

Feature-Space Bayesian Adversarial Learning Improved Malware Detector Robustness

Bao Gia Doan¹, Shuiqiao Yang², Paul Montague⁴, Olivier De Vel³, Tamas Abraham⁴, Seyit Camtepe³, Salil S. Kanhere², Ehsan Abbasnejad¹, Damith C. Ranasinghe¹

¹The University of Adelaide, Australia

²The University of New South Wales, Australia

³Data61, CSIRO, Australia

⁴Defence Science and Technology Group, Australia

{giabao.doan, ehsan.abbasnejad, damith.ranasinghe}@adelaide.edu.au, {shuiqiao.yang, salil.kanhere}@unsw.edu.au, seyit.camtepe@data61.csiro.au, olivierdevel@yahoo.com.au, {paul.montague, tamas.abraham}@defence.gov.au

Abstract

We present a new algorithm to train a robust malware detector. Malware is a prolific problem and malware detectors are a front-line defense. Modern detectors rely on machine learning algorithms. Now, the adversarial objective is to devise alterations to the malware code to decrease the chance of being detected whilst *preserving the functionality and realism of the malware*. Adversarial learning is effective in improving robustness but generating functional and realistic adversarial malware samples is non-trivial. Because: i) in contrast to tasks capable of using gradient-based feedback, adversarial learning in a domain without a *differentiable mapping function* from the *problem space* (malware code inputs) to the *feature space* is hard; and ii) it is difficult to ensure the adversarial malware is realistic and functional. This presents a challenge for developing scalable adversarial machine learning algorithms for large datasets at a production or commercial scale to realize robust malware detectors. We propose an alternative; perform adversarial learning in the *feature space* in contrast to the problem space. We *prove* the projection of perturbed, yet valid malware, in the problem space into feature space will always be a subset of adversarials generated in the feature space. Hence, by generating a robust network against feature-space adversarial examples, we inherently achieve robustness against problem-space adversarial examples. We formulate a Bayesian adversarial learning objective that captures the distribution of models for improved robustness. To *explain* the robustness of the Bayesian adversarial learning algorithm, we *prove* that our learning method bounds the difference between the adversarial risk and empirical risk and improves robustness. We show that Bayesian neural networks (BNNs) achieve state-of-the-art results; especially in the False Positive Rate (FPR) regime. Adversarially trained BNNs achieve state-of-the-art robustness. Notably, adversarially trained BNNs are robust against stronger attacks with larger attack budgets by a margin of up to 15% on a recent production-scale malware dataset of more than 20 million samples. Importantly, our efforts create a benchmark for future defenses in the malware domain.

Introduction

We are amidst a meteoric rise in malware incidents worldwide. Malware is responsible for significant damages, both

financial—in billions of dollars (Anderson et al. 2019)—and human costs in loss of life (Eddy and Perlroth 2020). According to statistics from Kaspersky Lab, at the end of 2020, there were an average of 360,000 pieces of malware detected per day (KasperskyLab 2020). The battle against such large incidents of malware remains an ongoing challenge and the need for automated and effective malware detection systems is a research imperative.

Advances in Machine Learning (ML) have led to state-of-the-art malware detectors (Arp et al. 2014; Peng et al. 2012; Harang and Rudd 2020; Raff et al. 2018; Anderson and Roth 2018). But, ML-based models are known to be vulnerable to *adversarial examples*; here, seemingly benign inputs with small perturbations can successfully evade detectors. Although adversarial examples were shown initially in the computer vision domain (Szegedy et al. 2013; Goodfellow, Shlens, and Szegedy 2015; Madry et al. 2018; Biggio and Roli 2018), malware is no exception. Recent attacks have crafted adversarial examples in the malware domain—so-called *adversarial malware*; now, a carefully crafted malware sample with minimal changes to malware code but still able to preserve the *realism* and *functionality* of the malware is able to fool ML-based malware detectors to misclassify them as benign-ware. These attacks pose an emerging threat against ML-based malware detectors (Grosse et al. 2017; Kolosnjaji et al. 2018; Kreuk et al. 2018; Suciu, Coull, and Johns 2019; Pierazzi et al. 2020; Demetrio et al. 2021).

Problem. In general, adversarial learning (Athalye, Carlini, and Wagner 2018) or training with adversarial examples is an effective method to build models robust against adversarial examples. However, generating adversarial malware samples for training, especially at the *production scale* necessary for deployable models, is non-trivial. Because:

- Generation of adversarial examples in the malware domain is confronted with the *inverse feature-mapping problem* where the function mapping from the *problem space* (the discrete space of software code binaries) to the feature space (vectorized features) is non-differentiable (Biggio et al. 2013; Biggio, Fumera, and Roli 2013; Quiring, Maier, and Rieck 2019). Hence, *fast*, gradient-driven methods to derive useful information to craft adversarial samples in the problem space are not suitable.

- The need to enforce malware domain constraints, *realism*, *functionality*, and *maliciousness* on generated perturbations in the problem space is a difficult proposition. Thus, arbitrary changes to the malware binaries are not possible because it could drastically alter the malware in a manner to break the malicious functionality of the binaries or even make it unloadable.

Although efforts to realize robust models on discrete spaces such as discrete image or graph data exist (Lee et al. 2019; Wang et al. 2021), the problem space of malware classification is significantly more challenging due to the imposed constraints in the problem space; the *realism* and *functionality* as well as *maliciousness* of the malware must be maintained. Unfortunately, a method to scale up adversarial training with samples in the problem space to production scale datasets, especially in the case of neural networks, does not exist.

Further, despite extensive work on adversarial ML in general, very few studies have focused on the problem in the context of malware as recently highlighted by Pierazzi et al. (2020), and a comprehensive investigation of robust defense methods in the area remains to be conducted.

Research Questions. Hence, in this study, we seek to answer the following research questions (RQs):

- **RQ1.** How can we overcome the challenging problem of adversarial learning for malware at a *production scale* to realize robust malware detectors against adversarial malware samples?
- **RQ2.** How can we formulate an adversarial learning problem for building robust malware detectors and how can we *explain* the robustness and *benefits*?
- **RQ3.** How robust are adversarially trained malware detectors, especially against problem-space (functional, realistic and malicious) adversarial malware samples?

