (i) our feature extraction method being close to Drebin; and (ii) the SVM, being a simple linear model, making its prediction easier to flip. We further analyze the robustness of the SVM model as a case study in §7. While on MaMaDroid-based detectors, considering that MaMaDroid utilizes call graph features, which are indirectly related to our API call features extracted from Android source code, it is expected that fewer test cases can bypass the detection. Specifically, the detection rates of MaMaDroid-based detectors decrease 52.50% on average (ranging from 49.76% on MaMaDroid-RF to 55.35% on MaMaDroid-SVM) on test cases. The detection rates decrease down to 36.54% (avg.) which is even lower than random guess. We conclude that all 8 detectors involved in our Android malware detector testing are vulnerable to the test cases generated by the explainability-guided feature space manipulation. Note, in white-box scenarios (*i.e.*, the 2 detectors using LGBM model), it is expected that the generated test cases are more effective against white-box models. To further explore the reason why our explainability-guided approach also works in black-box scenarios, we present further test results on transferability in §6.3.

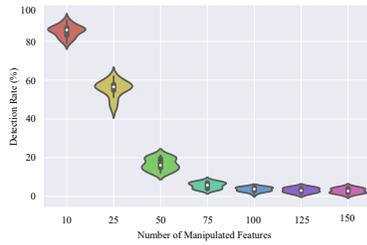


Figure 3: Detection rates of manipulating different numbers of features.

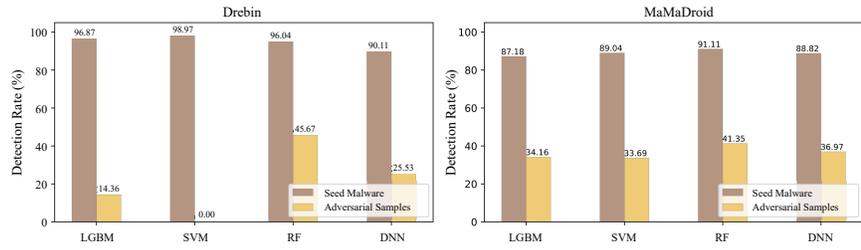


Figure 4: Testing results on machine learning-based detectors.

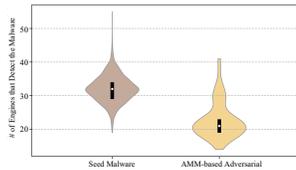


Figure 5: Number of VirusTotal anti-virus engines that can detect original malware samples and test cases.

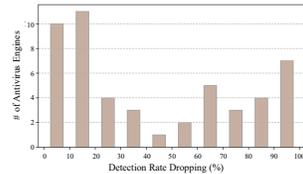


Figure 6: Histogram of antivirus engines with respect to the detection rate decreasing.

Testing VirusTotal engines. We next turn our attention to the 60 antivirus engines accessible via VirusTotal. Figure 5 plots the number of virus engines that successfully flag the seed malware and test cases as malicious. The average detection rate of seed malware is 53.40%, *i.e.*, on average 32.04 VT detectors detect a sample as malicious. In contrast, only an average of 22.65 detectors flag test cases as malicious. This lowers the detection rate of VT to 37.75% on average. This result indicates that not all detectors in VT are robust against these test cases.

To illustrate the performance of each antivirus engine and explain the overlapping results between samples, Figure 6 presents a histogram of the detection rate reductions for each antivirus engine. The x-axis represent the range of detection rates decreasing and the y-axis refers to the numbers of antivirus engines. To obtain a meaningful result, 10 engines that have detection rates less than 50% on seed malware samples are excluded. According to Figure 6, 10 engines are relatively resistant to the test cases, with detection rate reductions of under 10%. This explains the overlapping results between the upper percentile of test cases and the mean of the seed malware samples. In contrast, there are 22 engines with detection rates that decrease by over 50%, indicating that these antivirus engines are vulnerable to the AMM-based test cases.

Takeaway 1:

- Machine learning-based detectors are vulnerable to the test cases generated by the AMM-based strategy, which manipulates the most evasive features.
- Test cases can evade detection by the black-box antivirus engines in VirusTotal. This indicates that the majority of anti-malware vendors could use our testing framework to examine and improve their detectors.

