

Adversarial EXEmples: A Survey and Experimental Evaluation of Practical Attacks on Machine Learning for Windows Malware Detection

LUCA DEMETRIO, Università degli studi di Cagliari, ITA

SCOTT E. COULL, FireEye, Inc.

BATTISTA BIGGIO, Università degli studi di Cagliari, ITA and Pluribus One, ITA

GIOVANNI LAGORIO, Università degli Studi di Genova, ITA

ALESSANDRO ARMANDO, Università degli Studi di Genova, ITA

FABIO ROLI, Università degli Studi di Cagliari, ITA and Pluribus One, ITA

Recent work has shown that adversarial Windows malware samples - referred to as adversarial *EXEmples* in this paper - can bypass machine learning-based detection relying on static code analysis by perturbing relatively few input bytes. To preserve malicious functionality, previous attacks either add bytes to existing non-functional areas of the file, potentially limiting their effectiveness, or require running computationally-demanding validation steps to discard malware variants that do not correctly execute in sandbox environments. In this work, we overcome these limitations by developing a unifying framework that does not only encompass and generalize previous attacks against machine-learning models, but also includes three novel attacks based on practical, functionality-preserving manipulations to the Windows Portable Executable (PE) file format. These attacks, named *Full DOS*, *Extend* and *Shift*, inject the adversarial payload by respectively manipulating the DOS header, extending it, and shifting the content of the first section. Our experimental results show that these attacks outperform existing ones in both white-box and black-box scenarios, achieving a better trade-off in terms of evasion rate and size of the injected payload, while also enabling evasion of models that have been shown to be robust to previous attacks. To facilitate reproducibility of our findings, we open source our framework and all the corresponding attack implementations as part of the `secml-malware` Python library. We conclude this work by discussing the limitations of current machine learning-based malware detectors, along with potential mitigation strategies based on embedding domain knowledge coming from subject-matter experts directly into the learning process.

CCS Concepts: • **Computing methodologies** → **Machine Learning**; • **Security and privacy** → *Malware and its mitigation*.

Additional Key Words and Phrases: adversarial examples, malware detection, evasion, semantics-invariant manipulations

ACM Reference Format:

Luca Demetrio, Scott E. Coull, Battista Biggio, Giovanni Lagorio, Alessandro Armando, and Fabio Roli. 2020. Adversarial EXEmples: A Survey and Experimental Evaluation of Practical Attacks on Machine Learning for Windows Malware Detection. In *Proceedings of ACM TO EDIT conference*. ACM, New York, NY, USA, 31 pages. <https://doi.org/TOEDIT>

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 ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

Manuscript submitted to ACM

1 INTRODUCTION

Machine learning (ML) has become an important aspect of modern cybersecurity due to its ability to detect new threats far earlier than signature-based defenses. While many cybersecurity companies use machine learning models¹²³⁴ in their respective product offerings, creating and maintaining these models often represents a significant cost in terms of expertise and labor in developing useful features to train on, particularly when we consider that each new file type may require a completely different set of features to provide meaningful classification. With this in mind, researchers have recently proposed end-to-end deep learning models that operate directly on the raw bytes of the input files and automatically learn useful feature representations during training, without external knowledge from subject-matter experts. Several byte-based malware detection models for Windows PE files, for example, have demonstrated efficacy that is competitive with traditional ML models [9, 29] (Sect. 2).

While the use of end-to-end deep learning makes it easy to create new models for a variety of file types by simply exploiting the vast number of labeled samples available to such organizations, it also opens up the possibility of attacking these models using *adversarial evasion* techniques popularized in the image classification space. In particular, recent work has shown how an attacker can create what we call here *adversarial EXEmple*s, i.e., Windows malware samples carefully perturbed to evade learning-based detection while preserving malicious functionality [2, 7, 11, 12, 20, 23, 31, 33].

Unlike adversarial evasion attacks in other problem areas, such as image classification, manipulating malware while simultaneously preserving its malicious payload can be difficult to accomplish. In particular, each perturbation made to the input bytes during the attack process may lead to changes in the underlying syntax, semantics, or structure that could prevent the binary from executing its intended goal. To address this problem, the attacker can take one of two approaches: apply invasive perturbations and use dynamic analysis methods (e.g. emulation) to ensure that functionality of the binary is not compromised [7, 32], or focus the perturbation on areas of the file that do not impact functionality [11, 12, 20, 23, 33] (e.g. appending bytes). Naturally, this leads to a trade-off between strong yet time-consuming attacks on one extreme, and weaker but more computationally-efficient attacks on the other.

In this work, we overcome these limitations by proposing a unifying framework, called **RAMEN** (Sect. 3), built on top of a family of *practical* manipulations to the Windows Portable Executable (PE) file format that can alter the structure of the input malware without compromising its semantics. Our framework encompasses and generalizes previously-proposed attacks against learning-based Windows malware detectors based on static code analysis, including both white-box attacks that exploit full knowledge of the target algorithm, and black-box attacks that only require query access to it. The practical, functionality-preserving manipulations defined in our framework are not limited to perturbing bytes at the end of malware programs and do not require computationally-demanding validation steps during the attack optimization, thereby overcoming the limitations of existing attacks. In particular, we encode three novel practical manipulations that exploit the ambiguity in the specifications of the Windows PE file format: *Full DOS* that edit all the available bytes inside the DOS header; *Extend*, which enlarges the DOS header, thus enabling manipulation of these extra DOS bytes; and *Shift*, which shifts the content of the first section, carving additional space for the adversarial payload.

Our experimental results (Sect. 4) show that these attacks outperform existing ones in both white-box and black-box attack scenarios against different machine-learning models, deep network architectures, activation functions (i.e. linear

¹<https://www.sophos.com/products/intercept-x/tech-specs.aspx>

²<https://www.fireeye.com/blog/products-and-services/2018/07/malwareguard-fireeye-machine-learning-model-to-detect-and-prevent-malware.html>

³<https://www.kaspersky.com/enterprise-security/wiki-section/products/machine-learning-in-cybersecurity>

⁴<https://www.avast.com/technology/ai-and-machine-learning>

vs. non-linear models), and training regimes. In particular, our *Extend* and *Shift* attacks enable evading some models that are not affected by previously-proposed attacks, while generally achieving a better trade-off in terms of evasion rate and size of the injected payload; they create fully-functional, evasive malware by perturbing roughly 2% of the input bytes against most of the considered classifiers.

An additional finding from our experimental analysis is that, while dataset size and activation functions do not seem to play a significant role in improving adversarial robustness, model architecture does, at least to some extent, with all attacks working well against Raff et al.'s MalConv classifier [29] and only content-shifting attacks working well against Coull et al.'s classifier [9], possibly due to the importance of spatial locality in its design. This identifies a promising line of research towards strengthening models against adversarial attacks through inclusion of additional structure in the training process. We conclude the paper by discussing related work (Sect. 5) along with the limitations of our methodology (Sect. 6), and promising research directions to improve robustness of learning-based Windows malware detectors against adversarial attacks (Sect. 7). Besides considering network architectures that exploit spatial locality, we discuss other potential strategies to embed external domain knowledge directly into the learning process (e.g., via suitable constraints and loss functions) with the goal of learning more meaningful and robust representations from data [24]. We believe that this novel learning paradigm may help significantly improve adversarial robustness of such models, while at the same time exploiting knowledge from domain experts in an efficient manner.

To summarize, we highlight our contributions below.

- We propose RAMEN, a general framework for expressing white-box and black-box adversarial attacks on learning-based Windows malware detectors based on static code analysis.
- We propose three novel attacks based on practical, functionality-preserving manipulations, named *Full DOS*, *Extend* and *Shift*, which improve the trade-off between the probability of evasion and the amount of manipulated bytes in both white-box and black-box attack settings.
- We release the implementations of all the aforementioned white-box and black-box attacks encompassed by our framework (including previous attacks, *Extend* and *Shift*), as an open-source project available at https://github.com/zangobot/secml_malware.
- We identify promising future research directions towards improving adversarial robustness of learning-based Windows malware detectors leveraging static code analysis.

2 BACKGROUND

Before diving into the details of our proposed attack framework, we first provide some necessary background on the Windows Portable Executable (PE) file format (Section 2.1) and the malware classifiers that we will examine in our experiments (Section 2.2).

2.1 Executable File Format

The *Windows Portable Executable (PE)*⁵ format specifies how executable programs are stored as a file on disk. The OS loader parses this structure and maps the code and data into memory, following the directives specified by the header of the file. We show the components of the format in Figure 1.

DOS Header and Stub. The DOS header contains metadata for loading the executable inside a DOS environment, while the DOS Stub is made up of few instructions that will print “*This program cannot be run in DOS mode*” if executed

⁵<https://docs.microsoft.com/en-us/windows/win32/debug/pe-format>

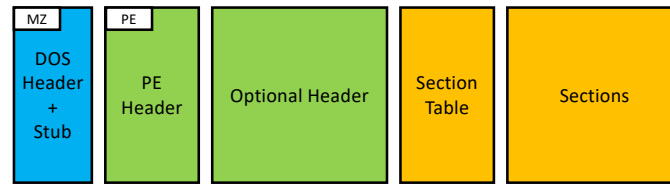


Fig. 1. The Windows PE file format. Each colored section describes a particular characteristic of the program.

inside a DOS environment. These two components have been kept to maintain compatibility with older Microsoft operating systems. From the perspective of a modern application, the only relevant locations contained inside the DOS Header are (i) the magic number MZ, a two-byte signature of the file, and (ii) the four-byte integer at offset $0x3c$, that works as a pointer to the actual header. If one of these two values is altered for some reason, the program is considered corrupted, and it will not be executed by the OS.

PE Header. It is the real header of the program, and it contains the magic number PE and the characteristics of the executable, such as the target architecture that can run the program, the size of the header and the attributes of the file.

Optional Header. Not optional for executables and DLLs, it contains the information needed by the OS for loading the binary into memory. Among these fields, the Optional Header specifies the (i) *file alignment*, that acts as a constraint on the structure of the executable since each section of the program must start at an offset multiple to that field, and the (ii) *size of headers* that specifies the amount of bytes that are reserved to all the headers of the programs, and it must be a multiple of the *file alignment*. Lastly, the optional header contains offsets that point to useful structures, like the Import Table needed by the OS for resolving dependencies, the Export Table to find functions that can be referenced by other programs, and more.

Section Table. It is a list of entries that indicates the characteristics of each section of the program. Each section entry has a name, an offset to the location inside the binary, a virtual address where the content should be mapped in memory, and the characteristics of such content (i.e. is read-only, write-only, or it is executable, and more).

Sections. These are contiguous chunks of bytes, loaded in memory by the loader after while parsing the Section Table. To name a few, there is the code of the program (*.text* section), initialized data (*.data*), read-only constants (*.rdata*), and counting. To maintain the alignment specified inside the Optional Header, these sections might be zero-padded to match the format constraint. It is clear that, even without executing the program contained inside a file, it is possible to infer some useful information from its headers, imports, exports, and sections.

2.2 Malware Classifiers

Here, we describe two recent byte-based convolutional neural network models for malware detection. Both take as input the raw bytes from the Windows PE file on disk, use an embedding layer to encode the bytes into a higher-dimensional space, and then apply one or more convolutional layers to learn relevant features that are fed to a fully-connected layer for classification with a sigmoid function. While they share common design concepts, they differ in their overarching architecture and, as we will see, this difference is the key to their respective robustness to the various adversarial evasion attacks described in this paper. In addition to these two deep learning models, we also consider a traditional ML model using gradient boosting decision trees on hand-engineered features, which we use as a baseline for purposes of comparison and to evaluate attack transferability from byte-based models to models with semantically-rich features. In

general, both neural networks and standard ML algorithms, can not really work in an end-to-end way, since bytes are categorical data that do not possess a defined metric. In practice, the pipeline for predicting the maliciousness from an input program is the same, as shown in Figure 2.

MalConv. This model, proposed by Raff et al. [29], is a convolutional neural network that combines an 8-dimensional, learnable embedding layer with a 1-dimensional gated convolution. The embedding layer acts as a non-differentiable feature mapping, as it maps each input byte (treated as a categorical value) onto a specific point in the embedded space. The goal of this step is to learn to represent bytes that exhibit a semantically-similar behavior as closer points in this space, thus obtaining a meaningful distance measure between bytes. The convolutional layer iterates over non-overlapping windows of 500 bytes each, with a total of 128 convolutional filters. A global max pooling is applied to the gated outputs of the convolutional layer to select the 128 largest-activation features, without considering the structure or locality of those features within the binary. The corresponding values are then used as input to a fully-connected layer for classification. While the original MalConv model considers a maximum input file size of 2MB, the model used in our experiments is that provided by Anderson et al. [3], trained on the EMBER dataset with a maximum input file size of 1MB. Files exceeding the maximum allowable size are truncated, while shorter files are padded using a special padding token separate from the standard bytes in the file (i.e., resulting in 257 unique tokens).

