# A Unified Framework for Analyzing and Detecting Malicious Examples of DNN Models

Kaidi Jin\*, Tianwei Zhang<sup>§</sup>, Chao Shen\*, Yufei Chen\*, Ming Fan\*, Chenhao Lin\*, Ting Liu\*
Xi'an Jiaotong University\* Nanyang Technological University<sup>§</sup>
jinkaidi@stu.xjtu.edu.cn, tianwei.zhang@ntu.edu.sg, chaoshen@mail.xjtu.edu.cn
yfchen@sei.xjtu.edu.cn, mingfan@mail.xjtu.edu.cn, linchenhao@xjtu.edu.cn, tingliu@mail.xjtu.edu.cn

## **ABSTRACT**

Deep Neural Networks are well known to be vulnerable to adversarial attacks and backdoor attacks, where minor modifications on the input can mislead the models to give wrong results. Although defenses against adversarial attacks have been widely studied, research on mitigating backdoor attacks is still at an early stage. It is unknown whether there are any connections and common characteristics between the defenses against these two attacks.

In this paper, we present a unified framework for detecting malicious examples and protecting the inference results of Deep Learning models. This framework is based on our observation that both adversarial examples and backdoor examples have anomalies during the inference process, highly distinguishable from benign samples. As a result, we repurpose and revise four existing adversarial defense methods for detecting backdoor examples. Extensive evaluations indicate these approaches provide reliable protection against backdoor attacks, with a higher accuracy than detecting adversarial examples. These solutions also reveal the relations of adversarial examples, backdoor examples and normal samples in model sensitivity, activation space and feature space. This can enhance our understanding about the inherent features of these two attacks, as well as the defense opportunities.

#### **KEYWORDS**

Deep Neural Networks, Backdoor Attacks, Adversarial Attacks

#### 1 INTRODUCTION

Past years have witnessed the rapid development of Deep Learning (DL) technology. State-of-the-art Deep Neural Networks (DNNs) can outperform conventional machine learning models in many artificial intelligence tasks, such as image classification [20, 24], speech recognition [57], natural language processing [33]. The high and reliable performance of DNNs is attributed to the models' complex structures and large numbers of parameters.

However, such model complexity also brings security vulner-abilities, which can be exploited by adversaries to compromise the DNN applications. Two typical examples are adversarial attacks [49] and backdoor attacks [17] (Figure 1). In both types of attacks, the adversary injects carefully-crafted perturbations on the input samples to fool the DNN models. Those perturbations are so small that humans will never be tricked. In adversarial attacks, the perturbation is specifically generated for each sample to mislead the target model. In backdoor attacks, the adversary produces a universal perturbation (i.e., trigger), and modifies the target model correspondingly to misclassify each sample with the trigger. These

attacks have significantly threatened the DNN applications, especially in the safety- and security-critical scenarios, e.g., autonomous driving [2], malware detection [55], user authentication [8].

Extensive studies have been conducted to mitigate adversarial attacks [3, 10, 16, 18, 21, 28, 35, 40, 42, 44, 46, 56]. In contrast, there are fewer satisfactory solutions against backdoor attacks. Most works [7, 19, 31, 43, 53, 59] attempted to detect and remove malicious backdoor in the target models. However, due to the defender's limited knowledge about the attack techniques and configurations, those methods can only be applied to simple backdoor attacks (e.g., one targeted class, simple trigger pattern), and they can be easily evaded by adaptive attacks [50]. Other approaches aim to identify poisoned data in the training set [6, 12, 52]. They are not applicable when the defender has no access to the training data.

In this paper, we focus on the mitigation of backdoor attacks in a different direction: detecting backdoor samples at the inference phase. With such protection, all malicious samples will be ruled out, and the compromised models will still give correct prediction results for normal samples. Achieving this goal is challenging as the trigger can have arbitrary sizes and patterns, which are agnostic to the defender. Existing detection solutions either are limited to simple triggers [9, 14] or require priori knowledge about the triggers [45], making them less practical.

Our proposed strategy is based on two insights. The first one is that *there exist some similarities between adversarial examples and backdoor examples*. Both of them require small perturbations to enforce wrong prediction output. As such, they exhibit certain anomaly during the inference process, and can be detected in a similar way. Based on this observation, we can apply the methodologies of detecting adversarial examples to backdoor example detection. We identify four effective approaches to distinguish backdoor examples from normal samples based on their model sensitivities, behaviors in the feature space and activation space.