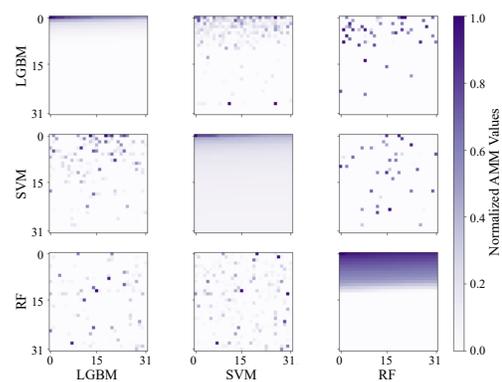


Figure 7: AMM values of top 1,024 features that have the highest AMM values in LGBM (row 1), SVM (row 2), and RF (row 3) models, separately.

6.3 Transferability Analysis

Transferability is the ability for a test case to be effective against multiple learning-based detectors. To study this, taking Drebin as an example, we next generate AMM-based test cases from each of the detectors and input them to different detectors. Here, we seek to understand how a manipulation guided by one detector performs against other detectors.¹ Specifically, we measure the transferability in two aspects: feature overlaps and detection rates.

Feature overlap. To explore the reason why the test cases can transfer across detectors, we present the top 1,024 features that have the highest AMM values in the generation models across each detector in a heatmap, shown in Figure 7. In each subplot, we present the features as 32 rows by 32 columns of dots (normalized to [0, 1]), where the darker dots represent higher values of AMM (which indicates a greater possibility to be selected as a feature to be manipulated). Further, we sort the features in the generation models according to the AMM values in descending order. Therefore, more darker dots scattered in the upper zone of the nine subplots indicate that there are more features having been selected across detectors. From the heatmap we can observe that (i) features with large AMM values in LGBM (dark dots) overlap with most of the counterparts of SVM; and (ii) many features with large AMM values in RF are out of the scope of the counterparts of LGBM. We further explain such overlaps with the result of detection rates across different models.

¹As the calculation of SHAP matrix on DNN model requires large amount of computing time using our facility, we conduct transferability analysis on 3 generation models and we believe such experimental setting is sufficient.

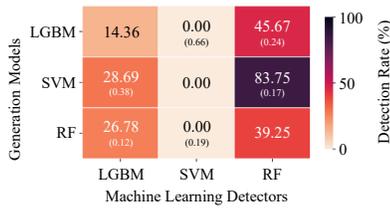


Figure 8: Detection rates of three Drebin-based detectors on test cases generated from LGBM, SVM, and RF.

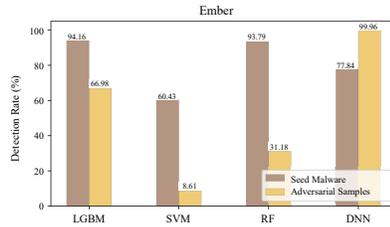


Figure 9: Detection rates of four detectors on WinPE seed malware and AMM-based test cases.

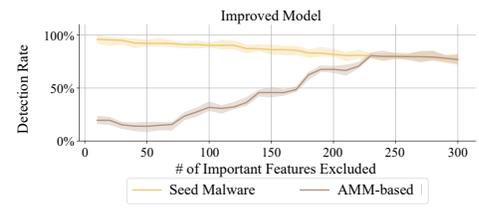


Figure 10: Detection rates of seed malware and AMM-based test cases detected by Drebin-based detectors excluding different amount of important features.

Detection rate. We further evaluate the transferability of test cases by inspecting the detection rates. Figure 8 shows the detection rates of test cases across machine learning detectors. The y-axis represents three generation models, and the x-axis shows the target Drebin-based detectors. A darker color indicates a higher detection rate (representing lower transferability).

We use the Cosine similarity (numbers in brackets) of AMM values of the top 1,024 features between models to quantify the overlaps, *i.e.*, a higher Cosine similarity value indicates a heavier overlap of high-AMM features between models. From the results in Figure 8, we observe that a higher similarity value reflects stronger transferability – for example, the transferability from LGBM to SVM outperforms the transferability from LGBM to RF – the overlaps resonate with the detection rate results. Thus, the overlaps explain why the test cases transfers across learning-based detectors. Simply put, if we manipulate enough features across different learning-based models (*i.e.*, feature overlaps), the evasion can be transferred.

Takeaway 2: The transferability of test cases depends on the overlaps of features with large Accrued Malicious Magnitude (AMM) values between different learning-based detectors.