DNN with Linear (DNN-Lin) and ReLU (DNN-ReLU) activations. Jeffrey Johns⁶ and Coull et al. [9] proposed a deep convolutional neural network that combines a 10-dimensional, learnable embedding layer with a series of five interleaved convolutional and max-pooling layers arranged hierarchically so that the original input size is reduced by one quarter (1/4) after each layer. The outputs of the final convolutional layer are globally pooled to create a fixed-length feature vector that is then provided as input to a fully-connected layer for the final classification. Since the convolutional layers are hierarchically arranged, locality information among the learned features is preserved and compressed as it flows upwards towards the final classification layer. The maximum length of this model is 100KB to account for the deep architecture and, as done by MalConv, files exceeding this length are truncated, while shorter files are padded with a special padding token. Several variations of this architecture are evaluated in this paper, including examining performance with both linear and Rectified Linear Unit (ReLU) activations for the convolutional layers, as well as performance when trained using the EMBER dataset and a proprietary dataset containing more than 10x the number of training samples. An analysis of the model by Coull et al. [9] demonstrates how the network attributes importance to meaningful features inside the binary, such as the name of sections, the presence of the checksum, and other structures.

Gradient Boosting Decision Tree (GBDT). A gradient-boosted tree ensemble model trained provided as part of the EMBER open-source dataset by Anderson et al [3]. The model uses a set of 2,381 hand-engineered features derived from static analysis of the binary using the LIEF PE parsing library, including imports, byte-entropy histograms, header properties, and sections, which generally represent the current state of the art in traditional ML-based malware detection. Given its use of a diverse set of semantically-meaningful static features, it provides an excellent baseline to compare the two above byte-based models against, and help demonstrate the gap between the features learned by byte-based neural networks and those created by subject-matter experts.

The MalConv architecture has been extensively studied by previous work and a wide variety of adversarial attacks have shown great success against the model [11, 12, 20, 23, 31, 33]. As pointed out by Suciu et al. [33], the lack of robustness in this model may be strongly tied to its weak notions of spatial locality among the learned features – meaning that the location of the injected adversarial noise does not matter as long as the activation on the noise bytes

⁶<https://www.camlis.org/2017/jeffreyjohns>

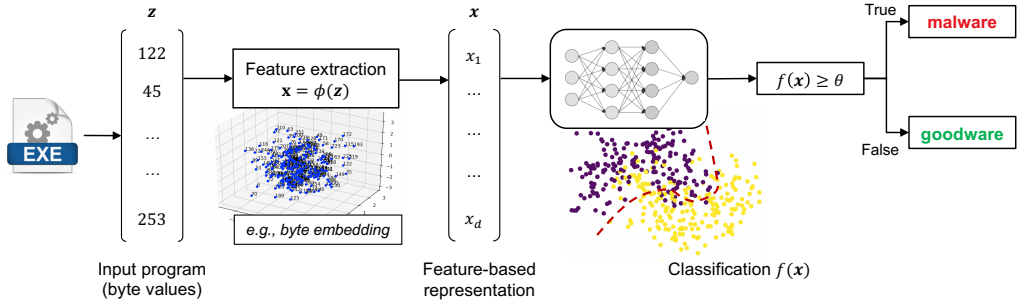


Fig. 2. Conceptual representation of malware detectors based on machine learning.

overwhelms those from the actual binary. By contrast, the deep convolutional net proposed by Johns and Coull et al. [9] enforces spatial locality among features, which means both the location and magnitude of adversarial noise play a role in the success of evasion attacks. Even the GBDT model has been shown to be vulnerable to evasion attacks [7, 12, 31], albeit with more advanced and computationally-intensive attacks. While each of these models has been previously evaluated in an ad-hoc manner, we are the first to treat attacks on machine-learning models in a holistic manner using a single unifying framework, and in doing so we uncover two new attack methods that apply to both MalConv and the Coull et al.’s model despite the unique architectural differences between the two.

3 ADVERSARIAL EXAMPLES: PRACTICAL ATTACKS ON WINDOWS MALWARE DETECTORS

In this section, we introduce RAMEN, our framework for the optimization of adversarial malware with practical manipulations. We first formalize the problem of optimizing adversarial examples with application-specific manipulation constraints under a unifying attack framework (Section 3.1), inspired from previous work in [5, 16, 28], and discuss how to implement gradient-based (white-box) and gradient-free (black-box) attacks within this framework (Sections 3.2-3.3). We then discuss the novel, functionality-preserving manipulation strategies identified in this work to manipulate Windows malware files (i.e., *Full DOS*, *Extend* and *Shift*), along with the other previously-proposed practical manipulations that are encompassed by our framework (Section 3.4). We conclude by discussing how to implement practical white-box (Section 3.5) and black-box (Section 3.6) attacks on Windows malware detectors based on the aforementioned manipulation strategies, detailing the implementation of our novel attacks and how to recast previously-proposed ones into our framework.

3.1 Attack Framework

We consider here machine-learning models that take a program as input and aim to classify it either as legitimate or malicious, as shown in Figure 2. Since code is written as arbitrarily-long strings of bytes, we define the set of all possible functioning programs in the input space as $\mathcal{Z} \subset \{0, \dots, 255\}^*$. Most classifiers process data through a *feature extraction* step, e.g., either embedding the input bytes on a representation space which reflects their semantic similarity, or using handcrafted feature values like byte histograms, n-grams, and other meaningful statistics extracted from the input file structure and content. We encode this step as $\phi : \mathcal{Z} \rightarrow \mathcal{X}$, being $\mathcal{X} \subseteq \mathbb{R}^d$ a d -dimensional vector space (i.e., the feature space). The prediction function is denoted with $f : \mathcal{X} \rightarrow \mathbb{R}$. We assume here that, without loss of generality, this function outputs a continuous value representing the probability of the input sample belonging to the malware

class. Then, we compute $f(\mathbf{x}) \geq \theta$, where θ is the detection threshold, to obtain a final decision in the label space $\mathcal{Y} = \{-1, +1\}$, being -1 and $+1$ the class labels for legitimate and malicious samples, respectively.

Under this setting, the attacker aims to craft adversarial malware by perturbing each malicious input program to achieve evasion with high confidence. To this end, the attacker is required to apply *practical manipulations*, i.e., transformations that alter the representation of the input program without disrupting its original behavior, by exploiting the redundancies offered by the executable file format. We encode these functionality-preserving manipulations as a function $h : \mathcal{Z} \times \mathcal{T} \rightarrow \mathcal{Z}$, whose output is a functioning program with the same behavior as the input, but with a different representation. The function h takes as input a program $z \in \mathcal{Z}$ and a vector $\mathbf{t} \in \mathcal{T}$, representing the parameters of the applied transformation. For example, the injection of K padding bytes at the end of the file, which is a functionality-preserving manipulation [12, 20], can be encoded in our framework as $\mathbf{t} = (t_1, \dots, t_K)$, being t_i the value of the i^{th} padding byte. Accordingly, the function $h(z, \mathbf{t})$ will return a manipulated program z' consisting of the input program z with the padding bytes \mathbf{t} appended at the end.

We are now in the position to present *RAMEN*, a general framework that reduces the problem of computing adversarial malware examples to optimization problems of the form:

$$\underset{\mathbf{t} \in \mathcal{T}}{\text{minimize}} \quad F(\mathbf{t}) = L(f(\phi(h(z, \mathbf{t}))), y). \quad (1)$$

where $L : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a *loss function* that measures how likely an input sample is classified as malware, by comparing the output of the prediction $f(\phi(z))$ on a malicious input sample z against the class label $y = -1$ of benign samples. By minimizing this loss function, the attacker aims to reduce the probability of the modified program being recognized as malware, i.e., increases the probability of evasion, while retaining malicious functionality. We discuss the two main strategies for solving this optimization problem in the following, which include gradient-based and gradient-free attacks.

3.2 Gradient-based (White-box) Attacks

The aforementioned optimization problem can be solved, at least in principle, using gradient-based approaches, similarly to attacks staged against machine-learning algorithms for image classification [4, 35]. This implicitly assumes that the attacker has white-box access to the target model to compute the gradient of the loss function L (as it requires knowledge of the model's internal parameters). However, when optimizing adversarial Windows malware, the loss gradient can not be typically used to directly update the parameter vector \mathbf{t} (e.g., the padding bytes), due to the non-differentiability of the inner feature-mapping function ϕ . For example, let us consider the case in which the input program bytes are embedded in a vector space, as done by MalConv, and our manipulation strategy $h(z, \mathbf{t})$ amounts to injecting padding bytes. In this scenario, it is not possible to compute the loss gradient with respect to the bytes in \mathbf{t} , as the embedding step performed by MalConv is not differentiable (in particular, bytes are treated as categorical variables, like words, rather than as numerical values). To overcome this issue, most of the gradient-based attacks on Windows malware proposed so far have considered optimizing the attack by performing gradient descent in the feature space, while iteratively trying to reconstruct the corresponding adversarial malware example in the input space, using different strategies. In the following, we formalize this process in the context of our framework, according to the steps detailed in Algorithm 1 and Figure 3.

Malware Manipulation. The first step consists of applying an initial perturbation $h(\cdot, \mathbf{t})$ to the input malware program z , obtaining a modified program $z' = h(z, \mathbf{t})$ (line 3). While we discuss in detail the practical manipulations that are

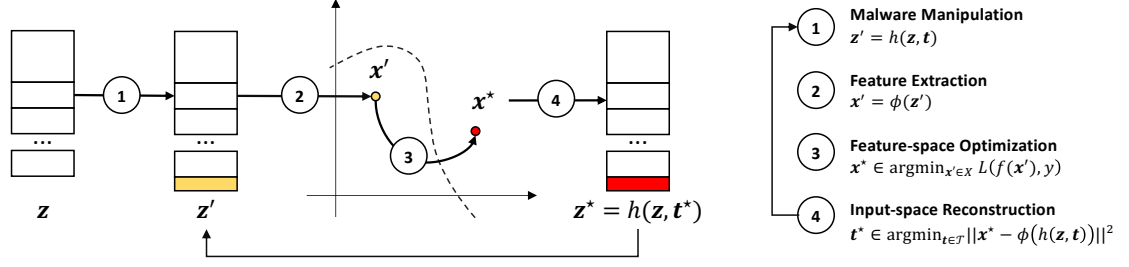


Fig. 3. The steps performed by RAMEN to optimize gradient-based (white-box) attacks, as also detailed in Algorithm 1.

Algorithm 1: Gradient-based (White-box) Attacks for Optimizing Adversarial Malware EXEmples in RAMEN.

Data: z , the initial malware sample; N , total number of iterations; y , the target class label; f , the target model.
Result: z^* , the adversarial EXEmple.

```

1  $t^{(0)} \in \mathcal{T}$ 
2 for  $i$  in  $[0, N - 1]$ 
3    $z' \leftarrow h(z, t^{(i)})$  # malware manipulation
4    $x' \leftarrow \phi(z')$  # feature extraction
5    $x^* \leftarrow \arg \min_{x' \in \mathcal{X}} L(f(x'), y)$  # feature-space optimization
6    $t^{(i+1)} \leftarrow \arg \min_{t \in \mathcal{T}} \|x^* - \phi(h(z, t^{(i)}))\|^2$  # input-space reconstruction
7  $t^* \leftarrow t^{(N)}$ 
8  $z^* \leftarrow h(z, t^*)$ 
9 return  $z^*$ 

```

encompassed by our framework in Section 3.4, the reader may consider here the injection of K randomly-chosen padding bytes at the end of the file (i.e., the yellow region in Figure 3) as an exemplary case.

Feature Extraction. The second step (line 4) amounts to encoding the manipulated program z' in the feature space as $x' = \phi(z')$ (e.g., through byte embedding in MalConv).

Feature-space Optimization. The third step (line 5) consists of modifying the feature-based representation x' of the manipulated input sample z' to minimize the loss function L , using gradient-based updates. It can be formalized as:

$$x^* \in \arg \min_{x' \in \mathcal{X}} L(f(x'), y). \quad (2)$$

The solution is obtained by iteratively updating the feature-based representation as $x' \leftarrow x' - \eta \nabla_{x'} L(f(x'), y)$, being η the gradient step size, until convergence or a maximum number of iterations is reached. As we will discuss in the remainder of this work, some attacks just consider a single gradient update in this step. It is worth remarking here that the feature-based representation x^* obtained after this step may not necessarily correspond to any input sample $z \in \mathcal{Z}$. For this reason, it is necessary to use proper reconstruction strategies to ensure not only that the reconstructed sample is a valid program, but that it also preserves the intended functionality of the initial malware program z .