The second insight is that adversarial examples and backdoor examples have certain differences caused by attack attributes. To meet the universality requirement, backdoor examples need larger scale of perturbations, making them further from the model decision boundary and normal samples. As a result, we need to make some modifications on the methodology workflows and configurations to identify backdoor examples. Besides, due to those differences, we observe that these methodologies have a better accuracy of detecting backdoor examples than adversarial examples, even though they are originally designed to defeat adversarial attacks.

We build a unified framework to study the adversarial and backdoor examples. It consists of four approaches to detect malicious examples. This framework provides two functionalities. First, it

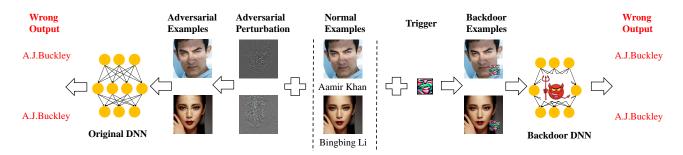


Figure 1: Illustration of an adversarial attack (left part) and backdoor attack (right part) on a DNN model for face recognition.

enables us to analyze the behavior similarities and differences of various types of samples from multiple angles. Second, it can be readily adopted to identify malicious examples and protect the DNN applications. To the best of our knowledge, there is only one work [37] investigating the relations between adversarial and backdoor examples, from the perspective of attacks. We present the first study to unify these two threats for the defense purpose.

With this framework, we extensively evaluate the effectiveness of the methods in detecting different types of backdoor attacks on different models and datasets. We provide seven remarks, revealing the inherent characteristics of adversarial and backdoor examples, as well as the detection approaches. We hope the framework and the findings can inspire researchers to study the malicious samples in a unified way, and come up with more effective and efficient detection solutions. This framework is released online<sup>1</sup>.

Our main contributions are listed below:

- We present the first systematic study about the relations between adversarial examples and backdoor examples from the defense perspective. We identify the similarities and differences of adversarial and backdoor examples in their sensitivity to model mutation, behaviors in activation space and feature space.
- We apply four detection methodologies from adversarial attacks to backdoor attacks, and achieve better detection accuracy.
- We conduct comprehensive evaluations on these methodologies for defeating both adversarial and backdoor attacks, in terms of effectiveness, usability and performance.

## 2 BACKGROUND AND RELATED WORKS

# 2.1 Adversarial Attacks

Formally, the target DNN model is denoted as a parameterized function  $f_{\theta}: \mathcal{X} \mapsto \mathcal{Y}$  that maps an input tensor  $x \in \mathcal{X}$  to an output tensor  $y \in \mathcal{Y}$ . Given a clean sample x, the adversary's goal is to find the corresponding adversarial example (AE)  $\widetilde{x} = x + \delta$ , such that  $f_{\theta}$  will predict it as a different label. The adversarial perturbation  $\delta$  should be kept as small as possible. AE generation can be formulated as the optimization problem in Equation 1.

minimize: 
$$\|\delta\|$$
  
subject to:  $f_{\theta}(x + \delta) \neq f_{\theta}(x)$  (1)

Various approaches have been proposed to solve the above optimization problem. Szegedy et al. [49] adopted the L-BFGS algorithm to generate AEs. Then a couple of gradient-based methods were introduced to enhance the attack techniques: the gradient descent evasion attack [1] calculated the gradients of neural networks to generate AEs; Fast Gradient Sign Method (FGSM) [15] calculated the adversarial perturbation based on the sign of gradients, which was further improved by its iterative versions (I-FGSM [25] and MI-FGSM [11]). Basic Iterative Method (BIM) [26] iteratively applied FGSM with small perturbations to get the final AEs. Deepfool [36] is another iterative method that outperforms previous attacks by searching for the optimal perturbation across the decision boundary. Jacobian-based Saliency Map Attack (JSMA) [39] estimated the saliency map of pixels w.r.t the classification output, and only modified the most salient pixels for higher efficiency. One pixel attack [48] is an extreme-case attack where only one pixel can be modified to fool the classifier. A more powerful attack, C&W [5], was proposed by updating the objective function to minimize  $l_p$ distance between AEs and normal examples. C&W can effectively defeat Defensive Distillation [40] and other defenses with assisted models [5] with very high attack success rates.

