# Query-Efficient Black-Box Attack Against Sequence-Based Malware Classifiers

Ishai Rosenberg, Asaf Shabtai, Yuval Elovici and Lior Rokach
Ben-Gurion University of The Negev

October 6, 2020

### Abstract

In this paper, we present a generic, query-efficient black-box attack against API call-based machine learning malware classifiers. We generate adversarial examples by modifying the malware's API call sequences and non-sequential features (printable strings), and these adversarial examples will be misclassified by the target malware classifier without affecting the malware's functionality.

In contrast to previous studies, our attack minimizes the number of malware classifier queries required. In addition, in our attack, the attacker must only know the class predicted by the malware classifier; attacker knowledge of the malware classifier's confidence score is optional. We evaluate the attack effectiveness when attacks are performed against a variety of malware classifier architectures, including recurrent neural network (RNN) variants, deep neural networks, support vector machines, and gradient boosted decision trees. Our attack success rate is around 98% when the classifier's confidence score is known and 64% when just the classifier's predicted class is known.

We implement four state-of-the-art query-efficient attacks and show that our attack requires fewer queries and less knowledge about the attacked model's architecture than other existing query-efficient attacks, making it practical for attacking cloud-based malware classifiers at a minimal cost.

## 1 Introduction

Next generation anti-malware products use machine learning and deep learning models instead of signatures and heuristics, in order to detect previously unseen malware. Windows-based API call sequence-based classifiers, such as [11] or [8], provide state-of-the-art detection performance [15]. However, those classifiers are vulnerable to adversarial example attacks. Adversarial examples are samples that are perturbed (modified), so they will be incorrectly classified by the target classifier.

In this paper, we demonstrate a novel *query-efficient* black-box attack against many types of malware classifiers, including RNN variants. We implement our attack on malware classifiers that are used to classify a process as malicious or benign. This classification is based on both API call sequences and additional discrete features (e.g., printable strings). We consider the challenging case of machine learning as a service (MLaaS). Examples of such services are Amazon ML [1] and GCP [5].

In this case, the attacker continuously queries the target malware classifier (e.g., [6]) and modifies the API call sequences the malware generates until the sequence is classified as benign. The attacker pays for every query of the target malware classifier and therefore aims to minimize the number of queries made to such cloud services when performing an attack. Another reason for minimizing the number of queries is that many queries from the same computer might arouse suspicion of an adversarial attack attempt, causing the cloud service to stop responding to those queries, e.g., using stateful defenses that keep track of queries to the system [22]. While the attacker may use a botnet to issue many queries, this approach would increase the attack's cost dramatically.

We develop an *end-to-end attack* [40, 41] (a.k.a, problem-space attack [38]) by recrafting the malware behavior so it can evade detection by such machine learning malware classifiers while minimizing the amount of queries to the malware classifier (meaning that the attack is *query-efficient*).

The main focus of most existing research (e.g., [30, 43]) is on the query-efficient generation of adversarial examples for images. This is different from our work, which is focused on generating adversarial API sequences, in two respects:

1. In the case of adversarial API sequences, one must verify that the original functionality of the malware remains intact. Thus, one cannot simply generate an adversarial feature vector but must generate an executable file containing the corresponding working malware behavior.

2. API sequences consist of discrete symbols of variable lengths, while images are represented as matrices with fixed dimensions, and the values of the matrices are continuous.

In this paper, we modify the malware's behavior by modifying the API call sequences it generates. We also modify static, non-sequential features: printable strings inside the process (Appendix D). We present eight novel attacks that are different in terms of three parameters:

1. Attacker knowledge - The attacker might have knowledge about the target classifier's confidence score (*score-based attack*) or only about the predicted class (*decision-based attack*). This knowledge can affect the way we modify both the positions and the types of modified API calls.

2. Modified API call types (values) - The attacker can either modify API calls randomly (random perturbation) or take them from a generative adversarial network (GAN) generating benign samples (benign perturbation).

3. Method to select the number of modified API calls - The attacker can either modify API calls until the sample successfully evades detection with minimal perturbation (linear iteration attack) or start with a maximum number of modifications and minimize them as long as evasion is maintained (logarithmic backtracking attack).

Each of the eight attacks we present is a combination of these three parameters, each of which has two different variants (score-based or decision-based; benign perturbation or random perturbation; logarithmic backtracking or linear iteration). As a benchmark, we adapted four state-of-the-art query-efficient attacks. We showed that our attacks are equal or outperform them, obtaining a higher success rate (98% using the classifier's confidence score and 88% without it) for a fixed number of queries.

The contributions of our paper are as follows:

1. This is the first *query-efficient, end-to-end, black-box* adversarial attack for *both sequential and non-sequential* input.

2. We provide two variants for our attack, to fit the knowledge available for the attacker (confidence score or label only), thus fitting practical cyber security domain use cases.

3. This is the first usage of GAN (previously used for other purposes) to generate benign API call trace samples to produce a query-efficient attack.

4. This is the first usage of self adaptive evolutionary algorithm (previously used for other purposes) to implement a query-efficient score-based attack.

## 2    Background and Related Work

The search for adversarial examples, such as those used in our attack, can be formalized as a minimization problem:

$$\arg_{\mathbf{r}} \min f(\mathbf{x} + \mathbf{r}) \neq f(\mathbf{x}) \; s.t. \; \mathbf{x} + \mathbf{r} \in \mathbf{D} \tag{1}$$

The input $\mathbf{x}$, correctly classified by the classifier $f$, is perturbed with $\mathbf{r}$, such that the resulting adversarial example $\mathbf{x} + \mathbf{r}$ remains in the input domain $\mathbf{D}$ but is assigned a different label than $\boldsymbol{x}$.

There are three types of adversarial example generation methods:

In *gradient-based attacks*, adversarial perturbations are generated in the direction of the gradient, that is, in the direction with the maximum effect on the classifier's output (e.g., fast gradient sign method [26]). Hu and Tan [28] used RNN GAN to generate invalid API calls and insert them into the original API sequences. Gumbel-Softmax, a one-hot continuous distribution estimator, was used to deliver gradient information between the generative RNN and substitute RNN. A white-box gradient-based attack against RNNs demonstrated against long short-term memory (LSTM) architecture for sentiment classification of a movie review dataset was shown in [36]. A black-box variant, which facilitates the use of a substitute model, was presented in [42]. The attack in this paper is different in a few ways:

1. As shown in Section 4.2, the attack described in [42] requires more target classifier queries and greater computing power to generate a substitute model.

2. We use a different adversarial example generation algorithm, which uses a stochastic approach rather than a gradient-based approach, making it harder to defend against (as mentioned in Section 4.3).

3. Our decision-based attack is generic and doesn't require a per malware pre-deployment phase to generate the adversarial sequence (either using a GAN, as in [28], or a substitute model, as in [42]). Moreover, generation takes place at run time, making it even more generic and easier to deploy.

*Score-based attacks* are based on knowledge of the target classifier's confidence score. Previous research used a genetic algorithm (GA), where the fitness of the genetic variants is defined in terms of the target classifier's confidence score, to generate adversarial examples that bypass PDF malware and image recognition classifiers, respectively [45, 16]. Those attacks used a computationally expensive GA compared to our approach and were only evaluated when performed against support vector

machine (SVM), random forest, and CNN classifiers using static features only, and was not evaluated against recurrent neural network variants using both static and dynamic analysis features, as we do. In [43], the simultaneous perturbation stochastic approximation (SPSA) method was used, since its gradient approximation requires only two queries (for the target classifier's score) per iteration, regardless of the dimension of the optimization problem. Alzantot et al. [17] presented a sequence-based attack algorithm that exploits population-based gradient-free optimization via GA; in this study, the attack was performed against a natural language processing (NLP) sentiment analysis classifier. While this attack can be used against RNNs, it requires more queries than our attack (see Table 3).

*Decision-based attacks* only use the label predicted by the target classifier. Ilyas et al. [30] used natural evolutionary strategies (NES) optimization to enable query-efficient gradient estimation, which leads to the generation of misclassified images as seen in gradient-based attacks. Dang et al. [24] used the rate of feature modifications from known malicious and benign samples as the score and used a hill climbing approach to minimize this score, evading SVM and random forest PDF malware classifiers based on static features in an efficient manner. Our approach, on the other hand, is more generic and can handle RNN classifiers and multiple feature types. The performance differences between our approaches and [24] are presented in Section 4.2.2.

All of the currently published score and decision-based attacks differ from our proposed attack in that:

1. They *only* deal with CNNs, random forest, and SVM classifiers, using non-sequential input, as opposed to *all* state-of-the-art classifiers (including RNN variants), using discrete or sequence input, as in our attack.

2. They deal primarily with images and rarely fit the attack requirements of the cyber security domain: while changing a pixel's color doesn't "break" the image, modifying an API call might harm the malware functionality. In addition, small perturbations, such as those suggested in [30, 43], are not applicable for discrete API calls: you can't change *WriteFile()* to *WriteFile()*+0.001 in order to estimate the gradient to perturb the adversarial example in the right direction; you need to modify it to an entirely different API. This is reflected in Table 3.

3. They did not present an end-to-end framework to implement the attack in the cyber security domain. Thus, the attack might be used for generating adversarial malware feature vectors but not for generating a working adversarial malware sample.

The differences between those attacks and our attacks are summarized in Table 1.

## 3   Methodology

### 3.1   Attacking API Call-Based Malware Classifiers

An overview of the malware classification process is shown in Figure 2 (in the appendix).