Input-space Reconstruction. The final step (line 6) optimizes the parameters t^* to generate a functionality-preserving malware sample z^* whose feature-based representation is as close as possible to the desired x^* . This can be expressed

as a reconstruction problem in the form of a minimization:

$$\mathbf{t}^* \in \arg \min_{\mathbf{t} \in \mathcal{T}} \|\mathbf{x}^* - \phi(h(\mathbf{z}', \mathbf{t}))\|^2. \quad (3)$$

For attacks that only consider a single gradient update in the previous step, the reconstruction step is similarly performed by finding the transformation \mathbf{t} that is best aligned with the loss gradient in feature space, i.e., $\mathbf{t}^* \in \arg \max_{\mathbf{t} \in \mathcal{T}} \nabla_x L(f(\mathbf{x}'), y)^\top (\mathbf{x}' - \phi(h(\mathbf{z}', \mathbf{t})))$ [31]. Note that these gradient-based attacks are only convenient if the reconstruction step can be computed efficiently; otherwise, gradient-free attacks may be preferable. For instance, let us assume that the input bytes are embedded as points in a bi-dimensional vector space; e.g., bytes 0 and 1 are encoded as $0 \mapsto (1, 2)$ and $1 \mapsto (1, 5)$, respectively. If we consider the injection of padding bytes (i.e., the yellow region in Figure 3), the feature-based representations of such bytes will be shifted by the attack along the gradient direction in the embedding space, and may result in feature values that do not correspond to any byte; e.g., the initial embedding of byte 0 may be modified from $(1, 2)$ to $(1, 4.5)$, which does not correspond to any valid byte. The reconstruction step in this case is quite efficient, and it simply corresponds to remapping the modified feature values $(1, 4.5)$ to the closest feature-based representation corresponding to a valid input byte, i.e., $(1, 5)$, which corresponds to byte 1. To summarize, the reconstruction step modifies here the padding bytes (i.e., the red region in Figure 3), providing a functional malware sample. However, as this process may not exactly match the optimal feature-based representation, the resulting malware sample may not evade detection or anyway exhibit a lower misclassification confidence. The feature-space optimization and the reconstruction step can be thus iteratively repeated to improve the attack effectiveness, gradually refining the input malware manipulations.

3.3 Gradient-free (Black-box) Attacks

Another suitable approach to solving Problem (1) is to use a gradient-free (black-box) optimizer, which only requires querying the target model and observing its outputs, without accessing its internal parameters or knowing how the model works precisely. This approach is also useful when dealing either with non-differentiable models (like decision trees and random forests), for which no gradient information is available, or with feature representations that make the input-space reconstruction step non-trivial or too computationally demanding. This typically happens with hand-crafted features, as those used by the GBDT classifier discussed in Section 2.2, n -grams, and other statistical-based features for which there is no clear and direct relationship with the input transformations.

Within the gradient-free scenario, the attacker must define the *loss function* to be minimized, the *malware manipulations* they will use, and the black-box optimizer that will combine them together. These steps are similar to the white-box case in Section 3.2, except for the *feature-space optimization* and *input-space reconstruction* steps, which are not required here as the attack is directly optimized in the input space. Algorithm 2 details a high-level implementation of gradient-free attacks. In each iteration, the attack solves the minimization problem with the chosen optimizer, perturbing the malware with the given practical manipulations (line 3). In the following, we discuss two different solution strategies, depending on the resources available to the attacker, referred to as *transfer* and *query* attacks.

Gradient-free (Black-box) Transfer Attacks. Within this setting, the attacker is assumed to craft the adversarial examples against a *surrogate model*, and then evaluate whether they successfully *transfer* to a different *target model* [4, 13, 26, 27, 34]. The surrogate model here is meant to provide a good approximation of the target model, e.g., obtained by training another model on the same classification task. The optimization of adversarial EXEmples can be done by following the steps described in Section 3.2, since the surrogate might be differentiable. The details on how to obtain such a surrogate

Algorithm 2: Gradient-free (Black-box) Attacks for Optimizing Adversarial Malware EXEmples in RAMEN.

Data: z , the initial malware sample; N , total number of iterations; y , the target class label; f , the target model function

Result: z^* , the adversarial EXEmple.

```

1  $t^{(0)} \in \mathcal{T}$ 
2 for  $i$  in  $[0, N - 1]$ 
3    $t^{(i+1)} = \arg \min_{t \in \mathcal{T}} L(f(\phi(h(z, t^{(i)}))), y)$ 
4  $t^* = t^{(N)}$ 
5  $z^* = h(z, t^*)$ 
6 return  $z^*$ 

```

model are beyond the scope of this paper, but it suffices to say that there are several methods for accomplishing this through model stealing, using open-source models, or simply training a new model on an open-source dataset (e.g., EMBER). Here we focus on the latter, by optimizing attacks on the networks considered in this work, and transferring them against all the others. For instance, let us assume that the attacker wants to target a remote service for malware detection (possibly implemented with a DNN-Lin network), and they have access only to the pretrained MalConv classifier. Then, Algorithm 6 can be used to attack MalConv, while the corresponding adversarial malware samples can be submitted to the target model. Clearly, the success of black-box transfer attacks is inherently connected to the similarity of the decision functions that are learnt by the surrogate and the target models, along with the amount of changes that may be applied to the input sample. We refer the reader to [13] for a more comprehensive analysis of the main factors affecting attack transferability between different models.

Gradient-free (Black-box) Query Attacks. Within this scenario, the attacker optimizes directly the malware manipulations by querying the target model and observing its outputs. To this end, many gradient-free (black-box) optimizers can be exploited, including genetic algorithms [7, 12], natural evolution strategies [17, 37], and zeroth-order optimizers [8]. Genetic algorithms optimize the objective by maintaining a population of N perturbed samples in each iteration, from which only the best individuals are selected and used to generate the next population. Natural evolution strategies are similar in principle, as they draw the N perturbed samples in each iteration from an underlying Gaussian distribution, and then optimize the parameters of that distribution to minimize the expected loss on the N samples. Zeroth-order optimization aims instead to estimate the gradient of the loss function via finite differences, by querying the target model several times. All these approaches are iterative, and enable defining a maximum number of queries that can be sent to the target (i.e., a query budget). As detailed in Section 3.6, we will consider in this work the approach proposed in [12], which optimizes the padding bytes using a genetic algorithm.

3.4 Practical Manipulations

The optimization algorithms detailed in the previous section can be used to find the best manipulation parameters t , which then have to be applied via $h(\cdot, t)$ to craft the actual adversarial malware samples. In this section, we enumerate and categorize the *practical manipulations* that can be applied to the input malware and discuss how they can be represented in our framework via the manipulation function $h(\cdot, t)$ and its parameters t .

These manipulations are functions that change the representation of a program, by exploiting redundancies and technicalities of the file format, while leaving their functionality intact. In particular, the attacker aims either to find suitable locations where they can freely alter bytes without breaking the structure, or to create space where they can

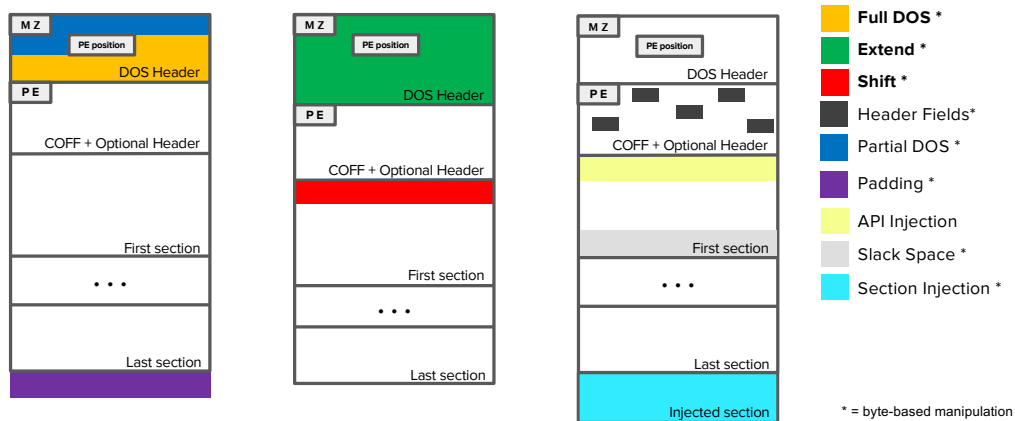


Fig. 4. Graphical representation of the locations perturbed by different attack strategies with the injection of adversarial payloads (shown in colors). The three manipulation strategies proposed in this work are highlighted in bold.

inject the adversarial payload. A simplified graphical representation of all these strategies is given in Figure 4. The colored areas highlight the locations where the adversarial payload can be injected, and the length of the boxes indicates how the content is shifted before injection.

3.4.1 Novel manipulations. We discuss here the three novel practical manipulations of Windows programs introduced in this work: *Full DOS*, *Extend*, and *Shift*. These manipulations are all *byte-based*, i.e., they can manipulate byte values independently. To this end, they exploit ambiguities of the PE format, abusing in-file offsets and adding content to specific locations. They are described in RAMEN by defining $h(\cdot, t)$ as the function manipulating offsets to create space where content can be injected, and t the byte values injected in the predefined locations given by $h(\cdot, t)$.

Full DOS. This manipulation edits all the bytes contained in the DOS header, which is only kept for compatibility with older operating systems as described in Section 2.1. Since the only two important fields in the DOS header are the magic number MZ and the 4 byte-long integer at offset $0x3c$, all the other bytes can be used by the attacker to inject the adversarial payload. Demetrio et al. [11] originally proposed a mutation applied between the magic number and the real header offset, even though, actually, the *whole* DOS header can be manipulated without corrupting the program. In particular, when the program is launched, the control is passed to the loader, which starts parsing the executable. After checking the magic number, it reads the PE offset and it jumps to the metadata of the PE header, skipping the parsing of the DOS header and stub, thus preserving the functionality of the input program. Algorithm 3 shows how to perform the editing while keeping the constraints mentioned above. After that, the algorithm proceeds by rewriting the initial bytes of the DOS header (line 3) and the bytes found after the pointer to the PE header and the PE header itself (line 4).

Extend. This manipulation aims to create new space inside the executable by enlarging the DOS header, where the adversary can inject the adversarial noise without breaking the structure of the executable. To this end, the attacker increases the offset to the PE header, forcing the loader to look up for it further ahead inside the binary, and then they can extend and manipulate this new enlarged DOS header at their will. Also, the loader will load into memory each header of the program, starting from the DOS header until the end of the Optional header, by looking at a field called *size of headers*. The latter specifies how many bytes are reserved to the metadata, and it must match their size, rounded up

Algorithm 3: Implementation of $h(z, t)$ for the *Full DOS* practical manipulation

Input : z , the original malware sample; t , the vector of parameters
Output: z' , the perturbed malware

```

1  $v = z.pe\_offset$ 
2  $z' = z$ 
3  $z'_{2,\dots,59} = t_{0,\dots,57}$ 
4  $z'_{64,\dots,v-1} = t_{58,\dots,|t|-1}$ 
5 return  $z'$ 

```

Algorithm 4: Implementation of $h(z, t)$ for the *Extend* practical manipulation

Input : z , the original malware sample; t , the vector of parameters
Output: z' , the perturbed malware

```

1  $inj = |t| - 58$ 
2  $v = z.pe\_offset$ 
3  $z' = z_{0,\dots,v-1} + 0x00 * inj + z_{v,\dots,|z|-1}$ 
4  $z'.pe\_offset = v + inj$ 
5  $z'.sizeof\_headers = z.sizeof\_headers + inj$ 
6  $S = z'.get\_sections()$ 
7 for  $s$  in  $S$  do
8   |  $s.physical\_offset = s.physical\_offset + inj$ 
9 end
10  $z'_{2,\dots,59} = t_{0,\dots,57}$ 
11  $z'_{64,\dots,v-1} = t_{58,\dots,|t|-1}$ 
12 return  $z'$ 

```

to the nearest multiple of the file alignment. Hence, to keep intact the structure of the input sample, the attacker must: (i) choose an amount that will not collide with the file alignment, (ii) enlarge the field containing the size of the headers, and (iii) increment the offset of each section entry. After this editing, the loader will again skip the adversarial content in search for the PE header, without altering its execution flow. Algorithm 4 shows the aforementioned approach. It first computes the amount of bytes to be added from the vector of parameters (line 1), and then it retrieves the PE header offset (line 2) to create the space that will host the adversarial payload (line 3). With a little abuse of notation, we write $0x00 * inj$, implying that the single byte $0x00$ is concatenated with itself inj times, and we use the symbol $+$ to denote the concatenation of byte strings and byte values. Then, the algorithm fixes all the constraints imposed by the format: (i) it increases the offset to the PE header (line 4), (ii) it increases the *size of headers* field (line 5), and (iii) it increases all the section entries (line 8). Lastly, the algorithm perturbs all the allowed bytes inside the enlarged DOS header (line 10 and 11), similarly to Algorithm 3.