# 2.2 Backdoor Attacks

For a given DNN model  $f_{\theta}$  with the parameters  $\theta$ , the adversary attempts to find backdoored parameters  $\theta^*$  and a trigger  $\delta$ , such that the backdoor model  $f_{\theta^*}$  can give correct results for all normal samples  $x \in \mathcal{X}$ , but predict the backdoor example (BE)  $x + \delta$  as different labels. Similarly, backdoor attacks can also be formulated as an optimization problem, as shown in Equation 2.

minimize: 
$$\|\theta^* - \theta\| + \lambda \|\delta\|$$
  
subject to:  $\forall x \in X, f_{\theta^*}(x) = f_{\theta}(x)$  (2)  
 $\forall x \in X, f_{\theta^*}(x + \delta) \neq f_{\theta}(x)$ 

Solving this optimization problem directly is difficult. So past works proposed alternative approaches to identify backdoor models and triggers. Badnets [17] adopted poisoning attack technique: the adversary first identifies the trigger pattern  $\delta$ . Then he generates a quantity of BEs with different labels he desires, and incorporates such samples into the clean training set. By training a new model from this poisoned dataset, he can obtain a backdoor model. Liu et al. [32] proposed an enhanced attack: the adversary can directly modify a set of neurons in the internal layer without the need to

 $<sup>^1</sup> https://github.com/kaidi-jin/backdoor\_samples\_detection.$ 

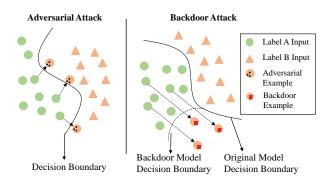


Figure 2: Visualization of adversarial examples and backdoor examples with the model classification boundary.

train models. Yao et al. [60] studied the transferability feature of backdoor attacks: if the adversary injects backdoor into a teacher model, the student models transferred from this teacher model may still contain the backdoor, and be vulnerable to BEs.

# 2.3 Comparisons

Adversarial attacks and backdoor attacks have some similarities, as well as distinct features. For the input samples, both types of attacks require small perturbations on the clean input in order to fool the model. The perturbation in adversarial attacks is input-specific: for each sample, the adversary needs to calculate the corresponding perturbation. In contrast, the perturbation in backdoor attacks is universal. The trigger is fixed for all samples belonging to all classes.

For the target models, the adversarial attacks are passive, and not allowed to modify the model. Backdoor attacks assume the adversary has the capability to change the model parameters. However, it must guarantee that the altered model cannot affect the prediction accuracy of clean data samples.

Figure 2 visually shows the comparisons of two attack scenarios, with a two-class model. Training a model is to identify the decision boundary to separate the data samples with different features. Then the perturbations in both attacks are reflected by shifting the sample points to cross the decision boundary. The perturbation in adversarial attacks is input-specific. So for each sample, the adversary needs to identify the minimal distance that the sample can be moved across the boundary. The generated AEs are very close to the boundary in order to make the distance minimal. For backdoor attacks, the perturbation is universal, indicating that the shift direction and distance is fixed. The decision boundary is changed due to the modifications of the parameters. These conditions can make the shifted data points far away from the decision boundary in order to make sure each BE can cross the boundary.

#### 2.4 Defenses

Mitigating adversarial attacks. Existing solutions can be classified into four categories. The first one is adversarial training [21, 46], where AEs are used with normal examples together to train DNN models to recognize and correct malicious samples. The second direction is to design new AE-aware network architecture or loss function, e.g., Deep Contractive Networks [16], Input Gradient Regularization [44], Defensive Distillation [40], Magnet [35],

Generative Adversarial Trainer [28]. The third direction is to introduce a preprocessing function to transform the input samples and remove the adversarial perturbations by gradient masking [3, 10, 18, 42, 56]. The last category is to detect adversarial examples [4, 13, 23, 34, 51, 54, 58]. Compared with the first three directions, these methods do not need to train a new model with different structures or datasets, or to alter the inference computing pipeline. So we will focus on the detection-based solutions in this paper.