Assume a malware classifier whose input is a sequence of API calls made by the inspected process. API call sequences can be millions of API calls long, making it impossible to train such a classifier on the entire sequence at once, due to training time and GPU memory constraints. Thus, the target classifier uses a non-overlapping sliding window approach [42]: Each API call sequence is divided into windows of $k$ API calls. Detection is performed on each window in turn, and if any

Table 1: Comparison to Previous Work

| Attack Type | Domain | Input Type | Query-Efficient? | Score/Decision-Based? |
|---|---|---|---|---|
| Rosenberg et al. [42] (Gradient-Based Attack) | Cyber Security | Sequence, Non-sequential, Mixed | No | Decision |
| Uesato et al. [43] | Image Recognition | Non-sequential | Yes | Score |
| Ilyas et al. [30] | Image Recognition | Non-sequential | Yes | Score |
| Alzantot et al. [16] | NLP | Sequence | Yes | Score |
| Our Score-Based Attack | Cyber Security | Sequence, Non-sequential, Mixed | Yes | Score |
| Our Decision-Based Attack | Cyber Security | Sequence, Non-sequential, Mixed | Yes | Decision |

window is classified as malicious, the entire sequence is considered malicious. Thus, even cases such as malicious payloads injected into goodware (e.g., using Metasploit), where only a small subset of the sequence is malicious, would be detected.

We use one-hot encoding for each API call type in order to cope with the limitations of scikit-learn's implementation of decision trees and random forests, as mentioned in [10]. The output of each classifier is either the predicted class (whether the API call trace is malicious or benign) or the confidence score of the prediction (a value between 0 for a benign process to 1.0 for a malicious process). Appendix B contains a description of the classifiers used in our study and their hyperparameters.

In order to attack such a malware classifier, we want to add API calls without changing the malware functionality. Removing API calls without modifying the malware functionality requires complex analysis of the malware code and thus cannot be scaled. Therefore, we use a *no-op mimicry attack* [44], that is, we add API calls with no effect or an irrelevant effect on the malware functionality. We use this method regardless of the perturbation used to generate the API call type, either random or benign perturbations (shown below). Almost every API call can become a no-op if provided with the right arguments, e.g., opening a non-existent file.

However, analyzing arguments would make our attack easier to detect, e.g., by considering only successful API calls and ignoring failed API calls or by looking for irregularities in the arguments of the API calls (e.g., invalid file handles, etc.). In order to address this issue, we use valid (non-null) arguments with a no-op effect, such as writing into a temporary file handle (instead of an invalid file handle) or reading zero bytes from a valid file. This makes detecting the no-op API calls much harder, since the API call runs correctly, with a return value indicative of success. It is extremely challenging for the malware classifier to differentiate between malware that is trying to read a non-existent registry key as an added adversarial no-op and a benign application functionality, e.g., trying to find a registry key containing information from previous runs and creating it if it doesn't
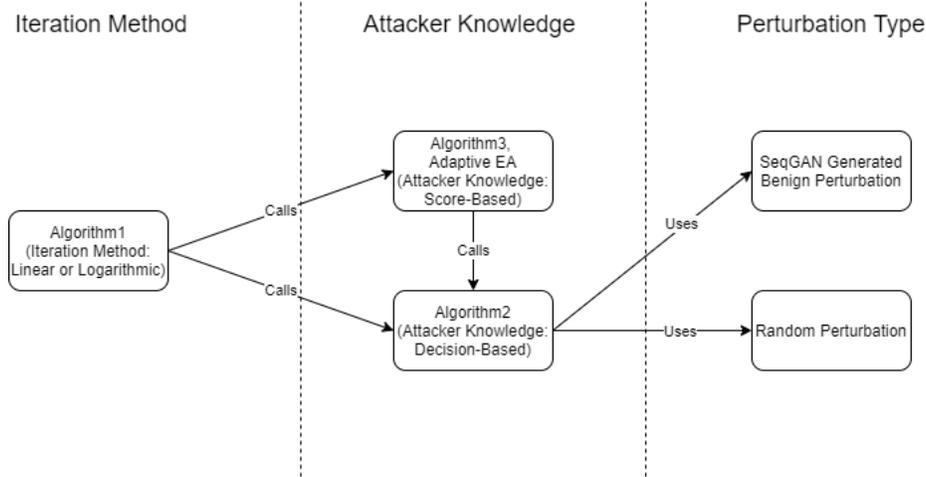
Figure 1: The Proposed Attack's Flow

exist (for instance, during the first run of the application). This makes our problem-space attack robust to preprocessing [38].

To conclude, we believe that any detector will suffer from a very high false positive rate and thus will not be a practical solution for detecting our attacks. We don't add API types which cannot be no-opped and thus always affect the malware functionality. In this study, we focus on the 314 API call types monitored by Cuckoo Sandbox. Of those, only two API call types are not added because they can't be no-opped: *ExitWindowsEx*() and *NtTerminateProcess*(). Every other API call have a no-op variant, generated by the abovementioned method.

## 3.2 The Proposed Black-Box Attack

Our proposed attack's flow is detailed in Figure 1. As mentioned before, our proposed attacks can be characterized by the knowledge the attacker has, the method for selecting which new API call to add (termed *perturbation type*), and the method to select the number of added API calls (termed *iteration method*). As can be seen in Figure 1, these characteristics affect the attack's flow. These characteristics are described in the subsections below.

### 3.2.1 Iteration Method

In this subsection, we describe the method to select the number of added API calls by the attack. In addition to the regular *linear iteration* method, we propose an efficient *logarithmic backtracking* transformation as a method for determining the number of API calls to add. This method starts with a large ratio of perturbation (that is, a larger number of API calls are added to the original sequence to fool the classifier) which rapidly decreases as long as the sequence remains misclassified.

In order to handle the entire API call sequence, we use Algorithm 1. In this algorithm, the attacker splits the malicious API call sequence $x^m$ into windows of $n$ API calls ($w_j^m$), similar to the division made by the malware classifier, and modifies each window in turn using Algorithm 3 (described below; line 8). The API calls "pushed out" from $w_j^m$ will become the beginning of

$\boldsymbol{w_{j+1}^m}$, so no API is ignored. The adversarial window size $n$ might be different from the malware classifier's window size $k$, which is not known to the attacker. As shown in Section 4.2, this has a negligible effect on the attack performance. In the case of a benign perturbation, the benign API call sequence $\boldsymbol{x^b}$ is similarly split into windows of $n$ API calls ($\boldsymbol{w_j^b}$).

The adversarial API call sequence length $l$ might be larger than $n$, the length of the sliding window API call sequence that is used by the adversary. Therefore, the attack is performed sequentially on $\left\lceil \frac{l}{n} \right\rceil$ windows of $n$ API calls (e.g., for $l = 1000$ and $n{=}100$ the malware classifier would run on ten windows, with API call indices: 1..100, 101..200, ..., 900..1000). Note that the knowledge of $k$ (the window size of the malware classifier) is not required, as shown in Section 4.2.

$D$ is the vocabulary of available features. In this case, these features are all of the API call types recorded by the malware classifier, e.g., $CreateFileW()$. Note that $D$ is not necessarily known to the attacker. The attacker knows $D^{'}$, which might be a subset or superset of $D$. This knowledge of $D^{'}$ is a commonly accepted assumption about the attacker's knowledge [29]. In fact, it is enough for the attacker to know the feature type used by the target classifier (API call types in this study), which is public information that is usually published by classifier implementers. With that knowledge, the attacker can cover *all* API call types (several thousands) to generate $D^{'}$, which is a superset of $D$.

In our research, we observed that API call types in $D' - D$ are not monitored by the classifier, and thus adding them does not assist in creating adversarial examples; those API calls just add API call overhead to the modified sequence and serve as wasted queries. API call types in $D - D'$, unknown to the attacker, are not generated by the attack and therefore decreasing the adversarial feature space and thus decreasing the possibilities for generating modified sequences that can evade detection. Thus, when $D^{'}$ is a superset of $D$, the attack has higher overhead but remains as effective. An attacker might also reverse engineer the features from the malware classifier program, as has already done in real-world adversarial example generation for malware classifier scenarios [3].

In order to decrease the number of malware classifier queries (the number of calls to $f(.)$), we can use *logarithmic backtracking*. This iteration method is similar to binary search. In this case, we only query the malware classifier in Algorithm 1 after modifying $M_w$ API calls in Algorithm 2 (instead of querying the model after modifying a single API call in a linear iteration), which should be a sufficiently large perturbation to evade the malware classifier. Then, we start reducing the number of modified API calls by half before querying the malware classifier (lines 13-14) until it detects the sample again. Finally, we keep restoring half of the API calls we previously removed before querying (line 22), until we achieve a perturbation that fools the malware classifier with a minimal number of additional API calls and malware classifier queries.

In Algorithm 2, the attacker chooses the API calls to add and remove randomly. Note that we do not remove the malware's original API calls (only the no-op API calls that were previously added by the adversary), in order to prevent harm to its functionality. Since we add or remove half of the API calls each time, we perform $O(\log n)$ queries per adversarial sliding window if $IterationMethod$ is logarithmic backtracking, instead of the $O(n)$ queries that are performed if $IterationMethod$ is linear iteration (where $n$ is the size of the adversarial sliding window), making this attack query-efficient, as can be seen in Tables 3 and 4. While the proposed attack is designed for API call-based classifiers, it can be used for any adversarial sequence generation.

The benign sequence $\boldsymbol{x^b}$ is generated by a specially crafted GAN, which is described below.