Shift. This manipulation aims to create new space inside the executable by shifting the content of the first section and injecting there the adversarial noise. We recall from Section 2.1 that each section entry specifies an offset inside the binary where the loader can find the content of that section. Each offset is multiple of the file alignment, specified in the Optional Header. The job of the loader is to read these entries, and load into memory the content of each section, by looking inside the binary at the offset specified by the entry. The attacker can not interfere with this behaviour, but they can increase the offset of each section entry to carve space before the beginning of each section, and fill the

Algorithm 5: Implementation of $h(z, t)$ for the *Shift* practical manipulation

```

Input :  $z$ , the original malware sample;  $t$ , the vector of parameters
Output:  $z'$ , the perturbed malware
1  $inj = |t|$ 
2  $v = z.get\_first\_section\_offset()$ 
3  $z' = z_{0,\dots,v-1} + 0x00 * inj + z_{v,\dots,|z|-1}$ 
4  $S = z'.get\_sections()$ 
5 for  $s$  in  $S$  do
6    $s.physical\_offset + s.physical\_offset + inj$ 
7 end
8  $z'_{v,\dots,v+(inj-1)} = t$ 
9 return  $z'$ 

```

new crafted holes with chunks of bytes. As discussed before, the injected content must be, in size, multiple of the file alignment, to keep the structure of the program intact. In this way, the loader will look for the content of each section, ignoring the adversarial perturbation introduced between them. For the scope of this paper, we only inject content before the first section, but in principle this technique could also be used for shuffling the order of the sections inside the binary, or also inject different-in-length chunks of bytes between sections. Algorithm 5 shows this approach. It retrieves the position of the first section inside the binary (line 2), and it inserts the desired number of bytes in that position (line 3). Then, the algorithm fixes each offset of the section entries of the program (line 8), and it copies the adversarial payload in the newly-created space (line 8).

3.4.2 Previously-proposed Practical Manipulations. We provide here an overview of the previously-proposed practical manipulations applicable to PE files which are encompassed by our framework [2, 7, 11, 12, 12, 20, 23, 31–33].

Byte-based Manipulations. We start by discussing *byte-based manipulations* that can alter byte values independently. *Header Fields* [2]. This manipulation technique modifies specific fields contained inside the PE and the Optional Header, e.g., changing the name of each section by editing the corresponding section entries. Each byte in these fields can be independently and arbitrarily changed.

Partial DOS [11]. This manipulation alters the first 58 bytes of the unused DOS header, starting from the byte after the magic number MZ to the offset to PE header.

Slack Space [23, 33]. This manipulation fills the space between sections. When the program is compiled, to keep the beginning of each section aligned to the value specified in the header, the compiler appends a trail of zero bytes to each section to fill the gap. This space can be used to inject the adversarial payload.

Padding [20, 23]. This manipulation appends padding bytes at the end of the input file.

Section Injection [2, 12]. This manipulation creates a new section that is injected inside the target executable, along with a new section entry inside the Section Table. Accordingly, this manipulation does not only change the byte distribution exhibited by the input program, but also its structure.

Other Manipulations. We finally discuss other practical manipulations which cannot manipulate byte values independently (i.e., they are not byte-based), as this would corrupt the input file functionality.

API Injection [2]. This manipulation aims to add entries inside the Import Table of an executable, causing the OS to include more APIs during the loading process. While the attacker can not remove APIs, as this would break the

Attack	Loss Function L	Practical Manipulations $h(\cdot, \mathbf{t})$	Feature-space Optimization	Input-space Reconstruction
Full DOS (this work)	malware score	manipulate all DOS header	single gradient step	closest positive (iterative)
Extend (this work)	malware score	extend DOS header	single gradient step	closest positive (iterative)
Shift (this work)	malware score	shift section content	single gradient step	closest positive (iterative)
Padding [20]	malware score	padding	single gradient step	closest positive (iterative)
Partial DOS [11]	malware score	partial DOS header	single gradient step	closest positive (iterative)
FGSM [23, 33]	malware score	padding + slack space	FGSM	closest (non-iterative)
Binary Diversification [31]	CW loss	equivalent instructions	single gradient step	gradient-aligned transformation (iterative)

Table 1. Implementing novel and state-of-the-art gradient-based (white-box) attacks within RAMEN, according to the steps detailed in Algorithm 1 and Figure 3.

functionality of the program, they can inject API imports that will not be used by the program. In this case, the injection is made up of a complete entry that can not contain arbitrary values, but rather it must comply with a specific format.

Binary Rewriting [31, 36]. These manipulations allow the attacker to alter the code of the program by different means, spanning from substituting a set of instructions with others that are semantically equivalent (e.g. replacing additions with subtraction, and changing signs of the values), to encode all the program inside another one (packing). It is thus clear that, also in this case, the input bytes corresponding to the location affected by this manipulation can not be changed independently from each other. We refer the reader to [12] for more details on how binary-rewriting techniques can be used to craft adversarial malware.

3.5 Implementing Practical Gradient-based (White-box) Attacks within RAMEN

In this section, we discuss how to implement gradient-based (white-box) attacks within RAMEN, by connecting the implementation of Algorithm 1 (and Figure 3) with the corresponding practical manipulations (described in Section 3.4) used by each attack. To this end, for each attack, we have to define the following components: (i) the loss function, (ii) the considered practical manipulations, (iii) how the gradient is exploited to perform the feature-space optimization, and (iv) how the adversarial malware is eventually reconstructed in the input space. Table 1 summarizes such specifications both for our novel attacks and for previously-proposed ones, as detailed more specifically in the following.

3.5.1 Byte-based Attack Algorithm. For attacks using byte-based manipulations against DNNs trained on byte embeddings (such as MalConv and the DNNs discussed in Section 2.2), Algorithm 1 can be concretely implemented as detailed in Algorithm 6. The main idea is to iteratively update the byte values \mathbf{t} that can be changed by the given byte-based manipulation (e.g., the padding bytes) according to the gradient computed in the embedding space.

The attack (Algorithm 6) works as follows. As in [20], our attack starts by creating a temporary matrix containing all the embedding values, using $\hat{\phi}$ that is the function that encodes a single byte inside the feature space (line 1). It initializes the initial vector of parameters \mathbf{t} with random values (line 2), and it starts the optimization. At the beginning of each iteration, the algorithm applies the practical manipulations optimized so far and it encodes the sample inside the feature space (line 4). Then, the optimizer computes the gradient w.r.t. the embedded bytes (line 5), which is negative here as we are considering a minimization problem. Note also that the gradient is a matrix in this case, given that each byte is represented as a vector in the embedding space. The algorithm ignores all the locations that cannot be modified by applying a binary mask \mathbf{m} to the gradient (line 5). In this way, the algorithm will discard every location that is not related to the computation of the attack, i.e., bytes which are not associated to the perturbed input bytes \mathbf{t} . After applying the mask, the algorithm computes the norm of each embedded value (line 6) and it takes the indexes of the first γ non-zero sorted entries (line 7). The γ parameter is a step-size constant that controls how many bytes

Algorithm 6: Implementation of Algorithm 1 for byte-based attacks on DNNs using byte embedding

Input : z the original malware sample; γ the number of bytes to optimise; N the total number of iterations
Output : z^* , the adversarial EXEmple

- 1 $E_i = \hat{\phi}(i), \forall i \in [0, 256]$
- 2 $t^{(0)} \in \mathcal{T}$
- 3 **for** i **in** $[0, N - 1]$
- 4 $X \leftarrow \phi(h(z, t^{(i)}))$ // feat. space
- 5 $G \leftarrow -\nabla_X f(X) \odot m$
- 6 $g \leftarrow (\|G_0\|, \dots, \|G_n\|)$
- 7 **for** k **in** $\text{argsort}(g)_{0, \dots, \gamma} \wedge g_k \neq 0$
- 8 **for** j **in** $[0, \dots, 255]$
- 9 $S_{k,j} \leftarrow G_k^t \cdot (E_j - X_k)$
- 10 $\tilde{X}_{k,j} \leftarrow \|E_j - (X_k + G_k S_{k,j})\|_2$
- 11 $t_k^{(i+1)} \leftarrow \arg \min_{j: S_{k,j} > 0} \tilde{X}_{k,j}$ //input space
- 12 $t^* \leftarrow t^{(N)}$
- 13 $z^* \leftarrow h(z, t^*)$
- 14 **return** z^*

are perturbed at each round, modulating how much the space is explored during the optimization. For each perturbed byte, the algorithm computes a line passing from the current value to be replaced, and whose direction is imposed by the gradient in that point. The algorithm proceeds by projecting all the 256 embedded byte values on this direction, computing the distance point-to-line and the alignment with such direction (line 10). Hence, our optimizer always performs one *single gradient step* at the time. Finally, the algorithm iterates over all the computed distances, and it chooses the byte associated with the closest embedding value that lies on the positive direction pointed by the gradient (line 11). Since each embedding value computed in this way has a one-to-one mapping with the bytes in input space, and only values aligned with the gradient will be chosen, the reconstruction is *closest positive*, since they must also be on the positive side of the privileged direction. Our approach is *iterative*, after N iterations the hybrid optimizer returns the adversarial EXEmples z^* optimized so far.

3.5.2 Novel attacks. We introduce here the three attacks proposed in this work, i.e., *Full DOS*, *Extend*, and *Shift*. All of them rely on Algorithm 6 for manipulating each single byte that is injected and manipulated, by specifying the practical manipulation they apply and how they shape the mask m .

Full DOS. To implement this attack, we use the homonym manipulation as the h function (line 4) described in Algorithm 3, and we customize the mask m by specifying as editable (i) the first 58 bytes after the magic number MZ, and (ii) all the bytes from the pointer to the PE header, to the beginning of the PE header itself.

Extend. To implement this attack, we use the homonym manipulation as the h function (line 4) described in Algorithm 4, that enlarge the DOS header as specified by the vector of manipulation t . Hence, the mask m includes all the locations manipulated by the *Full DOS* attack, plus all the content that has been injected by the practical manipulation.

Shift. To implement this attack, we use the homonym manipulation as the h function (line 4) described in Algorithm 5, that shifts the content after the PE header as specified by the vector of manipulation t . Hence, the mask m specify as editable all the locations that are added by the practical manipulation.

3.5.3 Recasting Existing Attacks into RAMEN. After having shown how RAMEN can encode many different types of practical manipulations, we reinforce the claim of its generality by showing how to encode also other optimization strategies proposed in the state of the art.

Padding. Kolosnjaji et al. [20] optimize the generation of adversarial EXEmples by minimizing the malicious score computed by the classifier, and they apply the *Padding* practical manipulation for injecting the content inside the input samples. The optimizer they use is gradient-based, and it is similar to Algorithm 6, as they choose the closest embedding value in the positive direction of the gradient (*closest positive* in Table 1). Lastly, the reconstruction inside the input space is done by inverting the embedding look-up that is applied on bytes (*inverse look-up* in Table 1). This is doable, since the optimizer only chooses embedding values that correspond to real bytes, and the one-to-one mapping is preserved.

Partial DOS. Similarly, also Demetrio et al. [11] use the malicious as loss function to minimize. The manipulation they use is the *Partial DOS*, and again they use the same optimizer and reconstruction algorithm as Kolosnjaji et al. [20].

FGSM. Kreuk et al. [23] and Suciu et al. [33] use the malware score as loss function to minimize, and both of them rely on the *Slack Space* and *Padding* manipulations. They leverage a gradient technique inspired to the Fast Gradient Sign Method (FGSM) [14]: they first map the sample inside the feature space, and they apply the classical FGSM method, moving the feature vectors without applying constraints. The latter is formalized as $\mathbf{x}_{i+1} = \mathbf{x}_i + \epsilon \text{sign}(\nabla_{\mathbf{x}_i} f(\mathbf{x}_i))$, where ϵ is the maximum distortion that can be applied on the sample. While Suciu et al. [33] use only one single iteration for this technique, Kreuk et al. [23] performs more steps inside the feature space, repeating the process until evasion. Since the non-iterative version is just the first step of the iterative one, we will discuss only the latter, since it is more generic. The reconstruction is done at the end of the optimization, where each embedding value is must be converted to a real byte inside the input space (*closest* in Table 1). This is done for each entry of the final feature vector, by taking the byte whose embedding value is closest to the one contained inside the embedded malware.

Binary Diversification. Sharif et al. [31] use the Carlini & Wagner loss (CW loss) [6], and they apply random manipulations that change instructions inside the *.text* section with semantics-equivalent ones, or they displace the code inside another section, with the use of *jump* instructions. At each iteration of the algorithm, the latest adversarial example is used as a starting point for the new one. Once randomly perturbed, the new and the old versions are projected inside the embedding space, where the two points are used for computing a direction. If this direction is parallel to the gradient of the function in that point, the sample is kept for the next iteration, otherwise it is discarded. In this case, the strategy does not optimize the sample inside the feature space, since each sample is constructed by applying random transformation.

Except for the manipulations proposed by Sharif et al. [31], all the strategies described so far can be found in the `secml-malware` library released along with this paper.⁷

3.6 Implementing Practical Gradient-free (Black-box) Attacks within RAMEN

We discuss here how to implement gradient-free (black-box) attacks within RAMEN based on detailing the implementation of Algorithm 2 using the practical manipulations described in Section 3.4.