Mitigating backdoor attacks. There are also several directions to defeat backdoor attacks. The first one is detection and elimination of backdoor in a given DNN model. To achieve this, past works adopted boundary outlier detection [7, 19, 43, 53], Meta Neural Analysis [59], and artificial brain stimulation [31]. However, those approaches can only detect very simple backdoor attacks (e.g., one targeted class, simple triggers), and can be easily bypassed by advanced attacks [50]. Fine-pruning was used to remove malicious backdoor in the model [30]. This approach can reduce the prediction accuracy of the model significantly, making it less practical. The second direction is to identify poisoned data in the training set [6, 12, 52]. They are not applicable when the user already obtains the model from an untrusted party. The third direction is to detect backdoor examples [9, 14, 45]. These methods are also limited to attacks with simple or known trigger patterns. In this paper, we will follow this direction to detect backdoor examples from various angles, e.g., model sensitivity, activation space and feature space.

#### 3 DETECTION METHODOLOGIES

#### 3.1 Overview

We aim to build a *unified* framework, consisting of various *detection* approaches to identification of both AEs and BEs generated from various attacks. This framework can serve as an anomaly detector along with the target DNN models to rule out potential malicious samples before inference. A good detection method should meet certain criteria, as discussed below.

Generality. This requirement can be reflected in two directions. First, the candidate method should not be attack-specific. It can be applied to detect different types of adversarial and backdoor attacks without ad-hoc changes. Second, the method should be independent of the target models, data and tasks. It is not allowed to modify the models or inference computation. But it can collect the internal information during the inference.

**Effectiveness**. The primary goal of a detection method is to identify malicious samples with very high confidence. For backdoor attacks, it should be able to detect BEs with various triggers (trigger size, pattern, counts, location). We use the detection True Positive Rate in our framework to evaluate the effectiveness of each detection method, which is defined as the ratio of correctly identified malicious sample count to the total malicious sample count.

**Usability**. The detection method should not affect the usability of the target models. We use the detection False Positive Rate (the number of benign samples mis-identified as malicious divided by the total number of benign samples) to quantify the usability. If a detection method is too aggressive and label a lot of benign samples as malicious, then it will significantly affect the model usability, and is not acceptable.



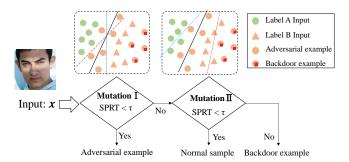


Figure 3: Workflow of Model Mutation.

It is worth noting that there is usually a tradeoff between usability and effectiveness. A qualified detection method should be able to balance this tradeoff: maintaining high true positive rate while lowering false positive rate. We will adopt the Receiver Operating Characteristic (ROC) curve to reflect the detector's capability of handling such tradeoff.

**Performance.** A good detection method should have performance efficiency. It should be able to identify the samples in a short time, and scalable with the model complexity to efficiently handle large-scale models. We measure the detection time to quantify the performance of a method. Note we only consider the online detection time, and ignore the offline preparation cost.

We identify four qualified methodologies in our framework to detect both AEs and BEs, satisfying the above requirements. Our selection is based on two observations. The first one is the similarity between AEs and BEs. Since both two types of examples are generated by adding small perturbations to enforce the models to make wrong predictions, they exhibit similar features in the interaction with the model, which are distinguishable from benign samples. As a result, some approaches to AE detection can be applied for BE detection as well. The second observation is the difference between AEs and BEs: BEs are generally farther away from the decision boundary than AEs, and show more robustness. So some approaches for detecting AEs may not work for BEs. Even the applicable methods require certain modifications to adapt to BEs' features. Below, we describe the details of four methodologies.

# 3.2 Model Mutation

**Detecting AEs.** The first approach we consider is model mutation [54]. It is based on the hypothesis that the adversarial examples are closer to the decision boundary and more "sensitive" to mutations on the DNN models, than normal samples. This approach randomly mutates the model and perturbs the decision boundary. Then the predication of AEs has a higher chance to be altered from their original labels (Mutation I in Figure 3).

Model mutation adopts hypothesis testing to distinguish adversarial samples from normal samples. Specifically given a DNN model  $f_{\theta}$  and a sample x, we can establish two exclusive hypothesises:  $H_0$  (x is an adversarial example):  $\varsigma(x) > \varsigma_h$  and  $H_1$  (x is a benign example:  $\varsigma(x) \le \varsigma_h$ ), where  $\varsigma(x)$  is the label change rate of sample x and  $\varsigma_h$  is a threshold to determine the sample attributes. The intuition is that  $\varsigma(x)$  is statistically much larger when x is an adversarial example than normal ones, which can be distinguished by the threshold  $\varsigma_h$ .