6.4 Generalizability Analysis

To examine whether our testing framework can generalize to other operating systems, we next conduct testing on Windows Portable Executable (WinPE) files. We use SOREL-20M [45] as the WinPE dataset in our experiment. SOREL-20M is a representative public dataset of malicious and benign WinPE samples used for malware classification, consisting of 2,381-dimensional feature vectors extracted from 9,470,626 benign and 9,919,251 malicious samples, as well as corresponding malicious binaries. It leverages the feature extraction function from Ember [22]. We randomly choose 10,000 benign and 10,000 malicious samples to train detectors with LGBM, SVM, RF, and DNN models.

Feature manipulation. In a WinPE file, many of the features are derived from the same underlying structures, potentially leading to contradictory that cannot be manipulated concurrently. For instance, each addition of a writable section, the overall section count inevitably rises. Previous work [61] shows that only 17 features can be modified directly and indirectly to preserve the functionality of WinPE binaries. The pilot experiment on WinPE dataset suggests that we should choose N as 17. We leverage LIEF [69] to

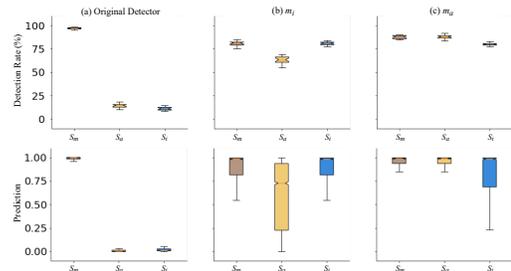


Figure 11: Detection rates and prediction values of the original model and two improved models.

extract features, and pefile [10] to apply feature manipulation on the WinPE binaries.

Testing results. Figure 9 shows the detection rates of 4 Ember-based detectors on original malware samples and test cases. The test cases have remarkable evasion performance on the LGBM, SVM, and RF detectors: On average, the detection rates are decreased to 35.59% across the three detectors. However, since most WinPE features correlate with each other, the method of parsing and generating WinPE binaries (*i.e.*, directly modifying values and adding empty sections) may negatively affect the evasion, illustrated by the DNN. In a nutshell, the test result shows that our proposed explainability-guided testing framework is also effective at identifying limitations in Windows malware detectors.

6.5 Revisiting AMM with Improved Detectors

Our evaluation shows that machine learning-based detectors are vulnerable to test cases generated by AMM. Therefore, we seek an interpretive approach to improve the robustness of existing detectors.

Improving detectors. We first follow the methodology in §4.4 and generate an improved Drebin-LGBM detector m_i by excluding SAGE-based important features in the training phase. Specifically, we first generate new machine learning detectors each with a different number of features excluded from the training set. The trend of detection rates on seed malware and AMM-based test cases are shown in Figure 10. From the result, we find that excluding 220 important features allows the improved detector to attain the highest detection rate (80.56%) on test cases. We therefore employ excluding 220 important features in the following experiment. We also generate another improved detector m_a by excluding *fragile* features. These are features that can be manipulated to flip the prediction, *i.e.*, features with high AMM values but low SAGE values, taking 0.2

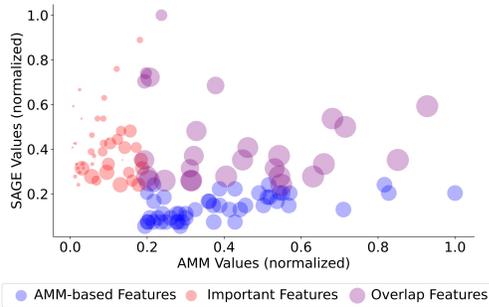


Figure 12: Distributions of evaded test cases.

as a threshold (the minimal AMM value of 75 features we selected in the pilot experiment is 0.2).

Measuring improved detectors. We further generate two groups of test cases, S_i (by SAGE) and S_a (by AMM). Meanwhile we introduce the seed samples S_m and the original Drebin-LGBM detector as a benchmark. In the experiment, we leverage 5,459 seed malware to generate test cases where we randomly select 300 samples for 20-round tests on each detector. Figure 11 presents the detection rates and prediction values. We present results for three malware datasets: seed malware set S_m , test case sets S_a and S_i . We define a sample as malicious when the prediction value is greater than or equals to 0.5. On the original detector, less than 15% of samples in both S_a and S_i are classified as malicious, indicating a poor robustness against test cases. On m_i , the average detection rate of the S_a samples increase to 67.3% while their prediction values range from 0.23 to 0.82; in contrast, the detection rates of S_i on m_a are around 78.2%, while their prediction values range from 0.69 to 0.91, which means that m_a can detect more test cases than m_i . Note, the detection rate of m_a on seed malware has been improved to 87.78%, compared with the 80.56% detection rate of m_i . From the experiment result, a trade-off to apply AMM on the improvement of anti-malware detector is presented: compared with the original detector, 7.34% detection rate degradation on seed malware vs. more than 55% increasing on the detection rate against test cases. This result further indicates that models with important features removed (m_i) is inefficient against test cases with AMM-based features manipulated (S_a). Meanwhile, improved detectors guided by AMM values (m_a) have better robustness on seed malware and adversarial detection.