To this end, one has to define the following components: (i) which loss function they minimize, (ii) which practical manipulations they use, and (iii) which optimizer they apply for solving the problem. We fill Table 2 with all the

⁷https://github.com/zangobot/secml_malware

Attack	Loss Function L	Practical Manipulations $h(\cdot, \mathbf{t})$	Optimizer	Validation
Full DOS (this work)	malware score	manipulate all DOS header	genetic	none
Extend (this work)	malware score	extend DOS header	genetic	none
Shift (this work)	malware score	shift section content	genetic	none
Partial DOS (this work)	malware score	partial DOS header	genetic	none
Padding (this work)	malware score	padding	genetic	none
GAMMA padding [12]	malware score + size penalty	padding with benign sections / benign section injection	genetic	none
RL Agent [2]	malware score	padding + section / API inj. + header fields + binary rewriting	reinforcement learning	none
AIMED [7]	malware score	padding + section/API inj. + header fields + binary rewriting	genetic	sandbox
AEG [32]	malware score	padding + section inj. + header fields + binary rewriting	random manipulations	sandbox

Table 2. Implementing novel and state-of-the-art gradient-free (black-box) attacks within RAMEN, according to the steps detailed in Algorithm 2.

specifications used for implementing either our novel attacks, either the existing techniques proposed in the previous literature within the RAMEN framework.

3.6.1 Gradient-free (black-box) byte-based optimizer. Symmetrically to Section 3.5, we start by introducing the optimizer we developed for this work. Since no gradient information is available, this optimizer relies on a genetic algorithm to compute bytes that mostly decreases the chosen loss function, injected by applying practical manipulations. We extend RAMEN with the genetic black-box optimizer used by Demetrio et al. [12] for computing their attack. This algorithm creates byte sequences that are injected inside the sample through the usage of practical manipulations, and the newly generated adversarial EXEmple is sent to the detector to be scored. Also, this optimizer is bounded by the total number of queries that can be sent to the detector, and the process halts when the algorithms has reached that limit. Since this genetic black-box optimizer works with real numbers, but the values that we inject are integer values between 0 and 255, we encode our vector of parameters as $\mathbf{t} \in [0, 1]^k$, where k is the number of values that will be perturbed (e.g. the *Partial DOS* attack sets k to 58), and we multiply this value by 255 when the vector of parameter is passed to the h function. Before applying the practical manipulation h , we need to multiply by 255 and rounding to the nearest value the vector of parameter, since Algorithm 3, 4 and 5 consider each entry of vectors \mathbf{t} as bytes to be placed inside the sample. We summarize the optimizer in Algorithm 7, where we have plugged the loss to minimize as dictated by RAMEN. The algorithm is initialized by randomly generating a matrix of bytes sequences to inject $S' = (\mathbf{t}_1, \dots, \mathbf{t}_N) \in \mathcal{T}^N \subset [0, 1]^{N \times k}$, which represents the initial population of N candidate manipulation vectors (line 2). The genetic optimizer iterates over three steps that ensures the differentiation of the upcoming generation of candidates: *selection*, *cross-over*, and *mutation*. The *selection* step (line 4) applies the objective function F to evaluate the candidates in S' . It then selects the best N candidates between the current population S' and the population generated at the previous iteration S . The resulting vectors are the candidates associated with the lowest values of F , i.e. the most promising for creating adversarial EXEmples. The *crossover* function (line 5) takes the selected candidates as input and returns a novel set of N candidates by mixing the values of pairs of randomly-chosen vector candidates. The *mutation* function (line 6) changes the elements of each input vector at random, with a fixed low probability. The combination of both *cross-over* and *mutation* ensures that each population has different traits w.r.t. to the previous one, allowing the algorithm to properly chose new potential solution, i.e. byte sequences, discarding the useless ones, and exploring the space of feasible solutions. After T queries, the optimizer extracts the vector of manipulation \mathbf{t} associated with a minimal value for the objective function F (line 9), and it is used to create and return the final adversarial EXEmple \mathbf{z}^* (line 11).

Algorithm 7: Implementation of Algorithm 1 for implementing byte-based attacks against Windows malware classifiers

Input : z , the initial malware sample; N , the population size; T , the query budget.
Output : z^* , the adversarial EXE file.

```

1  $q \leftarrow 0, S \leftarrow \emptyset$ 
2  $S' \leftarrow (t_1, \dots, t_N) \in \mathcal{T}^N$ 
3 while  $q < T$  do
4    $S \leftarrow \text{selection}(S \cup S', F, t)$ 
5    $S' \leftarrow \text{crossover}(S)$ 
6    $S' \leftarrow \text{mutate}(S')$ 
7    $q \leftarrow q + N$ 
8 end
9  $t^* \leftarrow$ , best candidate from  $S$  with minimum  $F$ 
10  $z^* \leftarrow h(z, t^*)$ 
11 return  $z^*$ 

```

3.6.2 *Novel attacks.* Since the algorithm optimizes directly t , the attacker only need to chose a practical manipulation of their choice, and plug it inside the objective function F . In this work, the attacks that we implement with Algorithm 7 are *Full DOS*, *Extend*, *Shift*, *Partial DOS*, and *Padding*, and all of them use the malware score as loss function to minimize.

3.6.3 *Recasting Existing Attacks into RAMEN.* Here we show how can we re-implement existing attacks into RAMEN, by specifying the main components of our formalization for each of them.

GAMMA [12]. Formulated as a black-box query attack, it tries to find adversarial EXE files by minimizing the malware score, plus a penalty term that controls the size of the injected content: $f(\phi(h(z, t)) + \lambda C(t))$, where λ is the regularization parameter, and C is a function that counts how many bytes have been injected using t . In this scenario, each vector t specifies how much content must be harvested from goodware programs, injected using either the *Padding* or the *Section Injection* practical manipulation. Since these manipulations do not break the functionality by design, the attack does not need a sandbox that verifies the functionality of the perturbed malware. Then, the attack applies a genetic optimizer to find a minimum to the aforementioned function, by taking into account both the score and the size constraint.

RL Agent [2]. Formulated as a black-box query attack, this strategy minimizes the malware score by using a reinforcement learning algorithm. It uses a mixture of many practical manipulations (*Padding*, *Section Injection*, *API Injection*, *Header Fields*) with random content, and the agent is trained to evade a local model (i.e. a baseline version of the GBDT classifier). They also use one single *Binary rewriting* manipulation, that creates a new entry point inside the code of the program. The new code then jumps back to the original entry point, continuing the normal execution. It learns the best sequence of actions that leads to adversarial EXE files, and then it is tested against other targets (e.g. VirusTotal). They do not use any sandbox, since most of the manipulations are functionality-preserving except for the entry point modification. The latter could break the functionality, since the program could be loaded at a custom address in memory, hence the jump could fail to find the original code, causing a crash.

AIMED [7]. This strategy is very similar to the one applied by the *RL agent* [2] discussed earlier. They minimize the malware score, and they use the exact same manipulations in the same way. The difference is in the way they optimize the chain of manipulations: here, they use a genetic optimizer that looks for the best sequence of actions to apply.

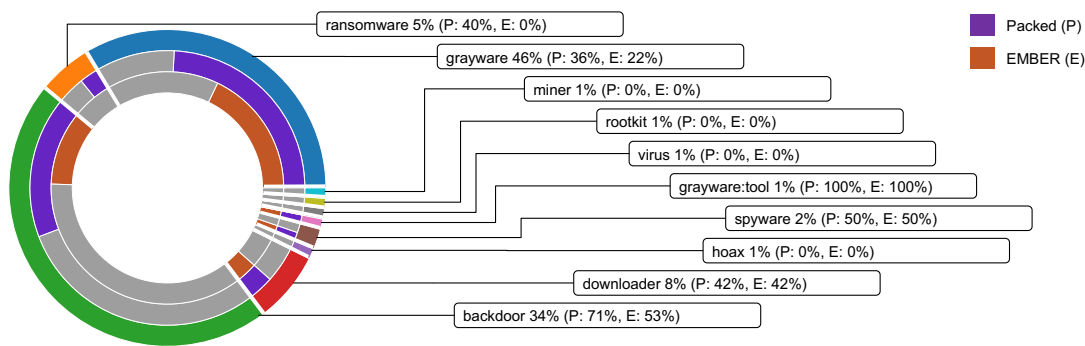


Fig. 5. The composition of the test dataset used during the attacks. The legend shows also the percentage of packed samples (**P**) in each malware category, and the percentage of malware present in the EMBER dataset (**E**).

At each step of the optimization process, they validate the newly created adversarial EXEmple by testing it inside a sandbox.

Adversarial Example Generation (AEG) [32]. This strategy minimizes the score of the target classifier, and it applies a mixture of different practical manipulations (*Padding*, *Section Injection*, *Header Fields*), and they also apply one *Binary rewriting* manipulation, that is the swapping of assembly instructions with equivalent ones (similarly to Sharif et al. [31]). Then, they apply these manipulations at random, looking for the best sequence of actions that leads to evasion. At each step, the attack checks the validity of the current sample inside a sandbox environment.

4 EXPERIMENTAL ANALYSIS

We used a Ubuntu 16.04.3 LTS server, with an Intel® Xeon® E5-2630 CPU, and 64 GB of RAM to perform our experiments. We also used a Windows 10 virtual machine during the development of the practical manipulations described in Section 3.4 to debug them during their development, and to validate that they indeed do not impact functionality. This virtual machine has not been used to test all the samples, but only for studying the manipulations and the file format with the current version of the Windows loader. To highlight the performance of our strategy, we encoded other attacks proposed in the state of the art [11, 20, 23, 33] and we ran them against the chosen targets. The network proposed by Johns⁸ and Coull et al. [9] has been trained with two different datasets. The first one is EMBER [3], which is an open-source dataset of goodware and malware hashes, including a set of pre-extracted features, while the second is a proprietary production-quality dataset used for training malware classifiers. The first dataset is smaller, counting 1.1 M samples, while the second is larger, counting 16.3M files. MalConv has been trained on EMBER [3], like the GBDT model proposed by Anderson et al. [3]. We encoded all the strategies inside a Python library we are developing and maintaining, called `secml-malware`,⁹ as an extension of `secml` [25]. Regarding the black-box setting, we omit from the analysis the FGSM proposed by Kreuk et al. [23] and Suci et al. [33], since they are similar to the Padding strategy. We also test GAMMA [12], a regularized black-box strategy whose practical manipulations consist in injecting content

⁸<https://www.camlis.org/2017/jeffreyjohns>

⁹https://github.com/zangobot/secml_malware

harvested from benign software to elude detection. For this experimental setup, we rely on the *padding* version of this attack, where the adversarial noise is appended at the end of the input malware.

Test Dataset. The malware set we used for the empirical evaluation is a richer version of the one used by Demetrio et al. [11], and we show its composition in Figure 5. Each slice of the plot resembles the percentage of one or more families of malware. The labels of the legend also depict the percentage of samples of that families that are packed, and how many of them are also contained inside the EMBER dataset. These samples were retrieved from *DasMalwerk*¹⁰ during late 2018. The total amount of samples is 104, and 37 of them are contained inside the EMBER dataset (roughly, one third of the overall dataset). This dataset has been labelled by querying VirusTotal,¹¹ and it includes 46% of *grayware*, 33% of *backdoor*, 5% of *ransomware*, and smaller fractions of *spyware*, *rootkits* and *miners*. According to the analysis, 46 of them are packed.

Malware Detection Performance. Before delving inside the performance of the different attacks against the target classifiers, we first compute the Receiver Operating Characteristic (ROC) of the four models, shown in Figure 6a. The score has been computed on the test set of EMBER v1 dataset [3]. For each classifier, we compute the detection threshold θ , and the corresponding Detection Rate (DR) at 0,1% of False Positive Rate (FPR). We report these findings in Table 6b. It is clear from Figure 6a that the GBDT model outmatches the convolutional networks, which aligns with previously reported results on this dataset [3]. This benefit might be connected to the manual feature engineering that is used by the GBDT model, instead of letting the network learn the relevant features itself. The non-linearity imposed by ReLU activation functions does not seem impact the overall score, implying that the majority of the examples in the dataset can be linearly separated. However, we do note that at lower false positive rates, the performance of the ReLU models does exceed that of the linear activation models, which may support the notion that there are some samples that are difficult to separate and where the non-linearity is useful. MalConv shows comparable though somewhat lagging results to both the DNN models trained on EMBER and the GBDT model. Clearly, the use of a larger and more diverse proprietary dataset does not necessarily improve the generalizability of the end-to-end models, which agrees with previous observations by Coull et al. [9] showing that overfitting may actually be beneficial in malware classification tasks. For each classifier, we compute the threshold such that the classifier has a 0.1% False Positive Rate (FPR).