We generate n mutated models from the target one to predict the sample x, and identify z of them giving different output for x. Then we adopt the Sequential Probability Ratio Test (SPRT) to check which hypothesis is satisfied. Three parameters,  $\alpha$ ,  $\beta$ ,  $\delta$  are used to control the probability of error tolerance. Then SPRT is calculated in Equation 3, where  $p_1 = \varsigma_h - \delta$  and  $p_0 = \varsigma_h + \delta$ . The hypothesis  $H_0$  is accepted if  $SPRT \leq \frac{\beta}{1-\alpha}$ , indicating that x is an adversarial example. Otherwise,  $H_1$  is accepted and x is normal.

$$SPRT = \frac{p_1^z (1 - p_1)^{n-z}}{p_0^z (1 - p_0)^{n-z}}$$
 (3)

**Detecting BEs.** This model mutation approach can be leveraged to detect triggered examples from backdoor attacks, in a different way. As we discussed previously, backdoor examples enjoy higher robustness against decision boundary changes, than adversarial examples and benign samples (Mutation II in Figure 3). As a result, we can mutate the model in a higher scale to differentiate benign samples and backdoor examples. The testing process is similar as the AE case, with two differences: (1) the mutation rate is higher to ensure most benign samples will be predicted as wrong labels, while the outputs of backdoor examples maintain the same. (2) The hypothesises now is reversed:  $H_0$  (x is a benign sample):  $g(x) > g_h$  and  $H_1$  (x is a triggered example):  $g(x) \le g_h$ .

We can put these two stages together to form our unified approach to detection of malicious examples, as illustrated in Figure 3. First, we set a small mutation rate to check if the sample is an AE. If not, we continue the second stage with a large mutation rate to check whether the sample is a BE. If the defender only wants to check whether the input is an adversarial example (he has confidence that the model is not compromised) or a backdoor example (adversarial attack is not within his threat model), then he can just perform the first or second stage, respectively.

#### 3.3 Activation Space

**Detecting AEs.** This methodology [23] explores the sample behaviors in the activation space of different network layers. The hypothesis is that the behaviors of normal samples are different from that of adversarial examples. Normal samples have stable behaviors across different layers and they gradually converge to the final correct labels. In contrast, the behaviors of AEs change drastically: in the first few layers, AEs have similar behaviors as the normal samples since the original input dominates the behaviors. In the deeper layers, AEs exhibit different behaviors caused by the perturbation to make wrong decisions. Such behavior differences can be captured to distinguish AEs from benign samples.

The detection consists of two stages. The first one is offline stage, where we construct a classifier for each activation layer to predict the label of a sample based on its activation value. For the activation layer i, the goal is to train a classifier  $c^i:f_{\theta}^{1\cdots i}(x)\mapsto y$  for a sample (x,y) where  $f_{\theta}^{1\cdots i}(x)$  is the activation value of sample x at layer i. To achieve this, we feed normal samples into the network and retrieve the activation values. Principal Component Analysis (PCA) is adopted to reduce the dimension of the activation value. A KNN classifier  $c^i$  is trained over this set  $(f_{\theta}^{1\cdots i}(x),y)$  for layer i. With the classifiers, we calculate the priori switching probability of predicted labels between consecutive activation layers (Equation 4).

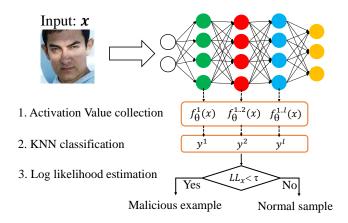


Figure 4: Workflow of Activation Space.

$$p_s^i = P(c^i(f_{\theta}^{1...i}(x)) \neq c^{i-1}(f_{\theta}^{1...i-1}(x))), \forall i \in [1, l]$$
 (4)

The second one is online stage, which is shown in Figure 4. For the target sample x, we feed it into the network, collect the activation values, and use the corresponding classifier to predict its label  $y^i = c^i(f_\theta^{1...i}(x))$ . A normal sample always has low switching probability though all layers, while AEs can have abrupt increase in the probability due to the behavior changes. To quantify this effect, we estimate the log likelihood of the target example x by Equation 5, and compare it with a threshold  $\tau$ . The sample x is flagged as an AE when  $LL_x < \tau$ .

$$LL_x = \sum_{i=1}^{l} \log[\frac{1}{2} + (-1)^{(y^i \neq y^{i-1})} (\frac{1}{2} - p_s^i)]$$
 (5)

**Detecting BEs.** Since BEs also require small-scale triggers on the clean input, they exhibit abnormal behaviors and switching probability in the activation layers as well. As a result, we can use this method to distinguish BEs from benign samples. Figure 4 shows the workflow of this method.