---

**ALGORITHM 1:** Full Sequence Attack

---

1  **Input**: $f$ - black-box model,

2  $\boldsymbol{x^m}$ - malicious sequence to perturb, $\boldsymbol{x^b}$ - benign sequence to mimic,

3  $n$ - size of adversarial sliding window, $D^{'}$ - adversarial vocabulary,

4  $M_w$ - maximum API modifications per window, $PerturbType$ - benign or random perturbation,

5  $IterationMethod$ - logarithmic backtracking or linear iteration, $AttackerKnowledge$ - decision or score-based.

6

7  **for** each sliding windows $(\boldsymbol{w_j^m}, \boldsymbol{w_j^b})$ of $n$ API calls in $(\boldsymbol{x^m}, \boldsymbol{x^b})$, respectively:

8    $(\boldsymbol{w_j^m}, addedAPIs) =$

    $Algorithm2(f, \boldsymbol{w_j^m}, \boldsymbol{w_j^b}, n, D^{'}, M_w, PerturbType, IterationMethod, AttackerKnowledge)\}$

9    **if** $IterationMethod$ is logarithmic backtracking:

10      $remainingAPIs = addedAPIs$

11      **while** $(f(\boldsymbol{w_j^m}) = benign)$:

12        # Remove added API calls until evasion is lost:

13        Randomly split $addedAPIs$ into two equally sized groups: $remainingAPIs, deletedAPIs$

14        remove $deletedAPIs$ from $\boldsymbol{w_j^m}$

15        **if** $f(\boldsymbol{w_j^m}) = malicious$ :

16          $\boldsymbol{w_j^m} = \boldsymbol{w_j^m} + deletedAPIs - remainingAPIs$ # Remove $remainingAPIs$ instead of $deletedAPIs$ from $\boldsymbol{w_j^m}$

17          Switch between $remainingAPIs$ and $deletedAPIs$

18      $recoveredAPIs = deletedAPIs$

19      **while** $(f(\boldsymbol{w_j^m}) = malicious)$:

20        # While there are still added API calls that were removed, add them back until evasion is restored:

21        $recoveredAPIs =$ Randomly pick half of the API calls remaining in $deletedAPIs$

22        Add $recoveredAPIs$ to $\boldsymbol{w_j^m}$

23  **return** (perturbed) $\boldsymbol{x^m}$

---

### 3.2.2 Perturbation Type

In this subsection, we describe the methods to select what API calls to add by our attack. As an alternative to choosing the API calls randomly from all available API calls (*random perturbation*), we propose the *benign perturbation* method. In this method, instead of adding random API calls, we add API calls selected from sequences generated by a generative adversarial network (GAN) that has been trained to mimic real benign sequences. This concept is inspired by the modus operandi of biological viruses (malware) which are sometimes composed of human ("benign") proteins in order to evade the immune system (malware classifier) of the host.

When the attackers add an API call to our adversarial sequence, they want to have the maximum impact on the classifier's score. Thus, Algorithm 1 can take $x^b$, a benign API sequence, as input. The idea is that adding a "benign" API call would make the trace "more benign" than adding a random API call. This is due to the fact that no API call is malicious or benign per se. The context and flow of API calls determine the functionality (and therefore the code's "maliciousness"). Using benign perturbation creates a "benign API call context" and thus improves the attack effectiveness and also makes it more query-efficient, because fewer queries are needed to fool the classifier, as can be seen in Table 4, which uses benign perturbations, as opposed to Table 3, which uses linear iteration.

One way to generate $x^b$ is by taking the API call sequence of an actual benign sample from our dataset. The downside of this approach is that those hard-coded API calls can be detected explicitly as an evasion attack.

A better approach is to generate a different benign sequence each time, using a generative model. One way to do this is to use a generative adversarial network (GAN), with a stochastic input seed and an output of an API call sequence that is indistinguishable (to the discriminator classifier) from actual benign sequences from the dataset. This approach is rarely used for API call sequence generation, but it has been used for text generation. Note that this approach does not require queries to the target classifier (unlike, e.g., building a substitute model, as done in [42]).

In comparison to other approaches (e.g., VAE), using a GAN tends to generate better output. Most other methods require that the generative model has some particular functional form (like a Gaussian output layer). Moreover, all of the other approaches require that the generative model puts non-zero mass everywhere.

However, a challenge with the GAN approach is that the discrete outputs from the generative model make it difficult to pass the gradient update from the discriminative model to the generative model. Another challenge is that the discriminative model can only classify a complete input sequence. We used SeqGAN [46] implementation. Following a pretraining procedure that follows the MLE (maximum likelihood estimation) metric, the generator $G$ is modeled as a stochastic policy in reinforcement learning (RL). The agent's state is the API call subsequence of the first $t$ API call types in the sequence to be generated by the GAN, $s_t = [x_0, x_1, ..x_{t-1}]$. The agent's action is the next API call type in the sequence to be generated by the SeqGAN model $x_t \sim (x|s_t)$. The reward is the feedback given to $G$ by $D$ when evaluating the generated sequence, bypassing the gradient update challenge by directly performing a gradient policy update.

In a stochastic parameterized policy, the actions are drawn from a distribution that parameterizes the policy. An action may be sampled from, e.g., a normal distribution whose mean and variance will be predicted by the policy. The objective of $G$ is to generate a sequence from the start state $s_0$ in such a way that maximizes the expected end reward. This action-value function is estimated by the discriminator. However, $D$ only provides a reward at the end of a finished sequence. Yet, it is important that at every time-step, the fitness of both previous tokens as well

as future outcome are considered. For this, the policy gradient used in SeqGANs employs a Monte Carlo search with a roll-out policy (P) to sample the unknown remaining tokens and approximates the state-action value in an intermediate step.

We trained our GAN using a benign hold-out set of 3,000 sequences that were taken from the same distribution as the test set, available to the attacker. This hold-out set was not used subsequently as part of the test set, to avoid data leak to the attacker, artificially increasing the attack effectiveness.

We also tried other GAN architectures (e.g., GSGAN), but SeqGAN outperformed all of them (Appendix C). SeqGAN outperformed the random perturbation as well, as shown in Section 4.2.

### 3.2.3 Attacker Knowledge

This characteristic entails a trade off: the more information the attackers have about the classifier, the more query-efficient their attack would be. In this paper, the attacker may have access to the confidence score of the malware classifier (*score-based attack*), or only to the predicted label (*decision-based attack*). For the first case, we propose a score-based attack that uses a gradient-free optimization algorithm that until now has never been used for adversarial learning, outperforming state-of-the-art attacks for low number of queries to the attacked classifier. This attack is designed for discrete values in sequences of variable length. Thus, it fits API call sequences, as opposed to image pixels, which were the focus of most previous research.

While some malware classifiers do expose their confidence score (e.g., MAX, a non-sequence based NGAV product in VirusTotal online scanning service [12]), others do not. Therefore, we also implemented a decision-based attack, which is less query-efficient, but requires only knowledge about the predicted label of the malware classifier.

**Decision-Based Attack** In Algorithm 2, we show how to generate an adversarial sequence for a single API call window (out of the entire API call sequence).

---

**ALGORITHM 2:** Single Iteration Decision-Based Window Sequence Generation

---

1 **Input**: $f$ - black-box model, $\boldsymbol{w^m}$ - malicious sequence to perturb, of length $l^m \leq n$,

2 $\boldsymbol{w^b}$ - benign sequence to mimic, of length $l^b \leq n$, $n$ - size of adversarial sliding window,

3 $D^{'}$ - adversarial vocabulary, $M_w$ - maximum API modifications per window, $PerturbType$ - benign or random perturbation,

4 $IterationMethod$ (logarithmic backtracking or linear iteration.

5

6 $addedAPIs = \{\}$

7 **while** (($IterationMethod$ is linear iteration) and ($f(\boldsymbol{w^m}) = malicious$)) or ($|addedAPIs| < M_w$):

8     Randomly select an API's position $i$ in $\boldsymbol{w^m}$

9     **if** $PerturbType$ is benign perturbation:

10         Add $\boldsymbol{w^b}[i]$ to $\boldsymbol{w^m}$ at position $i$

11     **else:** $PerturbType$ is random perturbation

12         Add a random API in $D^{'}$ to $\boldsymbol{w^m}$ in position $i$

13     Add the new API and its position to $addedAPIs$

14 **if** ($f(\boldsymbol{w^m}) = malicious$) and ($|addedAPIs| = M_w$): **return** Failure

15 **return** ($\boldsymbol{w^m}, addedAPIs$) # $\boldsymbol{w^m}$ includes the perturbation

---

The perturbation added is either random API calls, a.k.a. *random perturbation* (line 12), or API calls of a benign sequence, a.k.a. *benign perturbation* (line 10). Since only the predicted class is available, the API calls are added in a random position $i$ in the API sequence (line 8). The adversaries randomly chooses $i$, since they do not have a better way of selecting $i$ without incurring significant statistical overhead. Note that the addition of an API call in position $i$ means that the API calls from position $i..n$ ($\boldsymbol{w^m}[i..n]$) are "pushed back" one position to make room for the new API call, in order to maintain the original sequence and preserve the original functionality of the code (in line 10 and 12). Since the sliding window has a fixed length, the last API call, $\boldsymbol{w^m}[n+1]$, is "pushed out" and removed from $\boldsymbol{w^m}$. This API call addition continues until the modified sequence $\boldsymbol{w^m}$ is classified as benign or more than $M_w$ API calls are added, reaching the maximum overhead limit (line 7). In this case the attack has failed (line 14). In the case of a *linear iteration* attack, the API calls are added one at a time, checking whether the perturbed sequence evades detection after each addition. The case of a *logarithmic backtracking* attack was explained in Section 3.2.1.