Our Approach. We argue that a defender is not confronted with the problems we mentioned. Because, we show that constraining the adversarial examples in the problem space to *preserve malware realism, functionality and maliciousness can be turned to an advantage for defenders. The constraints make the perturbed malware in the problem space a subset of the adversarial examples in the feature space. Therefore, designing a robust method against feature-space adversarial examples will inherently be robust against constrained problem-space adversarial examples encapsulating the threats from adversarial malware.*

To construct a formulation to improve the robustness against feature-space adversarial malware examples, and ultimately problem space malware, we propose a Bayesian formulation for adversarially training a neural network: i) with the capability to capture the distribution of models to improve robustness (Liu and Wang 2016; Liu et al. 2019; Ye and Zhu 2018; Wicker et al. 2021; Carbone et al. 2020; Doan et al. 2022); and ii) prove our proposed method of diversified Bayesian neural networks hardened with adversarial training bounds the difference between the adversarial risk and the conventional empirical risk to *theoretically explain* the improved robustness.

Moreover, just recently, security researchers with domain expertise placed significant effort into providing features¹ for malware samples at a production scale of more than 20 million samples (Harang and Rudd 2020; Anderson and Roth 2018)—the SOREL-20M dataset. However, *the robustness of networks built on these extracted features in the face of evasion attacks are yet to be understood.* Therefore, our study to investigate production scale adversarial learning is timely and we focus our efforts to investigate methods using the SOREL-20M dataset.

Our Contributions. To address the problem of building robust malware detectors, we make the following contributions:

1. We *prove* the projection of perturbed yet, valid malware, in the problem space (the discrete space of software code binaries) into the feature space will be a subset of feature-space adversarial examples. Thus, a robust network against feature-space attacks is inherently robust against problem-space attacks. Our work provides *a theoretically justified basis for adversarially training malware detectors in the feature space.* Further, to corroborate our proof, we empirically demonstrate networks trained on feature-space adversarials are robust against functional and realistic problem-space adversarial malware (**RQ1**).
2. Hence, to improve robustness in the *problem space* we propose performing adversarial learning in the *feature space* and formulate a Bayesian Neural Network (BNN) adversarial learning objective that captures the distribution of models for improved robustness. The algorithm is capable of learning from production scale feature-space datasets of up to *20 million samples (RQ1 and RQ2)*.
3. We also *prove* hardening BNNs with adversarial examples bounds the difference between the adversarial risk and the empirical risk to explain the improved robustness (**RQ2**).
4. We empirically demonstrate Bayesian Neural Networks capturing model diversity to improve the performance of malware classifiers and adversarially trained BNNs to generate more robust models against the threat of adversarial malware. Adversarially trained BNNs achieve new benchmarks for state-of-the-art robustness—especially against unseen, stronger, attack samples (**RQ3**).

Scope. Notably, in our study, we focus on Windows Portable Executable (PE) malware for two reasons: i) Windows is the most popular operating system for end-users worldwide, and PE-file malware is the earliest and most studied threat in the wild (Schultz et al. 2001), making a robust method to detect adversarial PE files a significant contribution to security research; and ii) the intuition and methodology behind Windows PE malware can be applied and transferred to other file formats and operating systems, such as PDF malware or malware for Linux and Android systems (see the **Appendix**).

Background and Related Work

Machine Learning Methods in the Malware Domain. Malware detection is moving away from hand-crafted approaches relying on rules toward machine learning (ML)

¹Notably, a negligible computation time of 160 ms, on average, is required to derive vectorized features as described in the **Appendix**.

techniques (Schultz et al. 2001; Saxe and Berlin 2015; Raff et al. 2018; Krčál et al. 2018). Recently, MalConv (Raff et al. 2018) adopted a Convolutional Neural Network (CNN) based architecture design with a learnable, but non-differentiable, embedding space for malware detection from raw byte sequences. The adoption of a CNN for malware detection was also proposed in (Krčál et al. 2018). However, training malware detectors on raw byte sequences (arbitrary number, often millions, of bytes) is computationally expensive and time-consuming. In addition, as we discussed earlier, it is non-trivial to craft realistic adversarial examples on raw byte sequences to realize a robust network on large-scale datasets. Consequently, recent work has employed problem space to feature space mapping functions together with feed-forward neural networks to build benchmark models for the large-scale SOREL-20M dataset (Harang and Rudd 2020).

LUNA (Backes and Nauman 2017) proposed a simple linear Bayesian model for an Android malware detector, which preserves the concept of uncertainty, and shows that it helps to reduce incorrect decisions as well as improve the accuracy of classification. The benefit of a Bayesian classifier is to handle ML tasks from a stochastic perspective, where all weight values of the network are probability distributions. More recently, Nguyen et al. (2021) investigated the application of uncertainty and Bayesian treatment to improve the performance of malware detectors on neural networks.

Adversarial Malware (Adversarial Examples in the Malware Domain). ML-based classifiers are shown to suffer from *evasion* attacks, via *adversarial examples* (Goodfellow, Shlens, and Szegedy 2015). Recently, adversarial examples were demonstrated in the *problem space* (Grosse et al. 2016; Xu, Qi, and Evans 2016; Grosse et al. 2017; Hu and Tan 2017; Kolosnjaji et al. 2018; Kreuk et al. 2018; Suciú, Coull, and Johns 2019). In particular, Kolosnjaji et al. (2018) proposed a method to append bytes to the end of the binary PE file, while Kreuk et al. (2018) exploited the regions within the executable which are not mapped to memory to construct adversarial malware. These methods intend to make modifications that do not affect the intended behavior of the executable. Suciú *et al.* (Suciú, Coull, and Johns 2019) adopted FGSM (Goodfellow, Shlens, and Szegedy 2015) to show the generalization properties and effectiveness of adversarial examples against a CNN-based malware detector, MalConv, trained with small-scale datasets. Suciú, Coull, and Johns (2019) highlighted the threat from adversarial examples as an alternative to evasion techniques such as runtime packing, but showed that models trained on small-scale datasets did not generalize to robust models; hence, *emphasizing the importance of training networks on production-scale datasets*.

Improving Model Robustness. Among methods for improving the robustness of models (Madry et al. 2018; Chen et al. 2020; Fischer et al. 2019), adversarial training (Madry et al. 2018) and its variants are shown to be one of the most effective and popular methods to defend against adversarial examples (Athalye, Carlini, and Wagner 2018). The goal of adversarial training is to incorporate the adversarial search within the training process and, thus, realize robustness against adversarial examples at test time. In particular,

recently, Bayesian adversarial learning has been investigated and adopted in the computer vision domain to propose to improve the robustness of models against adversarial examples (Liu and Wang 2016; Ye and Zhu 2018; Liu et al. 2019; Wicker et al. 2021; Carbone et al. 2020; Doan et al. 2022).