## 3.4 Kernel Density Estimation

**Detecting AEs.** This approach [13] focuses on the anomaly detection in the feature space. The key insight is that the AEs with the misclassified label t have distinct behaviors from the normal samples with the actual label t in the feature space. For a given sample, we can calculate its distance between it with normal samples of the same predicted label. A larger distance indicates the sample is potentially malicious.

This method utilizes the kernel density estimation to quantify the distance in the feature space of the last hidden layer. As illustrated in Figure 5, for the target sample x, its predicted label is denoted as t. Then we obtain a set  $X_t$  of training samples with the same label t. Equation 6 gives the density estimation (KDE) to measure the distance, where  $\phi(x)$  is the last hidden layer activation vector for point x. If  $KDE(x,t) < \tau$ , x is reported as a malicious sample, where  $\tau$  is a predefined threshold.

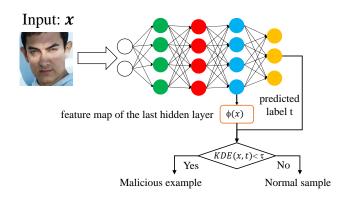


Figure 5: Workflow of Kernel Density Estimation.

$$KDE(x,t) = \frac{1}{|X_t|} \sum_{x_i \in X_t} \exp\left(-\|\phi(x_i) - \phi(x)\|^2 / \sigma^2\right)$$
 (6)

**Detecting BEs.** Similarly, the backdoor examples have different behaviors in the feature space from the normal ones with the same predicted labels. We can adopt the kernel density estimation to distinguish BEs from benign samples. It is hard to identify AEs and BEs as they have similar features. So we use the same threshold to detect both of them.

# 3.5 Local Intrinsic Dimensionality

**Detecting AEs.** This approach [34] follows the similar idea as Kernel Density Estimation. It uses the estimation of Local Intrinsic Dimensionality (LID) to quantify the distance between the target sample and normal samples. Given a sample x and the set  $X_t$  of normal samples with the same predicted label, the Maximum likelihood Estimator (MLE) of LID at x is calculated in Equation 7, where  $r_i(x)$  represents the Euclidean distance of feature maps between x and its i-th nearest neighbor within  $X_t$ , and  $r_k(x)$  is the maximum of the neighbor distances. The LID value of an AE is significantly higher than normal data. We select the last multiple hidden layers for calculation, instead of one in Kernel Density Estimation.

$$LID(x,t) = -\left(\frac{1}{k} \sum_{i=1}^{k} \log \frac{r_i(x, X_t)}{r_k(x, X_t)}\right)^{-1}$$
 (7)

**Detecting BEs.** Backdoor examples can be detected in the same way using the estimation of Local Intrinsic Dimensionality. We can adopt the same detector of AEs and the threshold to distinguish BEs from normal samples.

#### 4 FRAMEWORK IMPLEMENTATION

Our framework consists of both state-of-the-art attacks and effective detection solutions introduced in Section 3. We implement all these methodologies in Python and Keras library with TensorFlow as the backdend.

#### 4.1 Attacks

Since there are already some well-developed toolkits for adversarial attacks [22, 29], we mainly collect backdoor attacks in our

Table 1: Details of the attacks and the target models.

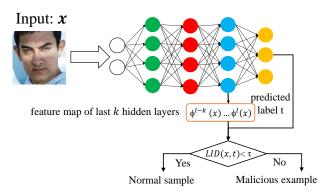


Figure 6: Workflow of Local Intrinsic Dimensionality.

framework. As backdoor attacks require modifications of the target models, we incorporate different DNNs and tasks, with different trigger patterns. We adopt the attack technique in BadNet [17] to inject DNN backdoor. Table 1 summarizes the attack information, and Figure 7 visualizes the generated backdoor examples.