4.1 White-box Attacks

We tested all the differentiable models with the attacks formulated in Section 3.4 (*Full DOS*, *Extend* and *Shift*), the header attack [11], the padding attack [20], an iterative implementation of the fast gradient sign method (FGSM) that address both padding and slack space [23, 33]. Since the *Full DOS* attack searches for the PE signature inside the sample, and it marks as editable all the bytes in between the two headers, the amount of bytes to perturb varies from 118 to 290 bytes (since such position it might change from file to file). We set the length of the adversarial injection for the *Extend* to 512, and this number is rounded to the nearest multiple of the file alignment specified by the sample: in our test set, it varies between 512 and 4096 bytes, resulting in a payloads whose length varies between 630 and 4386. Similarly, we chose for *Shift* attack an injection of 1024 bytes before the first section of the sample, and it must align to the nearest multiple of the specified file alignment, resulting in adversarial payloads with a length between 1024 and 4096 bytes. The *Partial DOS* attack proposed by Demetrio et al. [11] alters only the first 58 bytes contained in the DOS header. The *Padding* attack appends bytes at the end of the file, with a default payload size of 10240 KB, motivated by the results obtained by Kolosnjaji et al. [20]. The iterative implementation of FGSM has an ϵ parameter that controls the amount

¹⁰<https://www.dasmalwerk.eu/>

¹¹<https://www.virustotal.com>

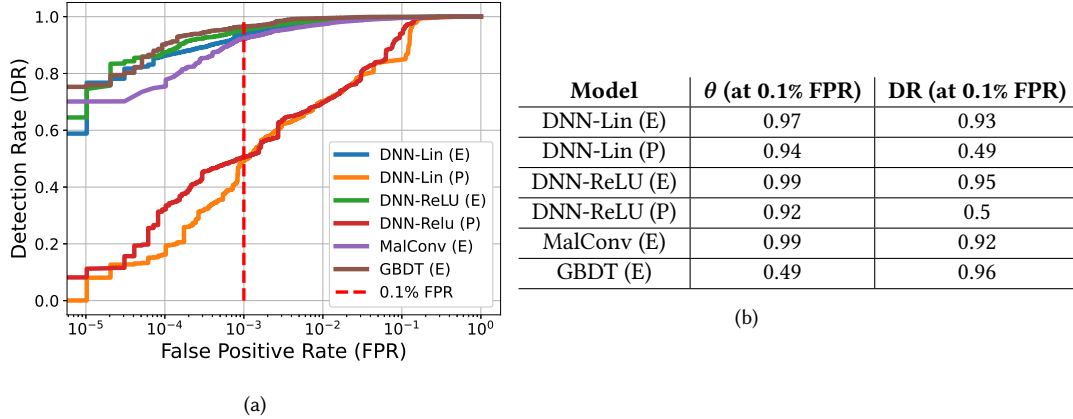


Fig. 6. On the left, the Receiver Operating Characteristic curve (ROC) of the classifiers under analysis, evaluated on the EMBER test set [3]. The letter inside the parenthesis specifies the dataset used for training the classifier: *E* means EMBER, while *P* implied the use of a larger proprietary dataset. The red dashed line highlights the performance of each classifier at 0.1% False Positive Rate (FPR). On the right, the detection thresholds of the classifiers (θ) at 1% FPR, with the corresponding Detection Rate (DR).

of injected noise, and we set this parameter to 0.1. For all the attacks, we have chosen a step size of 256 bytes optimized at each iteration.

While the *Partial DOS* technique is generally ineffective against all classifiers except for MalConv (as already pointed out by Demetrio et al. [11]), the *Full DOS* attack does substantially lower the detection rate of the networks proposed by Coull et al. [9]. This might be caused by spurious correlations learnt by the network, and altering these values cause the classifier to lose precision.

Both our novel strategies, *Extend* and *Shift* show great attack performance against all evaluated networks. Since these attacks replace a portion of the real header of the program, it might be possible that the adversarial noise interferes with the local patterns learned by the networks at training time, like the position of the meaningful metadata of the program. The *Extend* attack, for instance, covers the original position of the PE offset plus many fields of the Optional Header, like the *checksum* and the locations of directories such as the Import Table and Export table. This content is preserved, since it is shifted, but it is no longer present in the position the network believed them to be. The *Shift* attack does the same, but with the content of the first section, that is usually the one containing the code of the program. Surprisingly, the *Shift* attack against MalConv is not as effective as it was against the other networks. The reason might be once again the wrong feature importance that MalConv attributes to certain bytes. Analyzing the norm of the gradient computed on the location altered by the attack, we found that it is mostly zero, and the attack is unable to optimize the payload. If the attention is focused on the header, the rest of the file has a low impact on the final score. This can be glimpsed by looking at the *Extend* attack, which manipulates an extended portion of bytes starting from the DOS header.

The *Padding* attack proposed by Kolosnjaji et al. [20], and the *FGSM* attack proposed by Kreuk et al. [23] and Suciuc et al. [33] do not decrease much the detection rates of the networks, since most of the manipulations applied are cut off by the limited window size of the network itself. For instance, if a sample is larger than 100 KB, it can not be padded, and all the strategies that rely on padding fail. To achieve evasion, these FGSM attacks can only leverage the perturbation of the slack space, but the number of bytes that can be safely manipulated is too few to have significant impact. Also, this strategy is incapacitated by the inverse-mapping problem: they compute the adversarial examples inside the feature

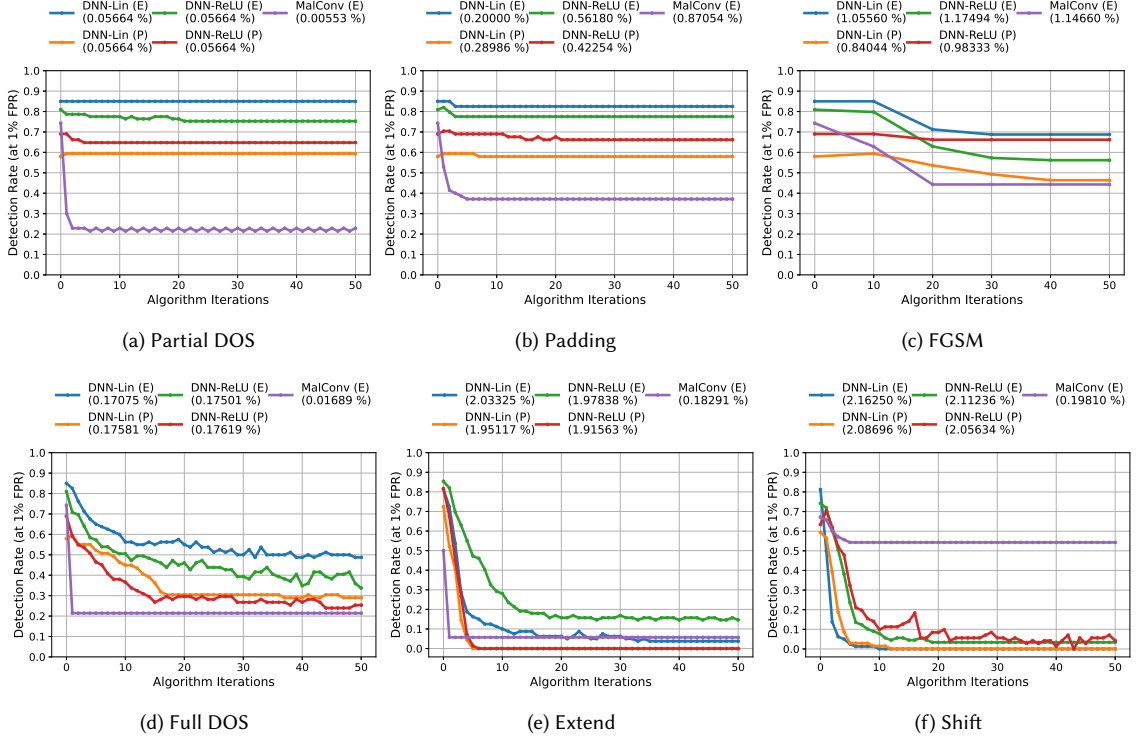


Fig. 7. The results of white-box attacks, expressed as the mean Detection Rate at 0.1% FPR, at each optimization step. Each plot sums up the degradation induced by a specific strategy against the classifiers we have considered, trained on different datasets (E for EMBER and P for proprietary). The number near the name of the classifier represents the size of the adversarial payload as a percentage of the input size.

space, and they project them back only at the end of the algorithm. This means that the attack might be successful inside the feature space, but not inside the input space, where there are a lot of constraints that are ignored by the attack itself. Against MalConv, the *Padding* attack proves to be quite effective, but it needs at most 10 KB to land successful attacks, as already highlighted by Kolosnjaji et al. [20]. The adversarial payload must include as many bytes as possible to counterbalance the high score carried by the ones contained inside the header.

Regarding the size of the adversarial payload injected by our novel attacks, we report the mean percentage size of the crafted noise w.r.t. to the input window of the target network near every label of the legend of Figure 7. This network has a window size of 100 KB, and each attack alter, on average, 2% of that quantity (approximately, 3 KB). Also, we can observe that both DNN-Lin and DNN-ReLU trained on the larger dataset are less robust w.r.t. to their counterparts trained on EMBER, and this pattern can be observed in almost every white-box attack we have showed.

4.2 Black-box Transfer Attacks

We show how the classifiers under analysis behave against black-box transfer attacks. The latter is crucial since an attacker might optimize attacks against a model they own, and then they can try to evade other systems in the wild. To this extent, we use the adversarial EXEmples crafted for the white-box attacks, we test them against all the other

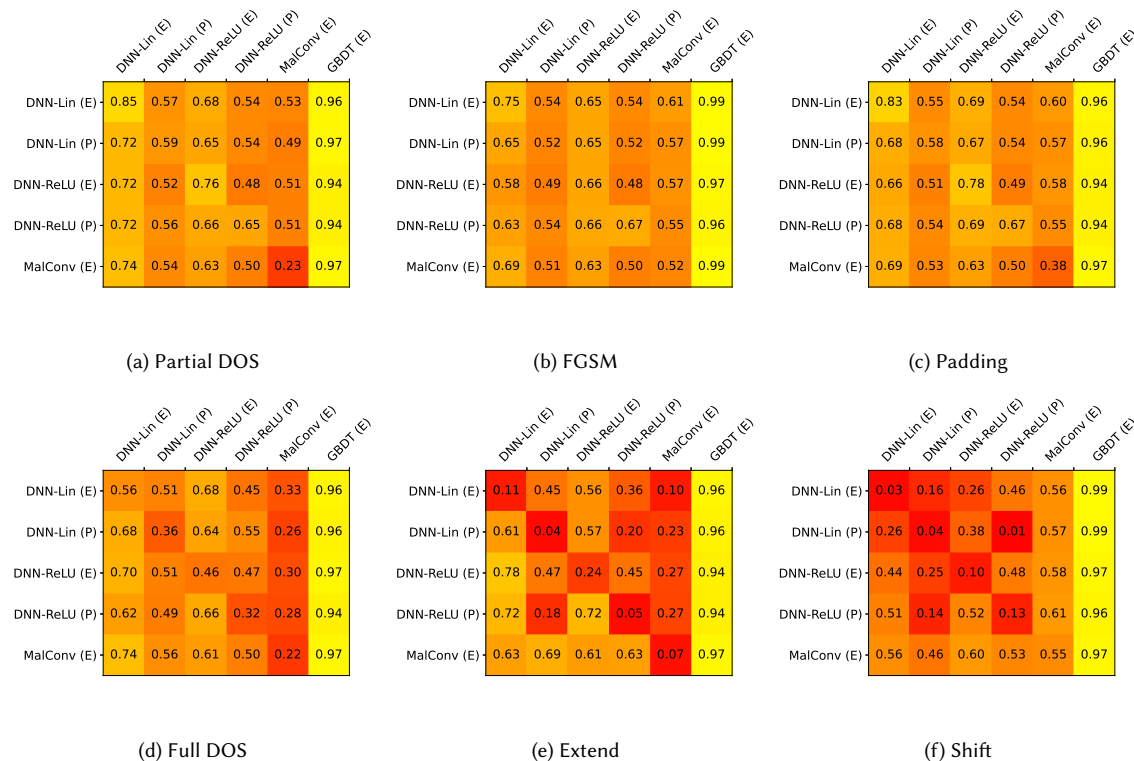


Fig. 8. Detection rate at 0.1% FPR for transfer attacks. Substitute models used to craft the attack are reported in rows, while the target models are reported in columns. Values on the main diagonal correspond to white-box attacks.

models, and we use again the threshold θ at 0.1% FPR to compute the DR. We plot the transfer results on Figure 8, where the rows of each square we report the model we used for computing the attack (i.e., the surrogate), while each columns represent the target of the attack. Since the GBDT model is not differentiable, we only use it as a target and not as a surrogate. From various plots in Figure 8, we can highlight different interesting aspects. The first one is that the new manipulations, *Extend* and *Shift* generally decrease the performance of every byte-based networks. This phenomenon probably happens because the networks learnt a specific location for particular patterns of bytes (e.g. the magic number PE of the PE header), or some section names (e.g. *.text*, *.data*, etc.), and association is disrupted by the injection of new content, leading these known values to be placed elsewhere, unable to be find again by the networks. The *Full DOS* attack has a good impact against MalConv, even in the transfer setting. Again, this is another empirical proof that MalConv focus on the header of a program, rather than its other components, and it is sufficient to optimize header attacks even against other network to produce successful adversarial EXEmple against it. The other attacks do not transfer as well as the novel strategies, and they reflect the results obtained in the ROC curve (Figure 6a). The *Partial DOS* transfer attack manipulates too few bytes compared to the other strategies. Also, this attack is not optimized directly against the target, so the effect is similar to almost-random noise. The *FGSM* performs slightly better than the *Padding* attack, since it relies also on the *Slack Space* practical manipulation, that it is less likely to have been cut by the input window of the target model. Hence, more adversarial noise is considered by the classifier, but not enough