One might claim that a simpler attack can be used instead: insert $n-1$ no-op API calls after each API call from the original binary. This attack effectiveness is 100%, and no queries are needed to implement it, making it extremely query-efficient. However, this trivial attack has two major issues:

1. It is easy to detect this trivial attack using anomaly detection, since no actual benign program call trace is composed like that.

2. Such malware would run much slower than the original malware due to the additional API overhead, allowing the intrusion prevention systems of the victim to mitigate such malware, e.g., by terminating perturbed ransomware, after encrypting only several files, due to its perturbation induced slowness.

**Score-Based Attack**   When the confidence score of the malware classifier is also returned, score-based attacks (e.g., gradient-free optimization algorithms) can be applied as well. The merged flow for both attacks is described in Algorithm 3. Assuming the attacker has a budget of $B$ queries per API call window, the call to Algorithm 2 in line 9 of Algorithm 3 can be replaced with $B-\log n$ iterations (line 11) of minimizing $f(\boldsymbol{w^m})$ (lines 13 and 15) by one of the score minimization algorithms presented in Section 4.2.2. In order to use the same budget for all attacks, we chose $B = M_w$ (lines 9, 11).

All random perturbation variants try to minimize the target classifier score by modifying only the values of the added API call types (while the API call positions are random but fixed, as in Algorithm 2). Trying to modify both API types and positions with the same budget results in inferior performance (this is not shown due to space limits). All benign perturbation variants try to minimize the score by modifying only the API positions (while the API types are taken from the GAN's output).

The maximum additional API calls allowed per sliding window was set to 70 (50%). The search space for this optimization would either be the $M_w$ added API call type values (out of $|D|$ values each) if this is a random perturbation, or the $M_w$ added API call indices in the adversarial window (each with $n$ possible values) if this is a benign perturbation.

**Score-Based Query-Efficient Attack for Discrete Input Sequence**   Most state-of-the-art query-efficient attacks for images assume continuous input ([43, 30]) and underperform when used for discrete input (e.g., API calls or position indices), as shown in Table 3. Genetic algorithms

(GAs), which use mutation (random perturbation) in existing adversarial candidates and crossover between several candidates (i.e., a combination of parts of several candidates), are an exception. GAs work well with discrete sequences [17]. However, while crossover makes sense in the NLP domain (e.g., for compound sentences), it makes little sense for API call sequences, where each program has its own business logic. Another issue is the poor performance of a fixed mutation rate, usually used by GAs. It is better to use an adaptive mutation rate, which fits itself to the domain without knowledge expertise [23].

We decided to use the self-adaptive uniform mixing evolutionary algorithm (EA) proposed by Dang et al. [23]. It starts with a population of adversarial candidates, and in every generation, a new population of candidates is produced, and the old generation dies. Besides the adversarial sequence, each candidate carries an additional property: its mutation rate. Each new candidate is produced in the same way:

1. The best of two uniformly selected individuals is selected (i.e., tournament selection).

2. The selected parent individual changes its mutation rate between two mutation rates: low and high, with a fixed probability $p$. We used the values proposed in [23].

3. The parent replicates, with mutations occurring at the new mutation rate.

Although the selection mechanism does not take into account the mutation rate, the intuition is that appropriate mutation rates are correlated with high fitness.

We assume that the EA attack would be query-efficient due to the following reasons:

1. The usage of the adaptive mutation rate helps reducing the number of queries, before the highest impact element is added to the sequence.

2. Unlike other attacks (including other self-adaptive heuristic strategies, e.g., [43]), EA is effective in a discrete feature space. This makes it more query-efficient, because fewer queries are needed before it fools the target classifier, as can be seen in Tables 3 and 4.

3. The crossover (combining API calls from malicious and benign parents, instead of mutating a malicious parent) used by other attacks (e.g., GA, also efficient in discrete feature spaces) makes no sense in the cyber domain, because it adds redundant API calls of the benign parent. Thus, an algorithm that forgoes the crossover (and thus the addition of redundant API calls) is more query-efficient in the cyber domain.

## 4  Experimental Evaluation

### 4.1  Dataset and Target Malware Classifiers

We use the same dataset used in [42], because of its size: it contains 500,000 files (250,000 benign samples and 250,000 malware samples), faithfully representing the malware families in the wild and providing a proper setting for an attack comparison. Details about the dataset are provided in Appendix A.

Each sample was run in Cuckoo Sandbox, a malware analysis system, for two minutes per sample. The API call sequences generated by the inspected code were extracted from the JSON report generated by Cuckoo Sandbox. The extracted API call sequences are used as the malware

---

**ALGORITHM 3:** Score-Based or Decision-Based Window Sequence Generation

---

**1 Input**: $f$ - black-box model,

**2** $\boldsymbol{w^m}$ - malicious sequence to perturb, of length $l^m \leq n$, $\boldsymbol{w^b}$ - benign sequence to mimic, of length $l^b \leq n$,

**3** $n$ - size of adversarial sliding window, $D'$ - adversarial vocabulary,

**4** $M_w$ - maximum API modifications per window, $PerturbType$ - benign or random perturbation,

**5** $IterationMethod$ - logarithmic backtracking or linear iteration, $AttackerKnowledge$ - decision or score-based.

**6**

**7 if** $AttackerKnowledge$ is decision-based:

**8**     **while** $(f(\boldsymbol{w^m}) = malicious)$ :

**9**         $(\boldsymbol{w^m}, addedAPIs) = Algorithm2(f, \boldsymbol{w^m}, \boldsymbol{w^b}, n, D', M_w, PerturbType, IterationMethod)\}$

**10 else**: $\#AttackerKnowledge$ is score-based

**11**     **if** $IterationMethod$ is logarithmic backtracking $optimIterationsCount = M_w - \lg n$ , else $optimIterationsCount = M_w$

**12**     **if** $PerturbType$ is benign perturbation:

**13**         $(\boldsymbol{w^m}, addedAPIs) =$
$ScoreMinimizationAlgorithm(f(\boldsymbol{w^m}), optimIterationsCount, addedAPIIndices)$

**14**     **else**: $\#$ $PerturbType$ is random perturbation:

**15**         $(\boldsymbol{w^m}, addedAPIs) =$
$ScoreMinimizationAlgorithm(f(\boldsymbol{w^m}), optimIterationsCount, addedAPIValues)$

**16 return** $(\boldsymbol{w^m}, addedAPIs)$ $\#$ $\boldsymbol{w^m}$ includes the perturbation

---

classifier's features. The samples were run on dozens of Windows 8.1 OS instances on the cloud, since most malware targets the Windows OS. Anti-sandbox malware was filtered to prevent dataset contamination (see Appendix A). After filtering, the final training set size is 360,000 samples, 36,000 of which serve as the validation set. The test set size is 36,000 samples. All sets are balanced between malicious and benign samples.

While some ML-based dynamic analysis cloud services are used by enterprises, e.g., [6], there are no trial versions of commercial products or open source API call-based deep learning malware classifiers available (such commercial products target enterprises and involve supervised server installation). Dynamic models are also unavailable on VirusTotal. In order to compensate for this, we used the malware classifiers detailed in Appendix B and simulated the cloud service use case by deploying Keras [7] models on Amazon cloud using SageMaker [4], and then we queried them by accessing the cloud service.

The API call sequences are split into windows of $k$ API calls each, and each window is classified in turn. Thus, the input of each of the classifiers is a vector of $k = 140$ (larger window sizes, such as $k = 1,000$, didn't improve the classifier's accuracy) API call types with 314 possible values (those monitored by Cuckoo Sandbox, mentioned in [2]).

The implementation and hyperparameters (loss function, dropout, activation functions, etc.) of the target classifiers are described in Appendix B.

On the test set, all of the DNN classifiers achieve over 95% accuracy, and all other classifiers reach over 90% accuracy. The classifiers' false positive rate ranged from 0.5 to 1%.

## 4.2 Attack Performance

### 4.2.1 Attack Performance Metrics

In order to measure the performance of an attack, we consider three factors (by a descending order of importance):

We consider the average number of malware classifier queries the attack performs per adversarial example before it is classified as benign by the malware classifier. The attacker aims to minimize this number, since in cloud scenarios, each query costs money and increases the probability of adversarial attempt detection.

We also consider the *attack effectiveness*, which is the percentage of malicious samples which were correctly classified by the malware classifier for which the adversarial sequence $x^{m*}$ generated by Algorithm 1 was misclassified as benign by the malware classifier. An attack is defined as query-efficient if it has the highest attack effectiveness for a given (fixed) number of queries. (A different approach of a fixed accuracy is computationally expensive to compute.)

Finally, we consider the *attack overhead*, that is, the fraction of API calls which were added (by Algorithm 1) to the malware samples, out of the total number of API calls.

The average length of the API call sequence is: $avg(length(x^m)) \approx 10,000$. We used a maximum of $M_w = 70$ additional API calls per window of $k = 140$ API calls, limiting the perturbation run time overhead (per window and thus per sample) to 50% in the worst case. While not shown here due to space limits, higher $M_w$ values cause higher average attack effectiveness and overhead, and more queries.

Adversarial attacks against images usually try to minimize the number of modified pixels in order to evade human detection of the perturbation. One might claim that such definition of minimal perturbation is irrelevant for API call traces: humans cannot inspect sequences of thousands or millions of APIs, so an attacker can add an unlimited amount of API calls. However, one should bear in mind that a malware aims to perform its malicious functionality as fast as possible. For instance, ransomware usually starts by encrypting the most critical files (e.g., the 'My Documents' folder) first, so even if the user turns off the computer and sends the hard-drive to IT - damage has already been done. The same is true for a key-logger - it aims to send the user passwords to the attacker as soon as possible, so they can be used immediately, before the malware is detected and neutralized. Moreover, adding too many API calls would cause the modified program's profile to become anomalous, making it easier for anomaly detection intrusion detection systems, e.g., systems that measure CPU usage [35] or contain irregular API call subsequences [27] to detect anomalies.