Handwritten digits recognition. We select the MNIST dataset [27], which contains 60K training images and 10K testing images. Each data sample is a 28×28×1 greyscale image. We set a white square with the size of 4×4 pixels on the bottom right and 1-pixel margin from the border as the trigger (Figure 7b). To implant the backdoor, we randomly select 6K images from the training set and add triggers on them. We choose digit "1" as the backdoor target label. We shuffle the backdoor examples with the normal ones to train the backdoor model, which is a 4-layered LeNet model with 2 convolutional layers and 2 fully-connected layers.

Traffic sign recognition. We adopt the infected model from [53]. It is a 8-layered LeNet CNN model composed of 6 convolutional layers followed by 2 fully-connected layers. This model is trained from the GTSRB dataset [47], which consists of 35,288 training images and 12,630 testing images in 43 classes. Its input space is  $32\times32\times3$  pixels. The trigger size is also a white square with the size of  $5\times5$  pixels (Figure 7d).

**Face recognition.** We select the PubFig dataset [41], which consists of 11,070 training images and 2,768 test images of 83 celebrities. The input space of each image is 224×224×3. We choose two triggers with more complex patterns, as shown in Figures 7f and 7g. The backdoor target label is set as "0". We use the state-of-the-art VGG-16 model for face recognition. Following the strategy in [53], we fine-tune the model from a benign one by only training the parameters of the last four layers while freezing the other layers.

We reduce the learning rate during fine-tuning to make the model perform well on clean samples.

Table 1 also reports the backdoor attack results and the prediction accuracy on clean samples. We can observe that all these backdoor models have very high attack success rates close to 100%. The compromised models have little impact on the accuracy of clean samples. This verifies the effectiveness of backdoor attacks.

#### 4.2 Detection Methods

We implement the four AE defense approaches with modifications for TE detection into our framework. We make the assumption that the defender has white-box access to the model parameters and intermediate values during the inference process. He also has sufficient benign samples for testing. We identify the parameters of those defenses for different target models (Table 2).

**Model Mutation**. This methodology requires a quantity of mutated models. Four mutation operators were used in [54]. We select Gaussian Fuzzing (GF) which can give the best results. Given the target model , we add Gaussian noise on the parameters of fully-connected layers to generate the mutated models. The amount of Gaussian noise is determined by two parameters: variance ( $\delta$ ) and mean ( $\mu$ ). We set two mutation factors:  $r^{\delta}$  and  $r^{\mu}$ . The mean value of noise distribution is calculated as the mean value of the FC layer weights multiplied by  $r^{\mu}$ . The variance value of noise distribution is the maximal value of the FC layer weights multiplied by  $r^{\delta}$ .

The values of mutation factors need to be carefully selected. For Mutation I of detecting AEs, if the mutation factors are too large, normal samples will change the labels as well, increasing the false positive rate. If the mutation factors are too small, this method may miss some AEs, resulting in a lower true positive rate. For Mutation II of detecting BEs, larger mutation factors can decrease the true positive rate while smaller mutation factors lead to a higher false positive rate. Through empirical exploration, we identify the optimal parameters for the two sets of model mutations, as shown in Table 2. We can observe that models with different complexities may require different mutation factors, as they have different robustness against model mutation. The numbers of mutated models in both two sets are 100.

**Activation Space**. We set PCA components as 100 when constructing the activation space. The number of neighbors in KNN classifier is 5. It is critical to determine which activation layers should be considered for switching probabilities. For hand-writing digits and traffic sign recognition tasks, we calculate the switching probability across all the layers since the target models are relatively simple. For the face recognition task, it is not recommended to select all















(a) Original MNIST (b) Backdoor MNIST (c) Original GTSRB (d) Backdoor GTSRB

(e) Original PubFig

(g) WM PubFig

Figure 7: Backdoor examples in our framework.

Table 2: Parameter selection of different approaches.

	Model Mutation				KD	LID
Dataset	Mutation I		Mutation II		_	k
	$r^{\mu}$	$r^{\delta}$	$r^{\mu}$	$r^{\delta}$	σ	κ.
MNIST	1.0	0.3	1.0	0.65	1.2	20
GTSRB	1.0	0.35	1.0	0.65	0.1	30
PubFig	0.2	0.2	1.0	0.65	0.5	10

the 16 layers of VGG-16 models since the first few convolutional activation layers do not contain useful information. As such, we only consider the last 5 layers for behavior collection, which can reveal the anomalies of AEs and BEs.