Comparison of AMM and SAGE. Next, we compare how fragile and important features impact the detection. Figure 12 illustrates the distribution of top-75 AMM and SAGE features (26 features are selected both by AMM and SAGE). The x-axis indicates normalized AMM values; the y-axis is normalized SAGE values; the size of the scatter point indicates the number of evaded samples manipulating the corresponding features. As shown in the figure, samples that manipulate important features alone account for only a small portion of the evaded samples while most evaded samples are generated by manipulating features with large AMM values (pink vs. blue), which indicates that AMM is more efficient than SAGE to generate test cases. In addition, this result also explains why AMM-based test cases can bypass the improved model, even they removed important features, since these models are still vulnerable to test cases that only manipulate fragile features with large AMM values.

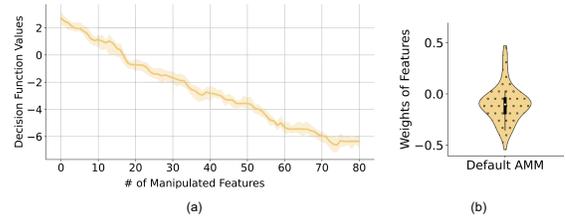


Figure 13: (a) Decision function values of samples with different numbers of features manipulated; (b) Distribution of weights of features in SVM model.

Takeaway 3:

- Machine learning-based antivirus products should consider improving the detectors with our AMM approach.
- AMM selects the fragile features and values that are efficient to flip the prediction results.

7 DISCUSSION

In this section, we conduct two case studies to understand why some test cases can or cannot evade detection, and further discuss threats to validity of our testing framework.

7.1 A Case Study on SVM Robustness

In prior experiments, the Drebin-SVM detector could not detect test cases from all three generation strategies. Here we conduct a case study to explore the detection robustness of SVM against test cases. First, we generate test cases by manipulating between 1 and 80 features (selected from Drebin-LGBM) on one seed. We then calculate the decision function value for each test cases, which represents the distance from the test case to the decision boundary. Positive decision function values represent malicious while negative ones represent benign. As shown in Figure 13(a) test cases generated by AMM-based strategy invert the prediction results after manipulating 16 features and push the generation towards benign with more features manipulated. From the result, we speculate that the features selected from Drebin-LGBM may occupy large weights in the Drebin-SVM, so that they can invert the prediction quickly. To verify this, we export the weights of each feature in the SVM detector and compare the weight values of selected features. Figure 13(b) illustrates the weight values of the selected features, where the y-axis represents the weight of each feature. We see that most weight values of features selected by the AMM strategy from Drebin-LGBM have relatively large negative values in the SVM model, making the prediction decision values negative (*i.e.*, benign), which matches the result in Figure 7. This case study confirms that features selected from LGBM occupy large negative weights in the SVM, making the prediction result of the adversarial sample benign.

7.2 A Case Study on Seed Selection

Our prior results have shown that not all test cases can evade detection. The reason could be either that the number of manipulable features in a seed is not enough to invert the prediction, or that the manipulated features have a limited impact on the prediction. To

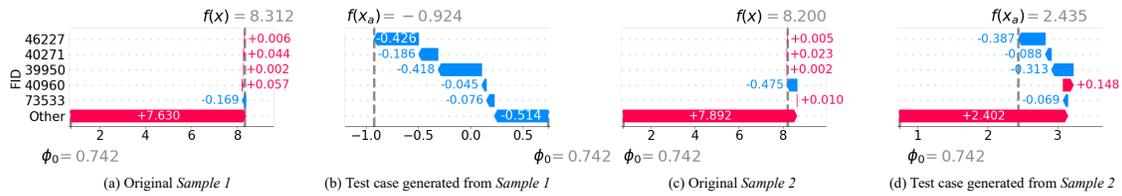


Figure 14: SHAP values of two APK samples before and after test case generation.

explore the reason, we choose two seed malicious APK examples, *Sample 1* and *Sample 2*, to generate their test cases. We further generate SHAP values of the original and test cases (with $N = 75$ features selected) of *Sample 1* and *Sample 2* to analyze the impact of the manipulated features.