### 4.2.2 Comparison to Previous Work

**Decision-Based Attack Performance**  A comparison of the attack effectiveness and attack overhead of our decision-based attack with logarithmic backtracking transformation and with benign perturbation (Algorithm 1) to BiRand, the more efficient attack used in Dang et al. [24], and to Rosenberg et al. for different attacked malware classifiers is presented in Table 2 (an average of five runs, with 100 queries).

Rosenberg et al. provides state-of-the-art performance for gradient-based attacks against a wide range of classifiers. Dang et al. [24] presented a decision-based attack similar to ours, but limited to non-sequential features. The attack effectiveness of our logarithmic backtracking attack, shown in Table 3 (logarithmic backtracking=yes columns), is identical to the BiRand algorithm presented

Table 2: Decision-Based Attack Performance (100 Queries)

| Classifier Type | Attack Effectiveness [%] Our Decision-Based Attack | Attack Effectiveness [%] BiRand, Dang et al. 2017 [24] | Attack Effectiveness [%] Rosenberg et al. 2018 [42] | Additional API Calls [%] Our Decision-Based Attack | Additional API Calls [%] BiRand, Dang et al. 2017 [24] | Additional API Calls [%] Rosenberg et al. 2018 [42] |
|---|---|---|---|---|---|---|
| LSTM | **62.21** | 39.49 | 51.15 | 22.22 | 31.42 | 17.22 |
| Deep LSTM | **63.62** | 40.38 | 50.80 | 22.71 | 32.12 | 29.51 |
| GRU | **63.35** | 40.21 | 51.16 | 21.47 | 30.36 | 16.09 |
| 1D CNN | **63.63** | 40.39 | 48.93 | 4.10 | 5.80 | 49.21 |
| Logistic Regression | **41.47** | 26.32 | 35.67 | 4.43 | 6.26 | 7.58 |
| Random Forest | **63.24** | 40.14 | 50.87 | 5.20 | 7.35 | 9.40 |
| SVM | **42.59** | 27.04 | 36.27 | 3.82 | 5.40 | 7.19 |
| Gradient Boosted Tree | **41.62** | 26.41 | 36.55 | 13.99 | 19.78 | 27.80 |

in [24] (because both attacks use a similar algorithm).

As can be seen in Table 2, our proposed decision-based attack has the highest effectiveness for all of the malware classifiers tested. Rosenberg et al. provides higher attack effectiveness than our decision-based attack for 2500 queries (not shown due to space limits), but underperforms when the number of queries is being reduced, because there aren't enough queries to build a substitute model with accurate decision boundary.

We see that our attack provides higher attack effectiveness than Dang et al., due to our attack's use of benign perturbation, which was not presented in [24]. When modifying BiRand to use benign perturbation, the results are identical to those obtained by our attack (because both attacks use a similar algorithm; this is not shown due to space limits). In addition, we see that our attack produces a smaller perturbation (adding 25-50% less API calls; see Table 2) than BiRand. This is due to the additional level of backtracking in our attack (lines 19-22 in Algorithm 1). This backtracking is not available in BiRand, which implements a binary search.

As mentioned in Section 4.1, $|TestSet(f)| = 36,000$ samples, and the test set $TestSet(f)$ is balanced, so the attack performance was measured on: $|\{f(\boldsymbol{x_m}) = Malicious | \boldsymbol{x_m} \in TestSet(f)\}| = 18,000$ samples.

We used $k = n$ for Algorithm 2, i.e., the non overlapping sliding window size of the adversary is the same as that used by the target classifier. However, even if this is not the case, the attack effectiveness is not significantly degraded. If $n < k$, the adversary can only modify a subset of the API calls affecting the target classifier, and this subset might not be diverse enough to affect the classification as desired, thereby reducing the attack effectiveness. If $n > k$, the adversaries would keep trying to modify different API calls' positions in Algorithm 2, until they modify the ones impacting the target classifier as well, thereby increasing the attack overhead without affecting the attack effectiveness. For instance, when $n = 100$, $k = 140$, there is an average decrease in attack

effectiveness from 87.96% to 87.94% for an LSTM classifier. Other classifiers have similar behavior, which is not shown due to space limits. The closer $n$ and $k$ are, the better the attack performance.

We used the adversarial vocabulary:
$D' = D - \{ExitWindowsEx(), NtTerminateProcess()\}$, where $D$ is all of the API call types recorded by the malware classifier, so $D'$ does not contain any API type that might harm the code's functionality.

**Score-Based Attack Performance**   There are no published query-efficient adversarial attacks against RNN variants. Attacks that minimize the number of queries exist, but they only work against CNNs [30, 16]. Those attacks aren't relevant, because they don't work with sequence input and discrete values. To address this gap, we used Nevergrad [39], a gradient-free optimization library, to implement discrete sequence input variants of a few state-of-the-art score-based adversarial attacks:

1. SPSA-based attack (Uesato et al. [43]).

2. NES-based attack (Ilyas et al. [30]).

3. GA-based attack (Alzantot et al. [17], Xu et al. [45]).

4. The gradient-based attack (Rosenberg et al. [42]).

We compare theses attacks to our score-based uniform mixing EA attack (described in Section 3.2.3) and to our decision-based attack. We didn't implement the ZOO attack of Chen et al. [21], because it has already been evaluated and was found to be less effective than both SPSA and NES attacks [43].

We used Nevergrad's default arguments for all attacks.

The attack performance (average of five runs) for the LSTM classifier with a fixed budget of 100, 200, and 2500 queries (the attack of Rosenberg et al. requires many queries to accurately build the substitute model required to estimate the gradients per API window [42]) is presented in Table 3 for random *perturbation type* attacks and in Table 4 for benign *perturbation type* attacks (Rosenberg et al. has the same performance in both tables because its perturbations are always determined by the maximum gradient).

The first two lines of each table pertain to our attacks with different *attacker knowledge*. When combining these lines with the *iteration method* values in Tables 3 and 4, one gets our eight previously described attacks (all combinations of: *iteration method*, *attacker knowledge* and *perturbation type*). The *iteration method* values are Linear (for *linear iteration*) and Log (for *logarithmic backtracking*). Other classifiers and budgets (not shown due to space limits) resulted in similar relative trends: a higher budget results in increased attack effectiveness. Note that here we use a fixed number of queries and try to maximize the attack effectiveness for the specified number of queries, since the reverse approach requires higher computational effort and yields the same results.

Our score-based attack variants (in Tables 3 and 4) provide a higher attack effectiveness for all classifiers (this is not discussed further due to space limits) because of the more efficient search algorithms used. The attack of Rosenberg et al. requires using 2255 queries per API call window to build the substitute model accurately enough to reach the performance mentioned in [42] (the rest of the queries, for a a total of 2500, are required to perform the attack itself). Trying to use fewer queries results in a substitute model with inaccurate decision boundary that affect the gradients, and thus the gradient based attack effectiveness.

Table 3: Random *Perturbation Type* Attack Effectiveness Comparison for a Fixed Number of Queries (LSTM Model)

| Number of Queries<br>Logarithmic Backtracking<br>(/BiRand)*Iteration Method* | 100<br>Linear | 200<br>Linear | 2500<br>Linear | 100<br>Log | 200<br>Log | 2500<br>Log |
|---|---|---|---|---|---|---|
| Our Score-Based Attack (Score-Based *Attacker Knowledge*) | **58.75** | **67.59** | **100.00** | **69.28** | **79.71** | **100.00** |
| Our Decision-Based Attack (Decision-Based *Attacker Knowledge*) | 19.86 | 21.25 | 31.43 | 39.49 | 42.25 | 62.50 |
| Rosenberg et al. [42] | 51.15 | 67.11 | 99.99 | 51.15 | 67.11 | 99.99 |
| Uesato et al. [43] | 2.37 | 2.73 | 2.87 | 5.17 | 5.95 | 6.25 |
| Ilyas et al. [30] | 37.50 | 43.14 | 74.94 | 43.78 | 50.37 | 87.50 |
| Alzantot et al. [17], Xu et al. [45] | 54.68 | 62.91 | 100.00 | 62.06 | 71.40 | 100.00 |

Table 4: Benign *Perturbation Type* Attack Effectiveness Comparison for a Fixed Number of Queries (LSTM Model)

| Number of Queries<br>Logarithmic Backtracking<br>(/BiRand)*Iteration Method* | 100<br>Linear | 200<br>Linear | 2500<br>Linear | 100<br>Log | 200<br>Log | 2500<br>Log |
|---|---|---|---|---|---|---|
| Our Score-Based Attack (Score-Based *Attacker Knowledge*) | **71.90** | **82.70** | **100.00** | **84.77** | **97.53** | **100.00** |
| Our Decision-Based Attack (Decision-Based *Attacker Knowledge*) | 41.34 | 44.24 | 46.34 | 62.21 | 63.97 | 87.96 |
| Rosenberg et al. [42] | 51.15 | 67.11 | 99.99 | 51.15 | 67.11 | 99.99 |
| Uesato et al. [43] | 6.56 | 7.55 | 11.47 | 14.31 | 16.46 | 25.00 |
| Ilyas et al. [30] | 66.23 | 76.19 | 81.25 | 77.32 | 88.96 | 94.87 |
| Alzantot et al. [17], Xu et al. [45] | 68.49 | 81.09 | 100.00 | 79.93 | 91.96 | 100.00 |

One might claim that the same substitute model can be used to camouflage more than ten malware samples, resulting in a lower average budget per sample. However, in most cases an attacker would try to modify only a single malware, so it can bypass the detector and perform its malicious functionality. Moreover, even if the average cost per example can be reduced by using the same substitute model, our attack presents a lower limit on the *absolute* number of queries, bypassing a cloud-service that blocks access for a host performing too many queries in a short amount of time in order to thwart adversarial efforts [22, 19]. The more efficient the attack, the less chances there are for it to be mitigated by this approach.