Adversarial learning was explored in the malware domain in (Al-Dujaili et al. 2018) to generate a robust detector for binary encoded malware. However, the computational cost to realize realistic, adversarial raw byte representations is prohibitively expensive (Suciú, Coull, and Johns 2019; Pierazzi et al. 2020) for adversarial learning.

Summary. We recognize that: i) a method capable of scaling up the adversarial training of neural networks in the problem space to production scale datasets does not exist; ii) a Bayesian adversarial learning objective that captures the distribution of models could provide improved robustness; however iii) such a formulation requires overcoming the challenging problem of generating problem-space adversarial examples at production scales.

In what follows, we begin with a problem definition, a theoretical basis for employing feature-space adversarial learning as an alternative to problem-space, followed by the formulation of a Bayesian adversarial learning objective and experimental results validating our claims and demonstrating state-of-the-art performance and robustness.

Problem Definition

Threat model. We assume an attacker with *perfect knowledge* (white-box attacker) (Biggio, Fumera, and Roli 2013), in which the attacker knows *all* parameters including feature set, learning algorithm, loss function, model parameters/hyperparameters, and training data. The reason for considering the strongest, perfect-knowledge adversary is because, even if access to the model is not possible, or the model is not publicly available, an adversary can employ a reverse engineering approach such as (Tramèr et al. 2016; Rolnick and Kording 2020; Carlini, Jagielski, and Mironov 2020) to extract the model. And, defending against such attacks is challenging. The attacker’s objective is to *evade detection*. Their capability is to modify the features at test time.

Problem-Space Attacks. We consider the **problem space** \mathcal{Z} which refers to the *input space* of real objects of a considered domain such as software code binaries. First \mathcal{Z} must be transformed into a compatible format such as numerical vector data (Anderson and Roth 2018; Harang and Rudd 2020) for ML to process. Then, a **feature mapping** is a function $\Phi : \mathcal{Z} \rightarrow \mathcal{X} \subseteq \mathbb{R}^n$ that maps a given problem-space software code binary $\mathbf{z} \in \mathcal{Z}$ to an n-dimensional feature vector $\mathbf{x} \in \mathcal{X}$ in the **feature space** such that $\Phi(\mathbf{z}) = \mathbf{x}$.

Normally, attackers have to apply a transformation on \mathbf{z} to generate \mathbf{z}' such that $\Phi(\mathbf{z}')$ is very close to \mathbf{x}' in the feature space. Formally, given a problem-space object $\mathbf{z} \in \mathcal{Z}$ with label $y \in \mathcal{Y}$, the goal of the adversary is to find the *transformation* $\mathbf{T} : \mathcal{Z} \rightarrow \mathcal{Z}$ (e.g. addition, removal, modification) such that $\mathbf{z}' = \mathbf{T}(\mathbf{z})$ is classified as a class $t \neq y$. In the malware domain, the adversary has to search in the problem space that approximately follows the gradient in the feature space. However, this is a major challenge that complicates the application of gradient-driven methods to the

problem-space attacks— so-called **inverse feature-mapping problem** (Quiring, Maier, and Rieck 2019; Biggio et al. 2013; Pierazzi et al. 2020) where the function Φ in the software domain—our focus—is typically not invertible and not differentiable, *i.e.* there is no one-to-one mapping from the adversarial examples in the feature space $\mathbf{x} + \delta$ to the corresponding adversarial problem-space object \mathbf{z}' . In addition, the generated object $\mathbf{T}(\mathbf{z})$ must be realistic and valid (Suciu, Coull, and Johns 2019). Thus, the search for adversarial examples in the problem space (software) cannot be a purely gradient-based method, hindering the adoption of well-known adversarial attacks in other domains such as computer vision. To achieve a realistic adversarial objective, the search for adversarial examples in the problem space has to be constrained in **problem-space constraints** denoted by Ω . We remark that the constraints on the problem space are well defined and can be found in (Biggio and Roli 2018; Quiring, Maier, and Rieck 2019; Xu, Qi, and Evans 2016; Pierazzi et al. 2020), we mentioned here, for completeness, that there are at least four main types of problem-space constraints including Preserved semantics, Plausibility, Robustness to Processing and Available Transformation explained in detail by Pierazzi et al. (2020).

Feature-Space Attacks. To alleviate the problems with problem space attacks, we propose an alternative that uses feature space. We note that all definitions of feature-space attacks are well defined and consolidated in related work (Biggio and Roli 2018; Carlini and Wagner 2017; Grosse et al. 2017). In this paper, we use a popular feature mapping function provided in the EMBER dataset (Anderson and Roth 2018) to map raw bytes of software to a vector of $n = 2381$ features. A **feature-space attack** is then to modify a feature-space object $\mathbf{x} \in \mathcal{X}$ to become another object $\mathbf{x}' = \mathbf{x} + \delta$ where δ is the added perturbation crafted with an *attack objective function* to misclassify \mathbf{x}' into another class $t \neq y$ where $y \in \mathcal{Y}$ is the ground-truth label of \mathbf{x} . We note that in the malware domain (a binary classification task), the intuition of the attackers is to make the malware be recognized as benign ware. These modifications has to follow **feature-space constraints**. We denote the constraints on feature-space modifications by Υ . Given a sample $\mathbf{x} \in \mathcal{X}$, the feature-space modification, or perturbation δ must satisfy Υ . This constraint Υ reflects the realistic requirements of problem-space objects. In the malware domain, feature perturbations δ can be constrained $\delta_{lb} \leq \delta \leq \delta_{ub}$ (Pierazzi et al. 2020).

Theoretical Basis For Feature-Space Adversarial Learning

We highlight that the realistic assumption of problem-space attacks makes the constraints imposed by Ω stricter or equal to those imposed by Υ (illustrated in Figure 1). Following the necessary condition for problem-space adversarial examples as stated in Pierazzi et al. (2020), we have:

Lemma 1. If there exists an adversarial example in the problem space (\mathbf{z}') that satisfies the constraints Ω , then there will be a corresponding adversarial example in the feature space (\mathbf{x}') under the constraints Υ . More formally, by abusing notation from model theory to use \models to indicate an instance

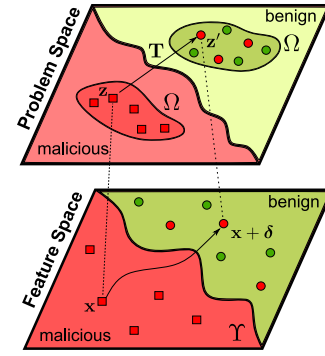


Figure 1: Illustrative example of adversarial examples. The adversarial example $\mathbf{x} + \delta$ is derived from \mathbf{x} in the feature space and its projection to problem-space constraints (which is more restrictive) determined by Ω is \mathbf{z}' . The color in the background illustrates the decision regions where red color is for malware and green is for benign programs. The solid arrow in Feature Space represents the gradient-based attack to transform a malware \mathbf{x} to $\mathbf{x} + \delta$, projected to the problem-space constraints as \mathbf{z}' to be misdetected as a benign program.