The x-axis in Figure 14 indicates the prediction value, and the y-axis is the feature ID. $f(x)$ and $f(x_a)$ are raw scores of the original and test cases given by the Drebin-LGBM detector. We use red bars to indicate positive SHAP values, and blue bars to highlight negative SHAP values of each feature (i.e., ϕ_j in Equation 2). Feature IDs shown in the figure are parts of manipulated features that had the greatest impact on the two samples. The test case of *Sample 1* inverts its prediction as benign, and that of *Sample 2* remains malicious. Figure 14(b) shows that the SHAP values of manipulated features change significantly towards negative (blue bars), thereby pushing the output, $f(x_a)$, towards negative. In contrast, the SHAP values of the features of *Sample 2*, shown in Figure 14(d), change far less. Specifically, only features 46227, 40271, 39950 and 73533 are manipulated towards negative with less magnitude comparing to *Sample 1*, while 40960 is not towards negative at all. This means that we cannot manipulate enough features to force the decision making towards benign for *Sample 2*. This result indicates that the manipulated features have limited impact on *Sample 2* to invert the result from malicious to benign. In practice, after we increasing the number of manipulated features to 290, *Sample 2* is identified as benign. Therefore, the capacity of a seed depends on how many features have malicious-oriented values that we can manipulate in the sample. However, infinitely increasing the number of selected features would lead to a heavy computational load and decrease the efficiency of measurement.

7.3 Threats to Validity

Adaptive attacks against our testing framework. Existing technology against adversarial generation, for example, adversarial training [55] and differential privacy (DP) [20, 38, 51], could be effective against our proposed testing framework. For adversarial training, a machine learning-based detector could be trained with AMM generated test cases to falsely raise the detection rate on test cases. However, this may also increase the false positive rate on benign samples as the model will be forced to learn benign features as malicious, leading to a degraded performance. On the other hand, DP-based robust machine learning techniques cannot bypass our tests, because unbounded random perturbations may break the generated samples' functionality.

Dynamic detection. Since we only insert static unreachable instructions into the malware sample, a detector with dynamic detection will easily pass our testing. Feature-space manipulation and

problem-space obfuscation rely on static syntactic and structural modification. These modifications can be used to test static machine learning-based detectors and rule-based antivirus engines. However, the malicious behaviors preserved in the test cases will still be exposed during run-time and identified by the detectors that adopt dynamic analysis. Considering that dynamic feature detection consumes more resources to monitor, this approach may be impractical on a large scale and static approaches are still widely used in practice. In this research, we focus on exposing the weaknesses of malware detectors with an explainable method.

8 CONCLUSION

This paper has proposed an explainability-guided malware detector testing framework. The framework performs test case generation, relying on Accrued Malicious Magnitude (AMM) to guide the feature selection and a binary builder to map feature-space manipulations onto problem-space binaries. We then use our framework to test the robustness of state-of-the-art malware detectors. Our research includes the following key findings: (i) commercial antivirus engines and state-of-the-art machine learning detectors are vulnerable to AMM-based test cases; (ii) the transferability of AMM test cases relies on the overlaps of features with large AMM values between different machine learning models; and (iii) AMM values can effectively measure the fragility of features and explain the capability of flipping classification results. According to our findings, we suggest that machine learning-based AV products should consider using the AMM values to improve their robustness. Exploring the latter constitutes our key line of future work, as we believe this could prompt a new approach to defending against malware evasion attacks.

DATA AVAILABILITY

The source code, excluding the test case generator, of this project is publicly available at <https://github.com/Immor278/AMM>. Considering the potential security issues, we will not release the test case generator and any test cases, as well as the information of commercial antivirus involved in our evaluation, except for academic uses that are approved by our institutional ethics committee.

ACKNOWLEDGMENTS

The work has been supported by the Cyber Security Research Centre Limited whose activities are partially funded by the Australian Government's Cooperative Research Centres Program. Minhui Xue and Ruoxi Sun are the corresponding authors of this paper.