The performance of our linear iteration attack, shown in Table 3 (logarithmic backtracking=no columns), is identical to the performance of the SeqRand algorithm presented in [24] (because both attacks use the same algorithm).

The attack overhead of all attacks is similar: about 30%, or 40 API calls, per window. Since a classifier with an API window size of $k = 100$ provides roughly the same accuracy as with $k = 140$

used here (96.76% vs. 97.61% with the same FP rate for the LSTM classifier), the success of these attacks is due to the perturbation and *not* because API sequences were split into two windows due to the added API calls.

As can be seen, the attacks of Uesato et al. [43] and Ilyas et al. [30] have low effectiveness. This is due to the fact that those attacks are not suitable for discrete values of API call types and indices.

In contrast, we see that our uniform mixing EA score-based attack has higher attack effectiveness, for a fixed number of queries, even when used for discrete input (API calls or position indices). This is due to the fact that the transformations used by EA work with discrete sequences: mutation (random perturbation) in existing adversarial candidates and crossover between several candidates. In our EA score-based attack, we don't use crossover, which might make sense for the NLP domain (e.g., for compound sentences) but not for API call sequences, where each program has its own business logic. The self-adaptive search used by our EA score-based attack also explains why it outperforms all other score-based attack variants and has better attack effectiveness than the gradient-based attack used in [42] with the same number of queries. Our proposed score-based attack outperforms existing methods because it maximizes the attack effectiveness for a fixed number of queries. Note that the number of queries is per sliding window and not per executable.

Based on the average malicious sequence length, $avg(length(\boldsymbol{x^m})) \approx 10,000$, and the adversarial sliding window size, $k = 140$, the average absolute number of queries per malware executable is ~10,000.

As expected, the benign perturbation effect on the decision-based attack effectiveness is the most significant, since without it, the API types are random.

While our decision-based attack effectiveness is 10% lower than the most effective score-based attacks when using the same budget, it doesn't require knowledge of the target classifier's confidence score, making it the only viable attack in some black-box scenarios.

## 4.3 Defenses and Mitigation Techniques

To the best of our knowledge, there is currently no published and evaluated method to make a sequence-based RNN model resistant to adversarial sequences, beyond a brief mention of adversarial training as a defense method [17, 33]. Adversarial training [26] is the method of adding adversarial examples, with their non-perturbed label, to the training set of the classifier. The reason is since adversarial examples are usually out-of-distribution samples, inserting them into the training set would cause the classifier to learn the entire training set distribution, including the adversarial examples.

Adversarial training has several limitations:

1. It provides a varying level of robustness, depending on the adversarial examples used.

2. It requires a dataset of adversarial examples to train on. Thus, it has limited generalization against novel adversarial attacks.

3. It requires retraining the model, incurring significant overhead.

We ran the adversarial attacks, both score-based and decision-based variants (Section 1), with and without benign perturbation (Section 3.2.2) on the training set, as suggested in [34]. For each column in Tables 3 and 4, we generated 14,000 malicious adversarial examples (50% generated by the black-box attack and 50% by the white-box attack), which replaced 14,000 malicious samples

in the original training set. Other sizes (smaller or larger) resulted in reduced detection rate of the pre-trained classifier for non-adversarial samples. The adversarial examples were generated using the same configuration (score/decision-based, random/benign perturbation, number of queries to generate) as the evaluated attack. The results were the same across all attack types: The attack effectiveness remains the same, while the attack overhead and number of queries were increased by 10-15%, on average. This is due to the fact that adversarial training is less effective against random attacks like ours, because a different stochastic adversarial sequence is generated every time, making it challenging for the classifier to generalize from one adversarial sequence to another.

More effective RNN defense methods, including domain specific methods, e.g., systems that measure CPU usage [35], contain irregular API call subsequences [27] (such as the no-op API calls used in this paper), or otherwise assess the plausibility of our attack [38], in order to detect adversarial examples, will be a part of our future work.

# 5    Conclusions and Future Work

In this paper, we presented the first black-box attack (based on the target classifier's predicted class, with and without its confidence score, to fit adversary's limited knowledge) that generates *adversarial sequences while minimizing the number of queries* for the target classifier, reducing the number of queries by more than 10 times with minimal loss of attack effectiveness in comparison to the state of the art attack ([42]). This query-efficient approach makes our attack suited to attack cloud models where a large amount of queries cost money and raise suspicion of an attack, failing previous attacks.

We demonstrated those attacks against API call sequence-based malware classifiers and verified the attack effectiveness for all relevant common classifiers: RNN variants, feedforward networks, and traditional machine learning classifiers. These are the first query-efficient attacks effective against RNN variants and not just CNNs.

We also evaluated our attacks against four variants of state-of-the-art score-based query-efficient attacks, modified to fit discrete sequence input, and showed that our attacks are equal or outperform all of them.

Finally, we demonstrated that our attacks are effective even when multiple feature types, including non-sequential ones, are used (Appendix D).

While this paper focuses on API calls and printable strings as features, the proposed attacks are valid for every modifiable feature type, sequential or not. Furthermore, our attack is generic and can be applied to other domains, like text analysis (using word sequences instead of API calls), as would be demonstrated in our future work.

Our future work will focus on developing domain-specific and domain-agnostic defense mechanisms against such attacks and analyzing additional self-adaptive algorithms to find more query-efficient attacks, while evaluating them on limited knowledge scenarios (e.g., unknown API calls window size, etc.).

# References

[1] Amazon Machine Learning. https://aws.amazon.com/machine-learning, 2019. Accessed: 2019-09-26.

[2] Cuckoo Sandbox Hooked APIs and Categories. https://github.com/cuckoosandbox/cuckoo/wiki/Hooked-APIs-and-Categories, 2019. Accessed: 2019-08-24.

[3] Cylance, I Kill You! https://skylightcyber.com/2019/07/18/cylance-i-kill-you, 2019. Accessed: 2019-08-24.

[4] Deploy trained Keras or TensorFlow models using Amazon SageMaker. https://aws.amazon.com/blogs/machine-learning/deploy-trained-keras-or-tensorflow-models-using-amazon-sagemaker/, 2019. Accessed: 2019-12-14.

[5] Google Cloud Prediction. https://cloud.google.com/prediction/, 2019. Accessed: 2019-09-26.

[6] Joe Sandbox ML. https://www.joesecurity.org/joe-sandbox-ML, 2019. Accessed: 2019-09-26.

[7] Keras. https://keras.io/, 2019. Accessed: 2019-09-26.

[8] Microsoft ATP. https://www.microsoft.com/security/blog/2018/02/14/how-artificial-intelligence-stopped-an-emotet-outbreak/, 2019. Accessed: 2019-09-26.

[9] SciKit Learn. http://scikit-learn.org/stable/, 2019. Accessed: 2019-09-26.

[10] Scikit Learn Decision Tree Categorial Variable. https://roamanalytics.com/2016/10/28/are-categorical-variables-getting-lost-in-your-random-forests/, 2019. Accessed: 2019-09-26.

[11] SentinelOne. https://www.sentinelone.com/insights/endpoint-protection-platform-datasheet/, 2019. Accessed: 2019-09-26.

[12] VirusTotal. https://www.virustotal.com/, 2019. Accessed: 2019-09-26.

[13] XGBoost. https://github.com/dmlc/xgboost/, 2019. Accessed: 2019-09-26.

[14] Yara Rules. https://github.com/Yara-Rules/rules, 2019. Accessed: 2019-09-26.

[15] Rakshit Agrawal, Jack W. Stokes, Mady Marinescu, and Karthik Selvaraj. Robust neural malware detection models for emulation sequence learning. *CoRR*, abs/1806.10741, 2018.

[16] Moustafa Alzantot, Yash Sharma, Supriyo Chakraborty, and Mani B. Srivastava. Genattack: Practical black-box attacks with gradient-free optimization. *CoRR*, abs/1805.11090, 2018.

[17] Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani B. Srivastava, and Kai-Wei Chang. Generating natural language adversarial examples. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2890–2896. Association for Computational Linguistics, 2018.

[18] Hyrum S. Anderson and Phil Roth. EMBER: an open dataset for training static PE malware machine learning models. *CoRR*, abs/1804.04637, 2018.

[19] Duen Horng Chau, Carey Nachenberg, Jeffrey Wilhelm, Adam Wright, and Christos Faloutsos. Polonium: Tera-scale graph mining for malware detection. In *Acm sigkdd conference on knowledge discovery and data mining*, 2010.

[20] Tong Che, Yanran Li, Ruixiang Zhang, R. Devon Hjelm, Wenjie Li, Yangqiu Song, and Yoshua Bengio. Maximum-likelihood augmented discrete generative adversarial networks. *CoRR*, abs/1702.07983, 2017.

[21] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security - AISec 17*. ACM Press, 2017.

[22] Steven Chen, Nicholas Carlini, and David Wagner. Stateful detection of black-box adversarial attacks, 2019.