“satisfies” constraints, and write $\mathbf{z}' \models \Omega$ and $\mathbf{x}' \models \Upsilon$, we have:

$$\begin{aligned} \exists \mathbf{z}' : \mathbf{z}' \models \Omega, \quad p(y | \Phi(\mathbf{z}'), \theta) &= p(y | \Phi(\mathbf{T}(\mathbf{z})), \theta), \\ p(y | \Phi(\mathbf{T}(\mathbf{z})), \theta) &< 0.5 \\ \Rightarrow \exists \mathbf{x}' = \mathbf{x} + \delta : \mathbf{x}' \models \Upsilon, \quad p(y | \mathbf{x}', \theta) &< 0.5 \end{aligned}$$

where \mathbf{T} is the transformation in the problem space to craft adversarial examples, $p(y | \mathbf{x}, \theta) = \text{sigmoid}(f(\mathbf{x}; \theta))$ is the output of a sigmoid function applied to the output of the neural networks f parameterized by θ , $p(y | \mathbf{x}, \theta) = 0.5$ is the threshold for malware detection where the predicting $p(y | \mathbf{x}, \theta) = 0$ is recognized as benign whilst $p(y | \mathbf{x}, \theta) = 1$ indicates a malware, Ω, Υ are, respectively, the problem-space and feature-space constraints, and $\Phi(\cdot)$ is the function that maps the problem space to feature space.

The proof of Lemma 1 is in **Appendix**. From **Lemma 1**, if there exists an attack in the problem space, then there exists a corresponding attack in the feature space. By contraposition, if there does not exist an attack in the feature space, there does not exist an attack in the problem space. However, we know that the opposite is not true: if there does not exist an attack in the problem space (e.g. due to functionality), there still exists an attack in the feature space. Thus, we can derive:

Corollary 1. The adversarial examples generated from constrained problem-space adversarial examples (imposed by Ω) are in a subset of feature-space adversarial examples (imposed by Υ).

Corollary 2. Detectors robust against feature-space adversarial examples (imposed by Υ) are robust against constrained problem-space adversarial examples (imposed by Ω).

Built upon these Corollaries, we propose to find a learning method robust against *feature-space adversarial malware*. On the one hand, adversarial training (Madry et al. 2018) and its variants are shown to be one of the most effective and popular methods to defend against adversarial examples (Athalye, Carlini, and Wagner 2018). On the other hand,

Bayesian neural networks (MacKay 1992; Ritter, Botev, and Barber 2018; Izmailov et al. 2021) with distributions placed over their weights and biases enabling the principled quantification of the uncertainty of their predictions are shown to be a robust method against adversarial examples. Thus, in this paper, demonstrating that robustness against feature-space adversarial examples is inherently robust against problem-space real malware. We propose to incorporate adversarial training with Bayesian neural networks to seek the first principled method of Bayesian adversarial learning to realize a robust malware detector without the difficulties of inverse feature-mapping and preserving semantics and functionalities of real malware samples. We name our method Adv-MalBayes, and the method is efficient enough to be scaled up to a large production scale of adversarial training data of 20 million adversarial samples with the pre-extracted feature set of SOREL-20M dataset (Harang and Rudd 2020).

Bayesian Formulation for Adversarial Learning

The goal of Bayesian adversarial learning is to find the posterior distribution using Bayes theorem:

$$p(\theta | \mathcal{D}_{\text{adv}}) = \prod_{(\mathbf{x}_{\text{adv}}, y) \sim \mathcal{D}_{\text{adv}}} p(y | \mathbf{x}_{\text{adv}}, \theta) p(\theta) / Z$$

where Z is the normalizer, \mathcal{D}_{adv} is the adversarial dataset obtained by generating adversarial examples from the benign dataset D using adversarial generation such as Eq. (1).

We consider $p(y | \mathbf{x}_{\text{adv}}, \theta) = \text{sigmoid}(f(\mathbf{x}_{\text{adv}}; \theta))$ to produce a binary prediction in malware detection. Notably, Eq. (1) is the Expectation-over-Transformation (EoT) PGD attack (Athalye et al. 2018; Zimmermann 2019), which is slightly different from the usual PGD attack (Madry et al. 2018). As has been highlighted in Zimmermann (2019), the EoT attack is better able to estimate the gradient of the stochastic Bayesian models:

$$\mathbf{x}^{t+1} = \Pi_{\varepsilon_{\text{max}}} \{ \mathbf{x}^t + \alpha \cdot \text{sign}(\mathbb{E}_{\theta} [\nabla_{\mathbf{x}} \ell(f(\mathbf{x}^t; \theta), y_o)]) \}. \quad (1)$$

where ε_{max} is the maximum attack budget, $\Pi_{\varepsilon_{\text{max}}}$ is the projection to the set $\{ \mathbf{x} | \|\mathbf{x} - \mathbf{x}_o\|_{\infty} \leq \varepsilon_{\text{max}} \}$, ℓ is the loss function (typically cross-entropy), f is the neural network, \mathbf{x} is the input, θ is the network parameter, and y is the ground-truth label. In this attack, an attacker starts from $\mathbf{x}^0 = \mathbf{x}_o$ and conducts projected gradient descent iteratively to update the adversarial example.

However, as highlighted in Izmailov et al. (2021), the posterior over a Bayesian neural network is extremely high-dimensional, non-convex and intractable. Thus, we need to resort to approximations to find the posterior distribution. In this work, we propose using Stein Variational Gradient Descent (SVGD) (Liu and Wang 2016) for two reasons. First, this approach learns multiple *network parameter particles* in parallel for faster convergence. Second, there is a *repulsive factor* in the method to encourage the diversity of parameter particles that helps to prevent mode collapse — a challenge of posterior approximation. To further demonstrate the robustness of our chosen Bayesian method, we compare Adv-MalBayes with previous BNNs (Liu et al. 2019) in the **Appendix Table 8**.