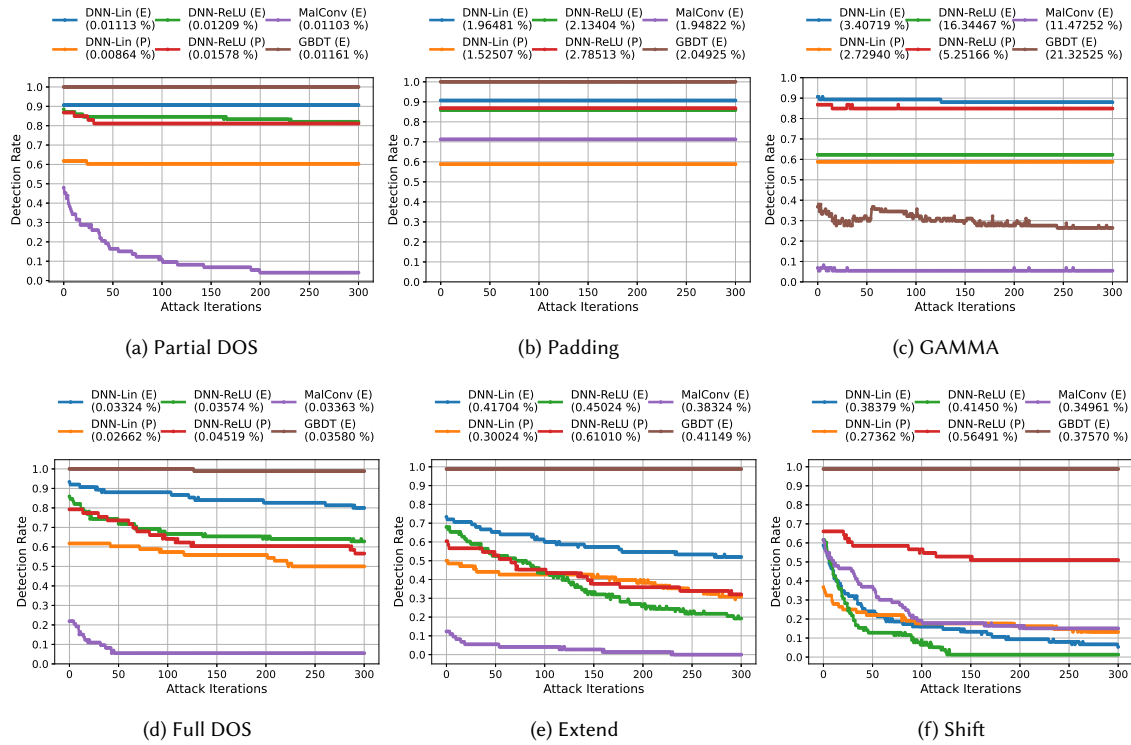


Fig. 9. The black-box attacks against all the models. We report the Detection Rate at each step of the black-box optimizer, using a DR at 0.1% FPR. The number near the name of the classifier represents the size of the adversarial payload w.r.t. the mean file size of malware test set.

for decreasing the detection rate of the targets. This highlights once again the problems of *Padding* manipulations, since they are easily stripped away. Another key observation is the similarity of transfer results for the DNN-Lin and DNN-ReLU networks. Since these networks are similar in architecture and they were trained on the same dataset, it is possible that they also learnt the same blind spots. The ones trained on the larger proprietary dataset are not more robust than the ones trained on EMBER, but rather they suffer more from transfer attacks. In the same way, the attacks optimized against them are not very effective against the EMBER versions. The explanation might consider that the decision boundary learnt over more data is much more complex, filled with more local minimum and maximum. This phenomenon has already been observed in the past, by Papernot et al. [26], where they show that more complicated models suffer from transfer attacks, caused by the high-non linearity of their decision boundary. On the other hand, the GBDT model is not affected by any adversarial transfer attack, as the byte-based features used by the decision-tree algorithm are only a small subset of all the characteristics considered by the classifier, such as the API imports, metadata and more. These attacks are not directly optimized against it, and the quantity of bytes that are altered is very little compared to the whole file size. For sure, the adversarial payload has a minimal effect on the byte-based features, but not enough to counterbalance all the other ones.

4.3 Black-box Query Attacks

We evaluate the effects of the black-box query attacks, and we plot the performances of these strategies in Figure 9. We use the mean Detection Rate as metric for this analysis, and we set the population size $N = 10$, and the number of queries to $q = 3000$. This produces 300 iterations of the genetic algorithm, since each round consumes a number of queries that matches the population size (10 in this case). Not surprisingly, the first thing we can notice is that the performances of these black-box query attacks are worse w.r.t. their white-box counterparts. This is intuitive, since the optimizer does not have any clue regarding a privileged direction to take, and it needs to explore the surrounding space to approximate such direction. Also, it seems that only the *Shift* attack meaningfully decrease the DR of the networks, while both the *Full DOS* and *Extend* attacks decreases this metric almost linearly in the number of the queries sent. This might implies that these attacks requires more queries to gain more effectiveness against the byte-based networks. It is interesting to see that, since these mutations origin from random perturbations, the DNN-Lin and DNN-ReLU trained on the proprietary dataset show more robustness w.r.t. their counterparts, opposed to the white-box results. This might be explained by robustness to random noise induced by the high volume of data used for at training time. Another interesting point, is that MalConv is evaded also by the *Shift* attack, when the same was failing in the white-box settings. As explained in Section 4.1, the norm of the gradients of the locations perturbed by that manipulations were zero, but the genetic optimizer do not rely on that. This success is probably caused by the combination of random exploration and optimization done by the optimizer, successfully avoiding regions with no gradient information. Another important key point is that the GBDT shows great stability against the novel proposed attacks. This can be caused by a combination of different failures. We have discussed one of these probable failures while analysing the transfer attacks, in Section 4.2: the byte-based features of the GBDT are a minority of all the characteristics that are considered, and these manipulations are not enough for being a relevant contribution during the classification. Also, the optimizer might not be able to craft patterns that are relevant for the GBDT, since it also has been trained on strings with meaning. On the other hand, as already proved by Demetrio et al. [12], the black-box GAMMA algorithm is effective against MalConv and the GBDT. We considered 100 *.data* sections taken from our goodware dataset, and we use $\lambda = 10^{-5}$, to match the settings of the original formulation. The attack, that leverage the *Padding* practical manipulation for injecting benign content, shows similar results to the one obtained by the authors, and it is the only functioning attack against this model. This is the effect of the injection of benign content: it is difficult for the optimizer to find patterns of bytes without guidance, while using content that already matches the target class is helpful, and the optimizer does not need to create these patterns. Maybe, by highly increasing the number of queries and the size of the noise, even the other attacks would be successful in the end, but the cost in terms of query budget could be unsustainable. On the other hand, this strategy is not effective against the DNN-Lin and DNN-ReLU networks, since they focus on a smaller input window than MalConv. Hence, the adversarial content is stripped away, losing its efficacy.

5 RELATED WORK

In this section we briefly review some concurrent work proposing a similar attack framework, provide additional details on attacks against learning-based malware detectors, and conclude by discussing relationships between the problem of adversarial malware to packing and obfuscation techniques.

5.1 Other Attack Frameworks

Recent concurrent work by Pierazzi et al. [28] propose a general formalization for defining the optimization of adversarial attacks inside the input space, spanning multiple domains like image and speech recognition and Android malware detection. They define sequences of practical manipulations that must preserve the original semantics. These manipulations need to be imperceptible to manual inspection, and they must be resilient to pre-processing techniques. The authors also explicitly define the resulting side-effects generated by applying such mutations to the original sample, as a summation of vectors. The attacker optimizes the sequences of practical manipulations that satisfy all the constraints mentioned above by exploring the space of mutations imposed by such practical manipulations. Our formalization shares the use of practical manipulations applied inside the input space, and we also minimize the applications of practical manipulations to compute adversarial examples. Furthermore, our practical manipulations are functionality-preserving by design, including all the constraints proposed by Pierazzi et al, and removing the need of the side-effect vector. We do not enforce the robustness to the pre-processing step, as the defender needs to know such manipulations in advance to decide what it must be sanitized. The latter is not trivial, as the novel attacks we propose act as a zero-day, and they can be patched properly only after their discovery. We generalize the objective function to optimize by including a loss function, and the attacker can choose how to implement it by using a function of their choice, while also adding constraints expressed as regularization parameters. As opposed to Pierazzi et al., we explicitly express the variables to optimize inside the optimization problem, since they are the parameters of the practical manipulations.

5.2 Attacks Against Malware Detectors

Many of the white-box attacks in the space of malware classifiers have been previously introduced in this paper, but we revisit here their respective contributions to the literature here, highlighting how they differ from our work.

Kolosnjaji et al. [21] propose an attack against MalConv that leverage padding bytes at the end of the input sample and chooses the best local approximation of each padding value. Demetrio et al. [11] propose the *Partial DOS* practical manipulation against the MalConv classifier. Both strategies are easy to apply, since the *Padding* manipulation does not require any particular effort for being applied, but they are not as effective as the novel manipulations we have proposed, as shown in Figure 7 and Figure 9, in both white-box and black-box settings. Also, the number of bytes altered by padding and partial dos are either too few or placed in locations with zero gradient, hence useless during the optimization.

Kreuk et al. [23] propose an iterative variant of the Fast Gradient Sign Method (FGSM) [14] by manipulating malware inside the feature space imposed by the target network MalConv using both padding and slack space bytes. Similarly, Suciú et al. [33] apply the classic non-iterative FGSM in feature space by adding bytes to slack space between sections. Both strategies alter the sample inside the feature space, reconstructing a real adversarial EXE only at the end of the iteration, which might reduce the adversarial payload effectiveness.

Sharif et al. [31] propose semantics-preserving practical manipulations that alter the code of the input executable, and evaluate against MalConv and the network proposed by Krvcál et al. [22]. For instance, they alter math operations, registers, operand and they add instructions without side-effect on the original behavior of the program. They apply such manipulations at random to each function of the executable, keeping them if the resulting feature vector is aligned with the gradient. Our approach is entirely guided by the gradient of the target function, and it does not leverage

randomness while searching for adversarial examples. Also, our practical manipulations target the structure of the executable rather than its code, minimizing the size of the perturbation.

5.3 Malware Detection Through Machine Learning

We review other techniques that have been produced in the literature to spot malware using machine learning technology. Saxe et al. [30] develop a deep neural network which is trained on top of a feature extraction phase. The authors consider type-agnostic features, such as imports, bytes and strings distributions along with metadata taken from the headers, for a total of 1024 input variables. Kolosnjaji et al. [21] propose to track which APIs are called by a malware, capturing the execution trace using the Cuckoo sandbox,¹² that is a dynamic analysis virtual environment for testing malware. Hardy et al. [15] statically extract which APIs are called by a program, and they train a deep network over this representation. David et al. [10] develop a network that learns new signatures from input malware, by posing the issue as a reconstruction problem. The network infers a new representation of the data, in an end-to-end fashion. These new signatures can be used as input for other machine learning algorithms. Incer et al. [18] try to tackle the issue of an adversarially-robust classifier by imposing monotonic constraints over the features used for the classification tasks. Krčál et al.[22] propose a similar architecture as the one developed by Johns¹³ and Coull et al. [9]: a deep convolutional neural network trained on raw bytes. Both architectures share a first embedding layer, followed by convolutional layers with ReLU activation functions. Krčál et al. use of more fully dense connected layers, with Scaled Exponential Linear Units (SeLU) [19] activation functions.

5.4 Lessons Learned with Packing

Another way to achieve evasion without applying any optimization is leveraging packing techniques, i.e. a binary rewriting technique that compress a program inside another program, and the latter is decompressed at run-time. Initially designed to save on disk space and protect intellectual property, packing is often used to obfuscate the representation of input programs. This causes an increase in the difficulty of studying packed samples, both malign or benign, using static analysis techniques. In this scenario, machine-learning techniques are not a disruptive technology for detecting threats based only on static information, as described by Aghakhani et al. [1]. The authors of the work studied how packing decreases the meaningfulness of different static feature sets, by destroying the original representation. On the other hand, content obfuscation by packing leads to the creation of a new program itself, and it can be seen as a very intrusive way of hiding the malicious content from static analysis-based classifiers.

From the perspective of adversarial machine learning, we are interested in sizing the worst case and evaluate adversarial robustness of machine-learning models by showing that even very small, non-invasive perturbations of the input program can successfully break detection, without the need of packing or obfuscating the whole input program. The goal of our analyses is to demonstrate how brittle learning-based models can be in face of perturbations carefully optimized against them, rather than showing that static code analysis can be bypassed by packing the whole program. We do believe that this is really important, as similar issues may also be found for learning-based models trained on features extracted from dynamic code analysis. In particular, a learning-based model trained on such features may anyway learn to discriminate between legitimate and malware programs based only a small subset of (very discriminant) feature values. In this case, even a small change to such feature values may allowing evading malware detection. For this reason, we believe that understanding and quantifying adversarial robustness of learning-based malware detectors

¹²<https://cuckoosandbox.org/>

¹³<https://www.camlis.org/2017/jeffreyjohns>

may not only unveil different, novel vulnerabilities, but that it also constitutes a very important research direction to improve and design more robust models in this space.