REFERENCES

- [1] [n. d.]. AI and machine learning. <https://www.avast.com/en-us/technology/ai-and-machine-learning>.
- [2] [n. d.]. Androguard documentation. <https://androguard.readthedocs.io/en/latest/>.
- [3] [n. d.]. Apktool - A tool for reverse engineering Android APK files. <https://androguard.readthedocs.io/en/latest/>.
- [4] [n. d.]. Cyber security statistics. <https://purplesec.us/resources/cyber-security-statistics/>.
- [5] [n. d.]. ESET smart security. <https://www.eset.com/>.
- [6] [n. d.]. How anti-virus software works. <https://cs.stanford.edu/people/eroberts/cs201/projects/viruses/anti-virus.html>.
- [7] [n. d.]. Kaspersky cyber security solutions for home & business. <https://www.kaspersky.com.au/>.
- [8] [n. d.]. Kiteshield. <https://github.com/GunshipPenguin/kiteshield>.
- [9] [n. d.]. More than free antivirus. <https://www.avast.com>.
- [10] . pefile. <https://github.com/erocarrera/pefile>.
- [11] [n. d.]. RelocBonus. <https://github.com/nickcano/RelocBonus>.
- [12] . Smali. <https://github.com/JesusFreke/smali>.
- [13] [n. d.]. Sonicwall research malware attacks 2019. <https://www.msspalert.com/cybersecurity-research/sonicwall-research-malware-attacks-2019/>.
- [14] . Virbox protector. <https://www.sense.com.cn/VirboxProtector.html>.
- [15] . VirusShare. <https://www.virusshare.com>.
- [16] . VirusTotal. <https://www.virustotal.com>.
- [17] [n. d.]. VMProject software protection. <https://vmpsoft.com/>.
- [18] 2021. Global antivirus software market report 2021: COVID-19 growth and change. <https://www.researchandmarkets.com/reports/5505347/global-antivirus-software-market-report-2021>.
- [19] 2022. As tanks rolled into Ukraine, so did malware. <https://www.nytimes.com/2022/02/28/us/politics/ukraine-russia-microsoft.html>.
- [20] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *ACM SIGSAC Conference on Computer and Communications Security*.
- [21] Hyrum S Anderson, Anant Kharkar, Bobby Filar, David Evans, and Phil Roth. 2018. Learning to evade static PE machine learning malware models via reinforcement learning. *arXiv preprint arXiv:1801.08917* (2018).
- [22] Hyrum S Anderson and Phil Roth. 2018. Ember: An open dataset for training static PE malware machine learning models. *arXiv preprint arXiv:1804.04637* (2018).
- [23] Simone Aonzo, Gabriel Claudiu Georgiu, Luca Verderame, and Alessio Merlo. 2020. Obfuscapk: An open-source black-box obfuscation tool for Android apps. *SoftwareX* (2020).
- [24] Daniel Arp, Michael Spreitzenbarth, Malte Hubner, Hugo Gascon, and Konrad Rieck. 2014. Drebin: Effective and explainable detection of Android malware in your pocket. In *Network and Distributed System Security Symposium (NDSS)*.
- [25] Erin Avllazagaj, Ziyun Zhu, Leyla Bilge, Davide Balzarotti, and Tudor Dumitras. 2021. When Malware Changed Its Mind: An Empirical Study of Variable Program Behaviors in the Real World. In *30th USENIX Security Symposium (USENIX Security 21)*. 3487–3504.
- [26] Federico Barbero, Feargus Pendlebury, Fabio Pierazzi, and Lorenzo Cavallaro. 2022. Transcending transcend: Revisiting malware classification with conformal evaluation. In *IEEE Symposium on Security and Privacy (SP)*.
- [27] Haipeng Cai. 2020. Assessing and improving malware detection sustainability through app evolution studies. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 29, 2 (2020), 1–28.
- [28] Yuxin Cao, Xi Xiao, Ruoxi Sun, Derui Wang, Minhui Xue, and Sheng Wen. 2023. StyleFool: Fooling Video Classification Systems via Style Transfer. In *44th IEEE Symposium on Security and Privacy (IEEE S&P)*. 818–835.
- [29] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. 2019. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705* (2019).
- [30] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy (SP)*.
- [31] Sen Chen, Minhui Xue, Lingling Fan, Lei Ma, Yang Liu, and Lihua Xu. 2019. How can we craft large-scale Android malware? An automated poisoning attack. In *International Workshop on Artificial Intelligence for Mobile (AI4Mobile)*.
- [32] Yizheng Chen, Shiqi Wang, Dongdong She, and Suman Jana. 2020. On training robust {PDF} malware classifiers. In *29th {USENIX} Security Symposium ({USENIX} Security 20)*. 2343–2360.
- [33] Ian Covert, Scott M Lundberg, and Su-In Lee. 2020. Understanding global feature contributions with additive importance measures. *Advances in Neural Information Processing Systems* 33 (2020), 17212–17223.
- [34] Luca Demetrio, Battista Biggio, Giovanni Lagorio, Fabio Roli, and Alessandro Armando. 2021. Functionality-preserving black-box optimization of adversarial windows malware. *IEEE Transactions on Information Forensics and Security* (2021).
- [35] Luca Demetrio, Battista Biggio, Giovanni Lagorio, Fabio Roli, and Alessandro Armando. 2021. Functionality-Preserving Black-Box Optimization of Adversarial Windows Malware. *IEEE Transactions on Information Forensics and Security* 16 (2021), 3469–3478. <https://doi.org/10.1109/TIFS.2021.3082330>