We consider n samples from the posterior (*i.e.* parameter particles). The variational bound is minimized when gradient

descent is modified as:

$$\begin{aligned} \theta_i &= \theta_i - \varepsilon_i \hat{\phi}^*(\theta_i) \\ \text{with } \hat{\phi}^*(\theta) &= \sum_{j=1}^n [k(\theta_j, \theta) \nabla_{\theta_j} \ell(f(\mathbf{x}_{\text{adv}}; \theta_j), y) \\ &\quad - \frac{\gamma}{n} \nabla_{\theta_j} k(\theta_j, \theta)]. \end{aligned}$$

Here, θ_i is the i th particle, $k(\cdot, \cdot)$ is a kernel function that measures the similarity between particles and γ is a hyper-parameter. The parameter particles are encouraged to be dissimilar to capture more diverse samples from the posterior thanks to the kernel function. This is controlled by a hyper-parameter γ to manage the trade-off between diversity and loss minimization. Following (Liu and Wang 2016), we use the RBF kernel $k(\theta, \theta') = \exp(-\|\theta - \theta'\|^2 / 2h^2)$ and take the bandwidth h to be the median of the pairwise distances of the set of parameter particles at each training iteration.

At the inference stage, given the test data point \mathbf{x}^* , we can get the prediction by approximating the posterior using the Monte Carlo samples as:

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathcal{D}_{\text{adv}}) &= \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathcal{D}_{\text{adv}}) d\theta \\ &\approx \frac{1}{n} \sum_{i=1}^n p(y^* | \mathbf{x}^*, \theta_i), \quad \theta_i \sim p(\theta | \mathcal{D}_{\text{adv}}), \end{aligned}$$

where θ_i is an individual parameter particle. Notably, we acknowledge that it is critical to have diverse parameter particles. Averaging over diverse and uncorrelated predictors was shown to improve network performance (Jacobs et al. 1991; Wolpert 1992; Breiman 1996). In the adversarial setting, when integrating out the parameters in our Bayesian formulation, we implicitly remove the vulnerabilities arising from a single choice of parameter existing in traditional neural networks, and hence improve the robustness.

Adversarial Risk is Bounded with the Bayesian Formulation

In this section, to explain the robustness of the Bayesian adversarial learning method that we propose, we prove that training the network with the Bayesian adversarial learning method bounds the difference between the adversarial risk and the empirical risk. This is important, because, now the risk of misclassification on adversarial examples is as the same as that of benign ones; hence eliminating the vulnerability of adversarial examples and reduce the risk of misclassification of adversarial examples to the generalization ability of the classifier. Notably, improving the generalization ability of the classifier is not our focus.

In this context, we make no specific assumption on the distribution of either the adversarial examples or the perturbations, to provide a generic defense approach. The only assumption we make is that the distribution of the data and the corresponding adversarial examples are sufficiently close. This is a mild and reasonable assumption because the idea of adversarial learning is that the added perturbation does not change the perceived samples or the distribution of the samples. Thus, we consider

the bound of $|R_{adv} - R|$ where the empirical risk $R = \mathbb{E}_{\theta} [\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathbb{E}_{y' \sim p(y|\mathbf{x}, \theta)} [\mathbb{I}(y = y')]]]$ and the adversarial risk $R_{adv} = \mathbb{E}_{\theta} [\mathbb{E}_{(\mathbf{x}_{adv}, y) \sim \mathcal{D}_{adv}} [\mathbb{E}_{y' \sim p(y|\mathbf{x}_{adv}, \theta)} [\mathbb{I}(y = y')]]]$

Proposition 1. *The difference between the adversarial risk (denoted by R_{adv}) and the empirical risk (denoted by R) of a classifier when trained on the observed training set and its adversarial counterparts is bounded, i.e.*

$$|R_{adv} - R| \leq \tau,$$

$$\text{where } \tau = 1 - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\exp \left(\mathbb{E}_{\theta} [r_{\theta}(\mathbf{x}, \mathbf{x}_{adv}, y)] \right) \right],$$

$$r_{\theta}(\mathbf{x}, \mathbf{x}_{adv}, y) = \sum_c^K p(y = c | \mathbf{x}, \theta) \log(p(y = c | \mathbf{x}_{adv}, \theta)).$$

Here, \mathbf{x}_{adv} denotes the adversarial example obtained from \mathbf{x} .

We can see that the difference between the empirical risk and the adversarial risk is minimized when the upper bound is minimized. Notably, as we know that $1 - \exp(-z)$ is a monotonically increasing function, and $1 - \exp(-z) \leq z$, to avoid computational instabilities and gradient saturation, we consider minimizing the upper bound without the exponential function. Thus, to minimize the upper bound, our main learning objective (in Algorithm 1 in the **Appendix**) is to:

Minimize cross entropy for the adversarial examples. This corresponds to matching the prediction from the adversarial data to that of the observations. Since (\mathbf{x}, y) is given in the training, we simply minimize the entropy of the adversarials.

Sketch of the Proof. We simplify the difference between the risks by considering that the difference between individual mistakes is smaller than their product, i.e.

$$\begin{aligned} & \mathbb{E}_{y_1 \sim p(y|\mathbf{x}, \theta)} [\mathbb{E}_{y_2 \sim p(y|\mathbf{x}_{adv}, \theta)} [\mathbb{I}[y \neq y_1] - \mathbb{I}[y \neq y_2]]] \\ & \leq \mathbb{E}_{y' \sim p(y|\mathbf{x}_{adv}, \theta)} [\mathbb{E}_{y'' \sim p(y|\mathbf{x}_{adv}, \theta)} [\mathbb{I}[y_1 \neq y_2]]] \\ & \leq 1 - \sum_{c=1}^K p(y = c | \mathbf{x}, \theta) p(y = c | \mathbf{x}_{adv}, \theta). \end{aligned}$$

We then use Jensen’s inequality when using $\exp(\log(\cdot))$ to obtain the upper bound. The complete proof is provided in **Appendix**. We empirically evaluate this difference of risk and illustrate the results in Figure 4 in the **Appendix**.

Experiments and Results

Classifiers. To validate our proposed method Adv-MalBayes, we conduct experiments on different neural networks. We employ the Feed Forward Neural Network (FFNN) classifier provided in the SOREL-20M dataset (Harang and Rudd 2020). This network architecture is also used for the experiments on the EMBER dataset (Anderson and Roth 2018). Our network implementation uses the default configuration provided in (Harang and Rudd 2020). We also adopt the architecture of FFNN to design the Bayesian Neural Network (BNN). The details of the network architecture are in **Appendix**. Then, we harden the FFNN and BNN with adversarial examples to generate the Adv-FFNN model and Adv-MalBayes. In addition, we also employ baseline networks including LightGBM (Anderson and Roth 2018) and

MalConv (Raff et al. 2018) for comparison. We compare their performance on malware datasets (no attacks) and its adversarial counterparts (adversarial malware designed to evade detectors) to evaluate the *detector performance* and *robustness*. The values of the attack budgets used for training and testing are detailed in Table 4 in the **Appendix**.