[23] Duc-Cuong Dang and Per Kristian Lehre. Self-adaptation of mutation rates in non-elitist populations. In Julia Handl, Emma Hart, Peter R. Lewis, Manuel López-Ibáñez, Gabriela Ochoa, and Ben Paechter, editors, *Parallel Problem Solving from Nature – PPSN XIV*, pages 803–813, Cham, 2016. Springer International Publishing.

[24] Hung Dang, Yue Huang, and Ee-Chien Chang. Evading classifiers by morphing in the dark. In Bhavani M. Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, pages 119–133. ACM, 2017.

[25] Jennifer G. Dy and Andreas Krause, editors. *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*. PMLR, 2018.

[26] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and Harnessing Adversarial Examples. *International Conference on Learning Representations (ICLR)*, December 2015.

[27] Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick D. McDaniel. On the (statistical) detection of adversarial examples. *ArXiv e-prints*, abs/1702.06280, 2017.

[28] Weiwei Hu and Ying Tan. Black-box attacks against RNN based malware detection algorithms. *ArXiv e-prints*, abs/1705.08131, 2017.

[29] Ling Huang, Anthony D. Joseph, Blaine Nelson, Benjamin I. P. Rubinstein, and J. D. Tygar. Adversarial machine learning. In Yan Chen, Alvaro A. Cárdenas, Rachel Greenstadt, and Benjamin I. P. Rubinstein, editors, *Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence, AISec 2011, Chicago, IL, USA, October 21, 2011*, pages 43–58. ACM, 2011.

[30] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In Dy and Krause [25], pages 2142–2151.

[31] Jeremy Z. Kolter and Marcus A. Maloof. Learning to detect and classify malicious executables in the wild. *J. Mach. Learn. Res.*, 7:2721–2744, 2006.

[32] Matt J. Kusner and José Miguel Hernández-Lobato. GANS for sequences of discrete elements with the gumbel-softmax distribution. *CoRR*, abs/1611.04051, 2016.

[33] Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. Textbugger: Generating adversarial text against real-world applications. *CoRR*, abs/1812.05271, 2018.

[34] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[35] Robert Moskovitch, Shay Pluderman, Ido Gus, Dima Stopel, Clint Feher, Yisrael Parmet, Yuval Shahar, and Yuval Elovici. Host based intrusion detection using machine learning. In *IEEE International Conference on Intelligence and Security Informatics, ISI 2007, New Brunswick, New Jersey, USA, May 23-24, 2007, Proceedings*, pages 107–114. IEEE, 2007.

[36] Nicolas Papernot, Patrick McDaniel, Ananthram Swami, and Richard Harang. Crafting adversarial input sequences for recurrent neural networks. In *MILCOM 2016 - 2016 IEEE Military Communications Conference*. IEEE, nov 2016.

[37] Feargus Pendlebury, Fabio Pierazzi, Roberto Jordaney, Johannes Kinder, and Lorenzo Cavallaro. TESSERACT: Eliminating experimental bias in malware classification across space and time. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 729–746, Santa Clara, CA, August 2019. USENIX Association.

[38] Fabio Pierazzi, Feargus Pendlebury, Jacopo Cortellazzi, and Lorenzo Cavallaro. Intriguing properties of adversarial ML attacks in the problem space. In *2020 IEEE Symposium on Security and Privacy, SP 2020, San Francisco, CA, USA, May 18-21, 2020*, pages 1332–1349. IEEE, 2020.

[39] J. Rapin and O. Teytaud. Nevergrad - A gradient-free optimization platform. https://GitHub.com/FacebookResearch/Nevergrad, 2018.

[40] Ihai Rosenberg, Asaf Shabtai, Yuval Elovici, and Lior Rokach. Adversarial learning in the cyber security domain, 2020.

[41] Ishai Rosenberg, Shai Meir, Jonathan Berrebi, Ilay Gordon, Guillaume Sicard, and Eli David. Generating end-to-end adversarial examples for malware classifiers using explainability. In *The 2020 International Joint Conference on Neural Networks (IJCNN 2020)*, 2020.

[42] Ishai Rosenberg, Asaf Shabtai, Lior Rokach, and Yuval Elovici. Generic black-box end-to-end attack against state of the art API call based malware classifiers. In Michael Bailey, Thorsten Holz, Manolis Stamatogiannakis, and Sotiris Ioannidis, editors, *Research in Attacks, Intrusions, and Defenses - 21st International Symposium, RAID 2018, Heraklion, Crete, Greece, September 10-12, 2018, Proceedings*, volume 11050 of *Lecture Notes in Computer Science*, pages 490–510. Springer, 2018.

[43] Jonathan Uesato, Brendan O'Donoghue, Pushmeet Kohli, and Aäron van den Oord. Adversarial risk and the dangers of evaluating against weak attacks. In Dy and Krause [25], pages 5032–5041.

[44] David Wagner and Paolo Soto. Mimicry attacks on host-based intrusion detection systems. In *Proceedings of the 9th ACM conference on Computer and communications security - CCS '02*. ACM Press, 2002.

[45] Weilin Xu, Yanjun Qi, and David Evans. Automatically evading classifiers: A case study on PDF malware classifiers. In *23rd Annual Network and Distributed System Security Symposium, NDSS 2016, San Diego, California, USA, February 21-24, 2016*. The Internet Society, 2016.

[46] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial nets with policy gradient. In Satinder P. Singh and Shaul Markovitch, editors, *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 2852–2858. AAAI Press, 2017.

[47] Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, and Lawrence Carin. Adversarial feature matching for text generation. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 4006–4015. PMLR, 2017.

[48] Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Texygen: A benchmarking platform for text generation models. In Kevyn Collins-Thompson, Qiaozhu Mei, Brian D. Davison, Yiqun Liu, and Emine Yilmaz, editors, *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pages 1097–1100. ACM, 2018.
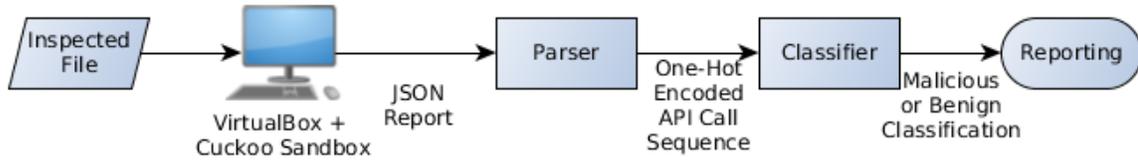
Figure 2: Overview of the Malware Classification Process

# A  Tested Dataset

We used identical implementation details (e.g., dataset, classifiers' hyperparameters, etc.) as Rosenberg et al. [42], so the attacks can be compared. The details are provided here for the reader's convenience.

An overview of the malware classification process is shown in Figure 2 (taken from [42]).
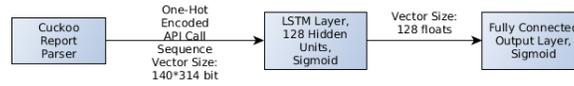
The dataset used is large and includes the latest malware variants, such as the Cerber and Locky ransomware families. Each malware type (ransomware, worms, backdoors, droppers, spyware, PUAs, and viruses) has the same number of samples, to prevent prediction bias towards the majority class. 20% of the malware families (such as the NotPetya ransomware family) were only used in the test set to assess generalization to an unseen malware family. 80% of the malware families (such as the Virut virus family) were distributed between the training and test sets, to determine the classifier's ability to generalize to samples from the same family. The temporal difference between the training set and the test set is six months (i.e., all training set samples are older than the test set samples, because using a non-time-aware split may cause a data leak [37]) based on VirusTotal's 'first seen' date.

The ground truth labels of the dataset were determined by VirusTotal [12], an online scanning service, which contains more than 60 different security products. A sample with 15 or more positive (i.e., malware) classifications from the 60 products is considered malicious. A sample with zero positive classifications is labeled as benign. All samples with 1-14 positives were omitted to prevent false positive contamination of the dataset. Family labels for dataset balancing were taken from Kaspersky Anti-Virus classifications.
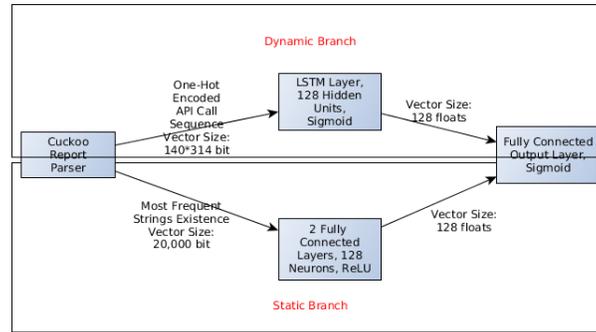
It is crucial to prevent dataset contamination by malware that detects whether the malware is running in a Cuckoo Sandbox (or on virtual machines) and if so, quits immediately to prevent reverse engineering efforts. In those cases, the sample's label is malicious, but its behavior recorded in Cuckoo Sandbox (its API call sequence) isn't, due to the malware's anti-forensic capabilities. To prevent such contamination of the dataset, two countermeasures were used:

1. Considering only API call sequences with more than 15 API calls, omitting malware that detects a virtual machine (VM) and quits.

2. Applying Yara rules [14] to find samples trying to detect sandbox programs, such as Cuckoo Sandbox, and omitting all such samples.