- [36] Ambra Demontis, Marco Melis, Battista Biggio, Davide Maiorca, Daniel Arp, Konrad Rieck, Igino Corona, Giorgio Giacinto, and Fabio Roli. 2017. Yes, machine learning can be more secure! A case study on Android malware detection. *IEEE Transactions on Dependable and Secure Computing* (2017).
- [37] Ambra Demontis, Marco Melis, Maura Pintor, Matthew Jagielski, Battista Biggio, Alina Oprea, Cristina Nita-Rotaru, and Fabio Roli. 2019. Why do adversarial attacks transfer? Explaining transferability of evasion and poisoning attacks. In *USENIX Security Symposium (USENIX Security)*.
- [38] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography, Third Theory of Cryptography Conference (TCC)*.
- [39] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. 2018. Robust physical-world attacks on deep learning visual classification. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- [40] Joshua Garcia, Mahmoud Hammad, and Sam Malek. 2018. Lightweight, obfuscation-resilient detection and family identification of android malware. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 26, 3 (2018), 1–29.
- [41] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. *International Conference on Learning Representations* (2015).
- [42] Kathrin Grosse, Nicolas Papernot, Praveen Manoharan, Michael Backes, and Patrick McDaniel. 2016. Adversarial perturbations against deep neural networks for malware classification. *arXiv preprint arXiv:1606.04435* (2016).
- [43] Kathrin Grosse, Nicolas Papernot, Praveen Manoharan, Michael Backes, and Patrick McDaniel. 2017. Adversarial examples for malware detection. In *European Symposium on Research in Computer Security*.
- [44] Mahmoud Hammad, Joshua Garcia, and Sam Malek. 2018. A large-scale empirical study on the effects of code obfuscations on Android apps and anti-malware products. In *International Conference on Software Engineering*. 421–431.
- [45] Richard Harang and Ethan M. Rudd. 2020. SOREL-20M: A large scale benchmark dataset for malicious PE detection. *arXiv preprint arXiv:2012.07634* (2020).
- [46] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. 2019. Adversarial examples are not bugs, they are features. *Advances in neural information processing systems* 32 (2019).
- [47] Junho Jung, Chanil Jeon, Max Wolotsky, Insu Yun, and Taesoo Kim. 2017. AVPASS: Leaking and bypassing antivirus detection model automatically. In *Black Hat USA Briefings (Black Hat USA)*.
- [48] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. LightGBM: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems* (2017).
- [49] TaeGuen Kim, BooJoong Kang, Mina Rho, Sakir Sezer, and Eul Gyu Im. 2018. A multimodal deep learning method for Android malware detection using various features. *IEEE Transactions on Information Forensics and Security* (2018).
- [50] Dhillung Kirat and Giovanni Vigna. 2015. MalGene: Automatic extraction of malware analysis evasion signature. In *ACM SIGSAC Conference on Computer and Communications Security*.
- [51] Mathias Lécuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. 2019. Certified robustness to adversarial examples with differential privacy. In *IEEE Symposium on Security and Privacy (SP)*.
- [52] Deqiang Li, Tian Qiu, Shuo Chen, Qianmu Li, and Shouhuai Xu. 2021. Can we leverage predictive uncertainty to detect dataset shift and adversarial examples in Android malware detection? *Annual Computer Security Applications Conference (ACSAC)* (2021).
- [53] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *International Conference on Neural Information Processing Systems*.
- [54] Wanlun Ma, Derui Wang, Ruoxi Sun, Minhui Xue, Sheng Wen, and Yang Xiang. 2023. The “Beatrix” resurrections: Robust backdoor detection via Gram matrices. In *The Network and Distributed System Security Symposium (NDSS)*.
- [55] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*.
- [56] Davide Maiorca, Davide Ariu, Igino Corona, Marco Aresu, and Giorgio Giacinto. 2015. Stealth attacks: An extended insight into the obfuscation effects on Android malware. *Computers & Security* (2015).
- [57] Enrico Mariconti, Lucky Onwuzurike, Panagiotis Andriotis, Emiliano De Cristofaro, Gordon Ross, and Gianluca Stringhini. 2017. MAMADROID: Detecting Android Malware by Building Markov Chains of Behavioral Models. In *Proceedings of the Annual Symposium on Network and Distributed System Security (NDSS)*.
- [58] Fabio Pierazzi, Feargus Pendlebury, Jacopo Cortellazzi, and Lorenzo Cavallaro. 2020. Intriguing properties of adversarial ML attacks in the problem space. In *IEEE Symposium on Security and Privacy (SP)*.
- [59] Matteo Pomilia. 2016. A study on obfuscation techniques for Android malware. *Sapienza University of Rome* (2016).