Datasets. In this paper, we use the two largest publicly available corpora for malware detection, namely:

- The production scale dataset Sophos AI **SOREL-20M** (Harang and Rudd 2020) containing 20 million pre-extracted samples.
- **EMBER** (Anderson and Roth 2018) dataset designed to be *more challenging* for ML-based classifiers.

We detail these datasets in **Appendix**.

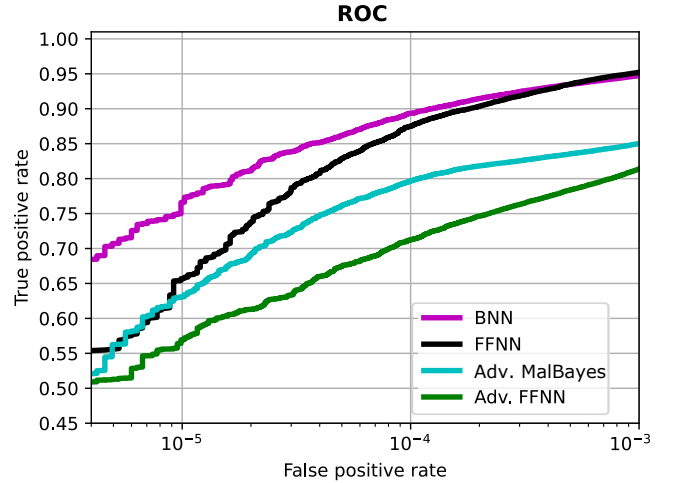


Figure 2: Performance of neural network classifiers in the absence of adversarial examples in the SOREL-20M dataset.

Results. We present our results by reporting: i) performance of the given classifiers on malware detection tasks (no attacks setting) using ROC (receiver operating characteristic curve); and ii) robustness (under evasion attacks with adversarial malware). We detail these metrics in the **Appendix**.

Performance (no attacks). Performance of the classifiers in the absence of attacks are shown in Figure 2 with additional details reported in Table 7 in the **Appendix**. The ROC curves in Figure 2 report the True Positive Rate (i.e. the percentage of correctly-classified malware samples) as a function of the False Positive Rate (FPR, i.e. the percentage of misclassified benign samples) for each classifier. From the figure, we can see that Bayesian neural networks of the same network architecture as FFNNs achieve better performance (compare BNN vs. FFNN and Adv-MalBayes vs. Adv-FFNN). Notably, the BNNs outperformed the FFNN counterparts with a large margin in the detection rate (of up to 20%) under low-FPR regimes. Notably, in Table 7 in the **Appendix**, we also show that BNNs built on feature-space samples achieve better performance compared with the popular ML-based malware detector built on problem-space samples (MalConv) (Raff et al. 2018) and its recently updated version in AAI-21 (MalConv w/ GCG) (Raff et al. 2021).

Robustness (against Feature-Space Adversarial Examples). To evaluate the robustness of the investigated classifiers, we apply the PGD attack from Equation (1) on malware samples with increasing attack budgets. Results for the robustness of given classifiers under different attack budgets are reported in Table 1. Notably, Adv-MalBayes outperforms the adversarially trained FFNN on both the production scale (SOREL-20M) and challenging (EMBER) datasets, especially under increasing attack budgets. This is significant because the problem with malware is that they are evolving extremely fast *e.g.* there are hundreds of thousands of new malware samples every day (KasperskyLab 2020). Further, results in Figure 3 in the **Appendix** illustrate, as expected and in line with the findings in the literature (Madry et al. 2018; Carlini and Wagner 2017; Goodfellow, Shlens, and Szegedy 2015), the adversarially trained networks are significantly more robust than non-adversarially trained counterparts.

Dataset	Networks	Attack budget					
		0	0.03	0.05	0.1	0.2	0.3
SOREL-20M	Adv-FFNN	95.38	93.31	89.92	47.74	17.34	13.3
	Adv-MalBayes	95.52	94.20	90.53	62.86	25.42	23.10
EMBER	Adv-FFNN	86.88	82.44	79.48	64.00	51.32	42.77
	Adv-MalBayes	89.17	86.73	84.79	78.03	63.06	52.63

Table 1: Robustness of networks against adversarial malware generated with increasing attack budgets.

Robustness (against Problem-Space Adversarial Malware). In this section, we evaluate the robustness of different networks against functional, malicious and real adversarial malware in the problem space. We employ two evaluation sets. Set A includes real malware collected from a previous study (Mantovani et al. 2020) and includes 7137 virus samples. We generate the real *adversarial* malware samples by utilizing the constant padding attack method proposed by (Fleshman 2019) used to win the machine learning static evasion competition (DEFCON 2019). In particular, 100,000 constant bytes valued `0xA9` were added to a new section of PE files to ensure the malicious functionality is not altered. The results in Table 2 show that this attack can significantly degrade the performance of the popular ML-based malware detector MalConv (Raff et al. 2018), however the LightGBM (Anderson and Roth 2018) model is still robust against this attack (confirming the previous result obtained by (Fleshman 2019)). Set B consists of the recent release by (Erdemir et al. 2021). This includes 1001 real adversarial malware samples generated using the Greedy Attack method shown to be stronger than the constant padding attack (Fleshman 2019). The results reported in Table 2 show that the Greedy Attack successfully fools the LightGBM model and downgrades its robustness to 11.2%. Notably, evaluations under both sets show the adversarially trained networks on *feature-space* adversarial samples (*i.e.* Adv-FFNN and Adv-

MalBayes) maintained their robustness. Importantly, Adv-MalBayes achieved very high *robustness under both attack datasets and is a clear demonstration of the effectiveness of our approach and the validity of the theoretical basis for training with feature-space adversarial samples.*

	Light-GBM	MalConv	FFNN	BNN	Adv-FFNN	Adv-MalBayes
Set A	92.5%	29.2%	69.5%	72.5%	92.6%	99.9%
Set B	11.2%	- ¹	74.9%	83.1%	91.8%	99.9%

¹ The released set is vectorized features, not applicable for MalConv.