One might argue that the evasive malware that applies such anti-VM techniques is extremely challenging and relevant, however in this paper we focus on adversarial attacks. Such attacks are generic enough to work for those evasive types of malware as well, assuming that other mitigation

(a) Dynamic Classifier Architecture



(b) Hybrid Classifier Architecture

Figure 3: Classifier Architecture Overview

techniques (e.g., anti-anti-VM), would be applied. After this filtering and balancing of the benign samples, about 400,000 valid samples remained. The final training set size is 360,000 samples, 36,000 of which serve as the validation set. The test set size is 36,000 samples. All sets are balanced between malicious and benign samples. Due to hardware limitations, a subset of the dataset was used (54,000 training samples and test and validation sets of 6,000 samples each). The dataset was representative and maintained the same distribution as mentioned above.

# B   Tested Malware Classifiers

As mentioned in Section 4.1, we used the malware classifiers from Rosenberg et al. [42], since many classifiers are covered, allowing us to evaluate the attack effectiveness against many classifier types. The maximum input sequence length was limited to $k = 140$ API calls, since longer sequence lengths, e.g., $k = 1,000$, had no effect on the accuracy, and shorter sequences were padded with zeros. A zero stands for a null/dummy value API in our one-hot encoding. Longer sequences are split into windows of $k$ API calls each, and each window is classified in turn. If any window is malicious, the entire sequence is considered malicious. Thus, the input of all of the classifiers is a vector of $k = 140$ API call types in one-hot encoding, using 314 bits, since there were 314 monitored API call types in the Cuckoo reports for the dataset. The output is a binary classification: malicious or benign. An overview of the LSTM architecture is shown in Figure 3a.

The Keras [7] implementation was used for all neural network classifiers, with TensorFlow used for the backend. XGBoost [13] and scikit-learn [9] were used for all other classifiers.

The loss function used for training was binary cross-entropy. The Adam optimizer was used for

all of the neural networks. The output layer was fully connected with sigmoid activation for all neural networks. For neural networks, a rectified linear unit, $ReLU(x) = max(0, x)$, was chosen as an activation function for the input and hidden layers due to its rapid convergence compared to $sigmoid()$ or tanh(), and dropout was used to improve the generalization potential of the network. A batch size of 32 samples was used. The classifiers also have the following classifier specific hyperparameters:

- DNN - two fully connected hidden layers of 128 neurons, each with ReLU activation and a dropout rate of 0.2.

- CNN - 1D ConvNet with 128 output filters, a stride length of one, a 1D convolution window size of three, and ReLU activation, followed by a global max pooling 1D layer and a fully connected layer of 128 neurons with ReLU activation and a dropout rate of 0.2.

- RNN, LSTM, GRU, BRNN, BLSTM, and bidirectional GRU - a hidden layer of 128 units, with a dropout rate of 0.2 for inputs and recurrent states.

- Deep LSTM and BLSTM - two hidden layers of 128 units, with a dropout rate of 0.2 for inputs and recurrent states in both layers.

- Linear SVM and logistic regression classifiers - a regularization parameter of C=1.0 and an L2 norm penalty.

- Random forest classifier - 10 decision trees with unlimited maximum depth and the Gini criteria for choosing the best split.

- Gradient boosted decision tree - up to 100 decision trees with a maximum depth of 10 each.

The classifiers' performance was measured using the accuracy ratio, which gives equal importance to both false positives and false negatives (unlike precision or recall). The false positive rate of the classifiers varied between 0.5-1%. The false positive rate was chosen to be on the high end of production systems. A lower false positive rate would mean lower recall either, due to the trade-off between false positive rate and recall, thereby making our attacks even more effective.

The performance of the classifiers is shown in Table 5. The accuracy was measured on the test set, which contains 36,000 samples.

Table 5: Classifier Performance

| Classifier Type | Accuracy (%) |
|---|---|
| LSTM | 98.26 |
| Deep LSTM | 97.90 |
| GRU | 97.32 |
| Bidirectional GRU | 98.04 |
| 1D CNN | 96.42 |
| Random Forest | 91.90 |
| SVM | 86.18 |
| Logistic Regression | 89.22 |
| Gradient Boosted Decision Tree | 91.10 |

Table 6: Benign Perturbation Attack Performance

| GAN Type | Attack Effectiveness [%] | Added API Calls [%] | Queries Used |
|---|---|---|---|
| None (Random Perturbation) | 21.25 | 27.94 | 119.40 |
| SeqGAN [46] | **89.39** | **12.82** | **17.73** |
| TextGAN [47] | 74.53 | 16.74 | 30.58 |
| GSGAN [32] | 88.19 | 14.06 | 20.43 |
| MaliGAN [20] | 86.67 | 15.12 | 22.74 |

As can be seen in Table 5, the LSTM variants are the best malware classifiers, in terms of accuracy, and, as shown in Section 4.2, BLSTM is also one of the classifiers most resistant to the proposed attack.

## C    Benign Perturbation GAN Comparison

To implement the benign perturbation GAN, we tested several GAN types, using Texygen [48] with its default parameters. We use maximum likelihood estimation (MLE) training as the pretraining process for all baseline models except GSGAN, which requires no pretraining. In pretraining, we first train 80 epochs for a generator, and then train 80 epochs for a discriminator. The adversarial training comes next. In each adversarial epoch, we update the generator once and then update the discriminator for 15 mini-batch gradients. Due to memory limitations, we generated only one sliding window of 140 API calls, each with 314 possible API call types, in each iteration (that is, generating $w^b$ and not $x^b$ as described in Algorithm 1).

We tested several GAN implementations with discrete sequence output. We trained our GAN using a benign hold-out set (3,000 sequences). Next, we run Algorithm 1 (logarithmic backtracking transformation with benign perturbation) on the 3,000 API call traces generated by the GAN. Finally, we used the benign test set (3,000 sequences) as the GAN's test set. The results for the LSTM classifier are shown in Table 6 (the results for other classifiers, which are not shown due to space limits, are similar).

We can see from the results presented in the table that SeqGAN outperforms all of the other models in all of the measured factors, due to its RL-based ability to pass gradient updates between the generator and discriminator parts of the GAN. We also see that, as expected, a random perturbation is less effective than a benign perturbation, regardless of the type of GAN used.

## D    Handling Multiple Feature Types and Hybrid Classifiers

Combining several types of features can make the classifier more resistant to adversarial examples against a specific feature type. For instance, some real-world next generation anti-malware products are hybrid classifiers, combining both static and dynamic features for a better detection rate. An extension of our attack, enabling it to handle hybrid classifiers, is straightforward: attacking *each feature type in turn* using Algorithm 1. If the attack against a feature type fails, we continue and

attack the next feature type with the modified binary until a benign classification by the target model is achieved or all feature types have been (unsuccessfully) attacked. We used the same hybrid malware classifier specified in Appendix C, for which the input consists of both an API call sequence and the most frequent 20,000 printable strings inside the PE file as Boolean features (exist or not).

While there are more complex static features (e.g., [18]), we chose printable strings, easy to modify features that have been used by many classifiers [31], as a concrete example of the multi-feature use case, to show that the suggested attack works not only against RNNs, but also against other classifiers, making it more generic.

We evaluated the performance of our decision-based, linear iteration, benign perturbation attack. When attacking only the API call sequences using the hybrid classifier, without modifying the static features of the sample, the attack effectiveness decreases to 23.76%. This is much lower than the attack effectiveness of 89.67% obtained for a classifier trained only on the dynamic features, meaning that the attack was mitigated by the use of additional static features. When attacking only the printable string features (again, assuming that the attacker has the knowledge of $D' = D$, which contains the printable strings being used as features by the hybrid classifier), the attack effectiveness is 28.25%. This is much lower than the attack effectiveness of 88.31% obtained for a classifier trained only on the static features. Finally, the multi-feature attack's effectiveness for the hybrid model was 90.06%. Other types of classifiers and attacks provided similar results. They are not presented here due to space limits.

To summarize, we have shown that while the use of hybrid models decreases the specialized attacks' effectiveness, our suggested hybrid attack performs well, with high attack effectiveness. While not shown due to space limits, the attack overhead isn't significantly affected.

# E  Tested Hybrid Malware Classifiers

As mentioned in Appendix D, we used the hybrid malware classifier used in [42], with printable strings inside a PE file as our static features. Strings can be used, e.g., to statically identify loaded DLLs and called functions, and recognize modified file paths and registry keys, etc. Our architecture for the hybrid classifier, shown in Figure 3b, is:

1. A static branch that contains an input vector of 20,000 Boolean values: for each of the 20,000 most frequent strings in the entire dataset, do they appear in the file or not? This is analogous to a similar procedure used in NLP which filters the least frequent words in a language.

2. A dynamic branch that contains an input vector of 140 API calls (each of which is one-hot encoded) inserted into an LSTM layer of 128 units and a sigmoid activation function, with a dropout rate of 0.2 for inputs and recurrent states. This vector is inserted into two fully connected layers with 128 neurons, a ReLU activation function, and a dropout rate of 0.2 each.

The 256 outputs of both branches are inserted into a fully connected output layer with a sigmoid activation function. Therefore, the input of the classifier is a vector containing 20,000 Boolean values and 140 one-hot encoded API call types, and the output is malicious or benign classification. The training set size was reduced by half in comparison to the training set specified in Appendix A (while keeping the same dataset structure) due to the larger memory requirements for training a classifier with more features. All other hyperparameters are the same as those mentioned in Appendix B.

A classifier using only the dynamic branch (Figure 3a) achieves 92.48% accuracy on the test set (this is different from the results presented in Table 5, due to the smaller training set), a classifier using only the static branch attains 96.19% accuracy, and a hybrid model that uses both branches (Figure 3b) obtains 96.94% accuracy, meaning that using multiple feature types improves the accuracy.