- [60] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “Why should I trust you?”: Explaining the predictions of any classifier. In *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [61] Giorgio Severi, Jim Meyer, Scott Coull, and Alina Oprea. 2021. Explanation-guided backdoor poisoning attacks against malware classifiers. In *USENIX Security Symposium (USENIX Security)*.
- [62] Maryam Shahpasand, Len Hamey, Dinusha Vatsalan, and Minhui Xue. 2019. Adversarial attacks on mobile malware detection. In *International Workshop on Artificial Intelligence for Mobile (AI4Mobile)*.
- [63] SkyLight. [n. d.]. Cylance, I kill you! <https://skylightcyber.com/2019/07/18/cylance-i-kill-you/>.
- [64] Wei Song, Xuezixiang Li, Sadia Afroz, Deepali Garg, Dmitry Kuznetsov, and Heng Yin. 2022. MAB-Malware: A Reinforcement Learning Framework for Blackbox Generation of Adversarial Malware. In *Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security*.
- [65] Octavian Suciu, Radu Marginean, Yigitcan Kaya, Hal Daume III, and Tudor Dumitras. 2018. When does machine learning FAIL? generalized transferability for evasion and poisoning attacks. In *27th USENIX Security Symposium*. 1299–1316.
- [66] Ruoxi Sun, Wei Wang, Minhui Xue, Gareth Tyson, Seyit Camtepe, and Damith C. Ranasinghe. 2021. An Empirical Assessment of Global COVID-19 Contact Tracing Applications. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. 1085–1097.
- [67] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*.
- [68] Microsoft Defender Security Research Team. 2019. New machine learning model sifts through the good to unearth the bad in evasive malware. <https://www.microsoft.com/security/blog/2019/07/25/new-machine-learning-model-sifts-through-the-good-to-unearth-the-bad-in-evasive-malware/>.
- [69] Romain Thomas. 2017. LIEF - Library to instrument executable formats. <https://lief.quarkslab.com/>.
- [70] Xabier Ugarte-Pedrero, Davide Balzarotti, Igor Santos, and Pablo G Bringas. 2015. SoK: Deep packer inspection: A longitudinal study of the complexity of run-time packers. In *2015 IEEE Symposium on Security and Privacy*. IEEE, 659–673.
- [71] Giovanni Vigna and Davide Balzarotti. 2018. When Malware is {Packin’} Heat. In *Enigma 2018 (Enigma 2018)*.
- [72] Ke Xu, Yingjiu Li, and Robert H. Deng. 2016. ICCDetector: ICC-based malware detection on Android. *IEEE Transactions on Information Forensics and Security* (2016).
- [73] Ke Xu, Yingjiu Li, Robert H. Deng, and Kai Chen. 2018. DeepRefiner: Multi-layer Android malware detection system applying deep neural networks. In *IEEE European Symposium on Security and Privacy*.
- [74] Weilin Xu, Yanjun Qi, and David Evans. 2016. Automatically evading classifiers. In *Network and Distributed Systems Symposium*.