6 LIMITATIONS

We discuss here the main limitations of our work, from the efficacy of such methods in other scenarios, to possible mitigation techniques.

Other Feature Sets. In this work, we showed that manipulations that affect byte-based feature vectors can overthrow end-to-end detectors, by editing only a small portion of the samples. However, we are aware that these approaches do not impact the performances of detectors that extract information from other objects inside the binary, e.g. the imported API, the control flow graph, etc. Also, as shown by the GBDT [3], these manipulations might not be effective against models that partially include byte-based features. However, the attacker can create suitable practical manipulations for addressing these feature sets, and RAMEN can encode them without loss of generality inside the h function (as shown in Section 3.4). This is the case of the strategy proposed by Demetrio et al. [12], where they also consider the injection of benign content through the addition of new sections. This attack strategy is not only able to decrease the performances of the GBDT classifier under 20% of DR, but it also achieve adversarial evasion against well-known antivirus programs hosted on VirusTotal.¹⁴ In particular, the authors shows that 12 of them are weak to adversarial attacks, and 9 of them appear in the Gartner Magic Quadrant for Endpoint Protection.¹⁵

Dynamic Classifiers. We focused on studying the robustness of classifiers that only address the structure of binary programs. However, the techniques we have developed would not affect a classifier that consider also the execution of such samples, extracting features like the sequence of system calls or else. This limitation arises since the content we inject and perturb is not executed at run-time from the program. The attacker would need to apply *binary re-writing techniques* [36], briefly discussed in Section 3.4, and they are practical manipulations specifically crafted for dealing with dynamic features. These manipulations leverage the editing of the code of the program, by inserting new branches or swapping instructions that are semantically equivalent, and RAMEN is in principle general enough for encoding them.

Mitigation. Both the *Full DOS* and *Extend* attacks can be easily patched before the classification step, by filtering out all the content between the magic number MZ and offset 0x3c, plus all the content between offset 0x40 and the header identifier PE. Also, the *Shift* attack can be sanitized by looking at whether the beginning of the first section and the end of the optional header match. The defender can get rid of these manipulations by either erasing the payload or shifting all the content backward and reverting altered section entries. In both cases, the adversarial payload is deleted, and the classification step can be applied without further complications. Since we are interested in finding minimal perturbations that alter the decision process, we focused on less invasive alterations, in contrast to the one proposed by Sharif et al. [31], whose application alters the code section of the input program. However, acknowledging the existence of such manipulations is essential for understanding how to defend against such attacks, and also to understand which feature sets are more useful and stable against adversarial attacks.

¹⁴<https://www.virustotal.com>

¹⁵<https://www.microsoft.com/security/blog/2019/08/23/gartner-names-microsoft-a-leader-in-2019-endpoint-protection-platforms-magic-quadrant/>

7 CONCLUSIONS

We propose RAMEN, a lightweight formalization that encapsulates all the needs of the attacker, with all the practical manipulations applicable in the domain of choice and with all the constraints expressed as a penalty term inside the optimization process. We define and apply new practical manipulations, namely *Full DOS*, *Extend* and *Shift*, crafted for the domain of malware detection of Windows PE programs, and, we offer an overview of all the practical manipulations that address such domain. We show how to recast all the white-box and black-box attacks proposed in the state of the art literature inside RAMEN without loss of generality, and implemented accordingly. We conduct experiments for assessing the results of our new techniques, taking into account state-of-the-art deep learning classifiers, and we produce evasive adversarial malware against them, in both gradient-based (white-box) and gradient-free (black-box) settings. The amount of noise added to the original malware samples is below 2% of the input size of the target network, showing that the attacker does not need to rely on complex manipulations to achieve evasion, but rather apply only a minimal perturbation. The novel white-box attacks shows that the newly-proposed *Extend* and *Shift* attacks are the most effective ones in our experimental analysis. We test the performance of transfer attacks, showing how an attacker can take advantage of only using surrogate models they own, and we highlight how the networks trained on a larger dataset are weak to these strategies. Lastly, we show the attacker is able to decrease the detection rate of remote classifiers, by only sending queries and optimizing the adversarial noise based on such responses.

Future work. We believe that future lines of research should include the study of which feature sets are easier to perturb with practical manipulations to craft adversarial malware and which are not. Since semantics have multiple syntactical representations, it can hardly be inferred by static features, as syntax can be easily twisted to shape adversarial malware that evades detection, hence it is interesting to have an hint on the possible moves of a possible attacker. The formalization we propose is general enough for including also attacks against dynamic and hybrid detectors, since RAMEN is highly modular, and it can comply with the attacker's needs. Also, as previously mentioned, it would be indeed useful to encode more manipulations in this formalization, ranging from perturbations of more complex feature sets to dynamic features. On a different note, it would be interesting to improve the robustness of machine learning malware classifiers by leveraging domain knowledge in the form of constraints and regularizers, e.g., to capture invariant transformations known to domain experts that may modify the input program bytes but preserve its semantics and functionality [24]. This would enable learning robust algorithms against such transformations in a very efficient manner, without the need of performing *adversarial training*, i.e., generating the actual adversarial EXEmples and retrain the model using them. In fact, such adversarial training procedure may anyway remain ineffective due to the high number of dimensions of the input space and variability of the transformations. For this reason, we believe that encoding domain knowledge directly into the learning process may substantially improve model robustness by bridging the gap between models that are learned entirely in a data-driven manner and the design of hand-crafted feature representations for them.

ACKNOWLEDGEMENT

This work was partly supported by the PRIN 2017 project RexLearn (grant no. 2017TWNMH2), funded by the Italian Ministry of Education, University and Research; and by the EU H2020 project ALOHA, under the European Union's Horizon 2020 research and innovation programme (grant no. 780788).

REFERENCES

- [1] H. Aghakhani, F. Gritti, F. Mecca, M. Lindorfer, S. Ortolani, D. Balzarotti, G. Vigna, and C. Kruegel. When malware is packin' heat; limits of machine learning classifiers based on static analysis features. In *Network and Distributed Systems Security (NDSS) Symposium 2020*, 2020.
- [2] H. S. Anderson, A. Kharkar, B. Filar, and P. Roth. Evading machine learning malware detection. *black Hat*, 2017.
- [3] H. S. Anderson and P. Roth. Ember: an open dataset for training static pe malware machine learning models. *arXiv preprint arXiv:1804.04637*, 2018.
- [4] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrncić, P. Laskov, G. Giacinto, and F. Roli. Evasion attacks against machine learning at test time. In H. Blockeel, K. Kersting, S. Nijssen, and F. Železný, editors, *Machine Learning and Knowledge Discovery in Databases (ECML PKDD), Part III*, volume 8190 of *LNCS*, pages 387–402. Springer Berlin Heidelberg, 2013.
- [5] B. Biggio, G. Fumera, and F. Roli. Security evaluation of pattern classifiers under attack. *IEEE Transactions on Knowledge and Data Engineering*, 26(4):984–996, April 2014.
- [6] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.
- [7] R. L. Castro, C. Schmitt, and G. Dreo. Aimed: Evolving malware with genetic programming to evade detection. In *2019 18th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/13th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, pages 240–247. IEEE, 2019.
- [8] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *10th ACM Workshop on Artificial Intelligence and Security, AISec '17*, pages 15–26, New York, NY, USA, 2017. ACM.
- [9] S. E. Coull and C. Gardner. Activation analysis of a byte-based deep neural network for malware classification. In *2019 IEEE Security and Privacy Workshops (SPW)*, pages 21–27. IEEE, 2019.
- [10] O. E. David and N. S. Netanyahu. Deepsign: Deep learning for automatic malware signature generation and classification. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2015.
- [11] L. Demetrio, B. Biggio, G. Lagorio, F. Roli, and A. Armando. Explaining vulnerabilities of deep learning to adversarial malware binaries. *Proceedings of the Third Italian Conference on CyberSecurity (ITASEC)*, 2019.
- [12] L. Demetrio, B. Biggio, G. Lagorio, F. Roli, and A. Armando. Functionality-preserving black-box optimization of adversarial windows malware, 2020.
- [13] A. Demontis, M. Melis, M. Pintor, M. Jagielski, B. Biggio, A. Oprea, C. Nita-Rotaru, and F. Roli. Why do adversarial attacks transfer? Explaining transferability of evasion and poisoning attacks. In *28th USENIX Security Symposium (USENIX Security 19)*. USENIX Association, 2019.
- [14] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [15] W. Hardy, L. Chen, S. Hou, Y. Ye, and X. Li. D4md: A deep learning framework for intelligent malware detection. In *Proceedings of the International Conference on Data Mining (DMIN)*, page 61. The Steering Committee of The World Congress in Computer Science, Computer ..., 2016.
- [16] L. Huang, A. D. Joseph, B. Nelson, B. Rubinstein, and J. D. Tygar. Adversarial machine learning. In *4th ACM Workshop on Artificial Intelligence and Security (AISec 2011)*, pages 43–57, Chicago, IL, USA, 2011.
- [17] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin. Black-box adversarial attacks with limited queries and information. In J. Dy and A. Krause, editors, *35th ICML*, volume 80, pages 2137–2146. PMLR, 2018.
- [18] I. Incer, M. Theodorides, S. Afroz, and D. Wagner. Adversarially robust malware detection using monotonic classification. In *Proceedings of the Fourth ACM International Workshop on Security and Privacy Analytics*, pages 54–63. ACM, 2018.
- [19] G. Klambauer, T. Unterthiner, A. Mayr, and S. Hochreiter. Self-normalizing neural networks. In *Advances in neural information processing systems*, pages 971–980, 2017.
- [20] B. Kolosnjaji, A. Demontis, B. Biggio, D. Maiorca, G. Giacinto, C. Eckert, and F. Roli. Adversarial malware binaries: Evading deep learning for malware detection in executables. In *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 533–537. IEEE, 2018.
- [21] B. Kolosnjaji, A. Zarras, G. Webster, and C. Eckert. Deep learning for classification of malware system call sequences. In *Australasian Joint Conference on Artificial Intelligence*, pages 137–149. Springer, 2016.
- [22] M. Krčál, O. Švec, M. Bálek, and O. Jašek. Deep convolutional malware classifiers can learn from raw executables and labels only. *Sixth International Conference on Learning Representations (ICLR) Workshop*, 2018.
- [23] F. Kreuk, A. Barak, S. Aviv-Reuven, M. Baruch, B. Pinkas, and J. Keshet. Deceiving end-to-end deep learning malware detectors using adversarial examples. *arXiv preprint arXiv:1802.04528*, 2018.
- [24] S. Melacci, G. Ciravegna, A. Sotgiu, A. Demontis, B. Biggio, M. Gori, and F. Roli. Can domain knowledge alleviate adversarial attacks in multi-label classifiers?, 2020.
- [25] M. Melis, A. Demontis, M. Pintor, A. Sotgiu, and B. Biggio. secml: A python library for secure and explainable machine learning, 2019.
- [26] N. Papernot, P. McDaniel, and I. Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016.
- [27] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, ASIA CCS '17*, pages 506–519, New York, NY, USA, 2017. ACM.
- [28] F. Pierazzi, F. Pendlebury, J. Cortellazzi, and L. Cavallaro. Intriguing properties of adversarial ml attacks in the problem space. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1332–1349. IEEE, 2020.

- [29] E. Raff, J. Barker, J. Sylvester, R. Brandon, B. Catanzaro, and C. K. Nicholas. Malware detection by eating a whole exe. In *Workshops at the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [30] J. Saxe and K. Berlin. Deep neural network based malware detection using two dimensional binary program features. In *Malicious and Unwanted Software (MALWARE), 2015 10th International Conference on*, pages 11–20. IEEE, 2015.
- [31] M. Sharif, K. Lucas, L. Bauer, M. K. Reiter, and S. Shintre. Optimization-guided binary diversification to mislead neural networks for malware detection. *arXiv preprint arXiv:1912.09064*, 2019.
- [32] W. Song, X. Li, S. Afroz, D. Garg, D. Kuznetsov, and H. Yin. Automatic generation of adversarial examples for interpreting malware classifiers. *arXiv preprint arXiv:2003.03100*, 2020.
- [33] O. Suciú, S. E. Coull, and J. Johns. Exploring adversarial examples in malware detection. In *2019 IEEE Security and Privacy Workshops (SPW)*, pages 8–14. IEEE, 2019.
- [34] O. Suciú, R. Marginean, Y. Kaya, H. D. III, and T. Dumitras. When does machine learning FAIL? generalized transferability for evasion and poisoning attacks. In *27th USENIX Security Symposium (USENIX Security 18)*, pages 1299–1316, Baltimore, MD, 2018. USENIX Association.
- [35] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- [36] M. Wenzl, G. Merzdovnik, J. Ullrich, and E. Weippl. From hack to elaborate technique—a survey on binary rewriting. *ACM Computing Surveys (CSUR)*, 52(3):1–37, 2019.
- [37] D. Wierstra, T. Schaul, T. Glasmachers, Y. Sun, J. Peters, and J. Schmidhuber. Natural evolution strategies. *Journal of Machine Learning Research*, 15:949–980, 2014.