Table 2: Comparing the robustness of detectors against *real* and *unseen* adversarial malware (problem-space attacks).

The Impact of Number of Parameter Particles. We investigate the contribution of the number of parameter particles to the robustness of the networks and report the results in Table 6 in the **Appendix**. The robustness of the BNNs is improved when more particles capable of modeling the multi-modal posterior are employed. Thus, increasing the number of parameter particles may further improve the network’s robustness.

Transferability of Robustness. We also evaluate the robustness of the BNN trained on PGD L_∞ , and its transferability to other attacks, such as FGSM (Goodfellow, Shlens, and Szegedy 2015). Table 3 shows that the network trained with PGD L_∞ is robust to other attacks, in line with (Madry et al. 2018), where PGD L_∞ is considered as the ‘universal’ attack. Consequently, we can expect our method to improve robustness against a wide range of other adversarial example generation methods (Suciu, Coull, and Johns 2019; Kolosnjaji et al. 2018; Kreuk et al. 2018) adopting the FGSM method to attempt to generate problem space malware samples.

Adv-MalBayes Networks	Attack budget					
	0	0.03	0.05	0.1	0.2	0.3
PGD L_∞	96.29	94.97	92.19	69.96	35.20	30.79
FGSM	-	95.28	94.87	95.24	93.78	92.20

Table 3: Results demonstrating the transferability of robustness to different attack methods. The evaluated model was trained on PGD L_∞ .

Conclusion

We proved and demonstrated that training a robust malware detector on feature-space adversarial examples inherently generates robustness against problem-space malware samples. Subsequently, we proposed a Bayesian adversarial learning objective in the feature space to realize a robust malware detector in the problem space. Additionally, we explain the improved performance by proving that our proposed method bounds the difference between adversarial risk versus empirical risk to improve robustness and show the benefits of a BNN as a defense method (see **Appendix**). Our empirical results, including a production scale dataset, demonstrates new state-of-the-art *performance* and *robustness* benchmarks.

Acknowledgments

This research was supported by the Next Generation Technologies Fund from the Defence Science and Technology Group, Australia. We also thank Dr. Sharif Abuadba for supporting us in collecting malware samples for the project.

References

- Al-Dujaili, A.; Huang, A.; Hemberg, E.; and O'Reilly, U.-M. 2018. Adversarial deep learning for robust detection of binary encoded malware. In *IEEE Security and Privacy Workshops (SPW)*.
- Anderson, H. S.; and Roth, P. 2018. Ember: an open dataset for training static PE malware machine learning models. *arXiv preprint arXiv:1804.04637*.
- Anderson, R.; Barton, C.; Böhme, R.; Clayton, R.; Ganán, C.; Grasso, T.; Levi, M.; Moore, T.; and Vasek, M. 2019. Measuring the changing cost of cybercrime. In *Workshop on the Economics of Information Security (WEIS)*.
- Arp, D.; Spreitzenbarth, M.; Hubner, M.; Gascon, H.; Rieck, K.; and Siemens, C. 2014. Drebin: Effective and explainable detection of android malware in your pocket. In *Network and Distributed System Security Symposium (NDSS)*.
- Athalye, A.; Carlini, N.; and Wagner, D. 2018. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International Conference on Machine Learning (ICML)*.
- Athalye, A.; Engstrom, L.; Ilyas, A.; and Kwok, K. 2018. Synthesizing robust adversarial examples. In *International Conference on Machine Learning (ICLR)*.
- Backes, M.; and Nauman, M. 2017. LUNA: Quantifying and Leveraging Uncertainty in Android Malware Analysis through Bayesian Machine Learning. In *IEEE European Symposium on Security and Privacy (Euro S&P)*.
- Biggio, B.; Corona, I.; Maiorca, D.; Nelson, B.; Šrncić, N.; Laskov, P.; Giacinto, G.; and Roli, F. 2013. Evasion attacks against machine learning at test time. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD)*.
- Biggio, B.; Fumera, G.; and Roli, F. 2013. Security evaluation of pattern classifiers under attack. *IEEE Transactions on Knowledge and Data Engineering*, 26(4): 984–996.
- Biggio, B.; and Roli, F. 2018. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84: 317–331.
- Breiman, L. 1996. Bagging predictors. *Machine learning*, 24(2): 123–140.
- Carbone, G.; Wicker, M.; Laurenti, L.; Patane, A.; Bortolussi, L.; and Sanguinetti, G. 2020. Robustness of Bayesian Neural Networks to Gradient-Based Attacks. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Carlini, N.; Jagielski, M.; and Mironov, I. 2020. Cryptanalytic extraction of neural network models. In *CRYPTO*.
- Carlini, N.; and Wagner, D. 2017. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy (S&P)*.
- Chen, Y.; Wang, S.; She, D.; and Jana, S. 2020. On training robust PDF malware classifiers. In *USENIX Security Symposium*.
- DEFCON. 2019. Machine Learning Static Evasion Competition. <https://www.elastic.co/blog/machine-learning-static-evasion-competition>. Accessed: 2022-08-09.
- Demetrio, L.; Biggio, B.; Lagorio, G.; Roli, F.; and Armando, A. 2021. Functionality-preserving black-box optimization of adversarial windows malware. *IEEE Transactions on Information Forensics and Security*, 16: 3469–3478.
- Doan, B. G.; Abbasnejad, E. M.; Shi, J. Q.; and Ranasinghe, D. 2022. Bayesian Learning with Information Gain Provably Bounds Risk for a Robust Adversarial Defense. In *International Conference on Machine Learning (ICML)*.
- Eddy, M.; and Perloth, N. 2020. <https://www.nytimes.com/2020/09/18/world/europe/cyber-attack-germany-ransomware-death.html>. Accessed: 2022-12-01.
- Erdemir, E.; Bickford, J.; Melis, L.; and Aydore, S. 2021. Adversarial robustness with non-uniform perturbations. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Fischer, M.; Balunovic, M.; Drachler-Cohen, D.; Gehr, T.; Zhang, C.; and Vechev, M. 2019. D12: Training and querying neural networks with logic. In *International Conference on Machine Learning (ICML)*.
- Fleshman, W. 2019. Evading Machine Learning Malware Classifiers. <https://towardsdatascience.com/evading-machine-learning-malware-classifiers-ce52dabdb713>. Accessed: 2022-08-09.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. In *International Conference on Learning Representations (ICLR)*.
- Grosse, K.; Papernot, N.; Manoharan, P.; Backes, M.; and McDaniel, P. 2016. Adversarial perturbations against deep neural networks for malware classification. *arXiv preprint arXiv:1606.04435*.
- Grosse, K.; Papernot, N.; Manoharan, P.; Backes, M.; and McDaniel, P. 2017. Adversarial examples for malware detection. In *European Symposium on Research in Computer Security (ESORICS)*.
- Harang, R.; and Rudd, E. M. 2020. SOREL-20M: A Large Scale Benchmark Dataset for Malicious PE Detection. *arXiv preprint arXiv:2012.07634*.
- Hu, W.; and Tan, Y. 2017. Generating adversarial malware examples for black-box attacks based on GAN. *arXiv preprint arXiv:1702.05983*.
- Izmailov, P.; Vikram, S.; Hoffman, M. D.; and Wilson, A. G. 2021. What Are Bayesian Neural Network Posteriors Really Like? In *International Conference on Machine Learning (ICML)*.
- Jacobs, R. A.; Jordan, M. I.; Nowlan, S. J.; and Hinton, G. E. 1991. Adaptive mixtures of local experts. *Neural computation*, 3(1): 79–87.
- KasperskyLab. 2020. The number of new malicious files detected every day increases by 5.2% to 360,000 in 2020. https://www.kaspersky.com/about/press-releases/2020_the

- number-of-new-malicious-files-detected-every-day-increases-by-52-to-360000-in-2020. Accessed: 2022-04-01.
- Kolosnjaji, B.; Demontis, A.; Biggio, B.; Maiorca, D.; Giacinto, G.; Eckert, C.; and Roli, F. 2018. Adversarial malware binaries: Evading deep learning for malware detection in executables. In *European Signal Processing Conference (EU-SIPCO)*.
- Krčál, M.; Švec, O.; Bálek, M.; and Jašek, O. 2018. Deep convolutional malware classifiers can learn from raw executables and labels only. In *International Conference on Learning Representations (ICLR) Workshop*.
- Kreuk, F.; Barak, A.; Aviv-Reuven, S.; Baruch, M.; Pinkas, B.; and Keshet, J. 2018. Deceiving end-to-end deep learning malware detectors using adversarial examples. *arXiv preprint arXiv:1802.04528*.
- Lee, G.-H.; Yuan, Y.; Chang, S.; and Jaakkola, T. 2019. Tight certificates of adversarial robustness for randomly smoothed classifiers. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Liu, Q.; and Wang, D. 2016. Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm. *Advances in Neural Information Processing Systems (NIPS)*.
- Liu, X.; Li, Y.; Chongruo, W.; and Cho-Jui, H. 2019. ADV-BNN: Improved Adversarial Defense Through Robust Bayesian Neural Network. In *International Conference on Learning Representations (ICLR)*.
- MacKay, D. J. 1992. A practical Bayesian framework for backpropagation networks. *Neural computation*, 4(3): 448–472.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations (ICLR)*.
- Mantovani, A.; Aonzo, S.; Ugarte-Pedrero, X.; Merlo, A.; and Balzarotti, D. 2020. Prevalence and Impact of Low-Entropy Packing Schemes in the Malware Ecosystem. In *Network and Distributed System Security Symposium (NDSS)*.
- Nguyen, A. T.; Raff, E.; Nicholas, C.; and Holt, J. 2021. Leveraging Uncertainty for Improved Static Malware Detection Under Extreme False Positive Constraints. In *International Joint Conferences on Artificial Intelligence (IJCAI) Workshop*.
- Peng, H.; Gates, C.; Sarma, B.; Li, N.; Qi, Y.; Potharaju, R.; Nita-Rotaru, C.; and Molloy, I. 2012. Using probabilistic generative models for ranking risks of android apps. In *ACM Conference on Computer and Communications Security (CCS)*.
- Pierazzi, F.; Pendlebury, F.; Cortellazzi, J.; and Cavallaro, L. 2020. Intriguing properties of adversarial ml attacks in the problem space. In *IEEE Symposium on Security and Privacy (S&P)*.
- Quiring, E.; Maier, A.; and Rieck, K. 2019. Misleading authorship attribution of source code using adversarial learning. In *USENIX Security Symposium*.
- Raff, E.; Barker, J.; Sylvester, J.; Brandon, R.; Catanzaro, B.; and Nicholas, C. K. 2018. Malware detection by eating a whole exe. In *AAAI Conference on Artificial Intelligence Workshop*.
- Raff, E.; Fleshman, W.; Zak, R.; Anderson, H. S.; Filar, B.; and McLean, M. 2021. Classifying sequences of extreme length with constant memory applied to malware detection. In *AAAI Conference on Artificial Intelligence*.
- Ritter, H.; Botev, A.; and Barber, D. 2018. A scalable laplace approximation for neural networks. In *International Conference on Learning Representations (ICLR)*.
- Rolnick, D.; and Kording, K. 2020. Reverse-engineering deep ReLU networks. In *International Conference on Machine Learning (ICML)*.
- Saxe, J.; and Berlin, K. 2015. Deep neural network based malware detection using two dimensional binary program features. In *International Conference on Malicious and Unwanted Software (MALWARE)*.
- Schultz, M. G.; Eskin, E.; Zadok, F.; and Stolfo, S. J. 2001. Data mining methods for detection of new malicious executables. In *IEEE Symposium on Security and Privacy (S&P)*.
- Suciu, O.; Coull, S. E.; and Johns, J. 2019. Exploring adversarial examples in malware detection. In *IEEE Security and Privacy Workshops (SPW)*.
- Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- Tramèr, F.; Zhang, F.; Juels, A.; Reiter, M. K.; and Ristenpart, T. 2016. Stealing machine learning models via prediction apis. In *USENIX Security Symposium*.
- Wang, B.; Jia, J.; Cao, X.; and Gong, N. Z. 2021. Certified robustness of graph neural networks against adversarial structural perturbation. In *ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD)*.
- Wicker, M.; Laurenti, L.; Patane, A.; Chen, Z.; Zhang, Z.; and Kwiatkowska, M. 2021. Bayesian Inference with Certifiable Adversarial Robustness. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Wolpert, D. H. 1992. Stacked generalization. *Neural networks*, 5(2): 241–259.
- Xu, W.; Qi, Y.; and Evans, D. 2016. Automatically evading classifiers. In *Network and Distributed System Security Symposium (NDSS)*.
- Ye, N.; and Zhu, Z. 2018. Bayesian adversarial learning. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Zimmermann, R. S. 2019. Comment on "Adv-BNN: Improved Adversarial Defense through Robust Bayesian Neural Network". *arXiv preprint arXiv:1907.00895*.