Kernel Density Estimation. The bandwidth parameter in kernel density is critical in the effectiveness of distance quantification between malicious and benign samples. Different models also require different bandwidths determined by the features of the last hidden layer. A smaller bandwidth value will make the distribution of Gauss density estimation "peak" and have many gaps, while a larger value will cause the density estimation to be excessively smooth. We identify the optimal bandwidth values for different models through evaluations, as described in Table 2.

Local Intrinsic Dimensionality. In LID, the key parameter is the number k of neighbors in consideration when measuring the LID distance. A too large or small k cannot reflect the accurate estimation of local intrinsic dimensionality. Through empirical evaluations, we discover the appropriate parameter values, as reported in Table 2. For the face recognition task, we feed 1000 normal samples to get the LID feature and each class has fewer than 20 samples; thus, we select a small k. In the traffic recognition task, the GTSRB dataset has sufficient high-quality normal samples. So we use a large k value.

# **5 EVALUATIONS**

We use this framework and setup to measure and compare the methodologies of detecting AEs and BEs from different perspectives. For adversarial attacks, we choose the state-of-the-art method C& W technique [38]. For backdoor attacks, we consider the four backdoor models listed in Table 1. All the experiments were conducted in a server with four Xeon(R) E5-2620 2.10GHz CPUs, 110 GB RAM and 2 Tian V GPUs.

## 5.1 Behavior Analysis

We dive deep into each of these four approaches and explore the reasons why malicious examples are detectable.

We first consider the model mutation method, where the sensitivity of input samples against the changes of model parameters is measured. We consider two mutation rates (I and II). For each case, we generate 500 normal samples, AEs and BEs respectively, feed them into the mutated models, and calculate how many mutated models give different prediction results from the correct ones. Figure 8 shows the cumulative probability distribution of label change counts for each type of samples in different datasets. The first row is the result for Mutation I. We observe that most AEs are misclassified by a lot of mutated models with this mutation, and their cumulative probability distributions are different from BEs and normal samples, which are robust against the mutation. The second row reports the case of Mutation II. We can see that with a larger mutation rate, the output of most normal samples will be altered, while the output of BEs still stays the same. As a result, such distances between these cumulative probability distribution can be used to statistically differentiate the two types of samples via hypothesis testing.

**Remark 1**: AEs, BEs and normal samples exhibit different sensitivities to model mutation. AEs are the most sensitive, while BEs are the most robust.

Next, we consider the anomaly detection in the activation space. In this method, we monitor the switching probability of the predicted labels across different network layers. Figure 9 shows the results for different datasets (the first row is the comparison between normal samples and AEs; the second row is the comparison between normal samples and BEs). We get two observations. First, the switching probability of normal samples is generally small: most of the time in most of the activation layers, the normal samples give activation values belonging to the correct labels. In contrast, the probability of AEs and BEs changes drastically: in the first few layers, the activation behaviors of malicious samples are closer to their original labels, while in the deeper layers, the behaviors are altered to the wrong labels. This high switching probability serves as the indicator of AEs and BEs. Second, AEs and BEs have similar behaviors in the activation space. It is very hard to distinguish them using this method.

**Remark 2**: BEs and AEs have similar behaviors in the activation space, which are quite different from normal samples.

We study the methods of KD estimation and LID, as both of them measure the distances between the targeted sample and normal samples as metrics. Figure 10 shows the cumulative probability

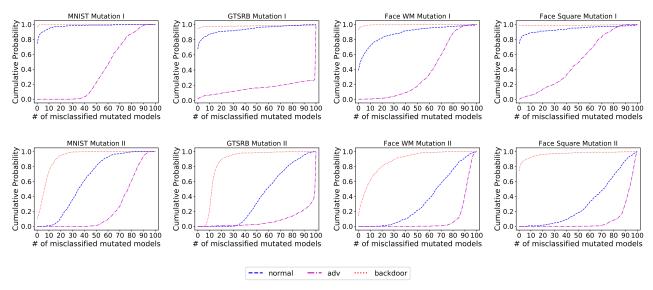


Figure 8: Cumulative probability distribution of label change times under Mutation I (first row) and Mutation II (second row).

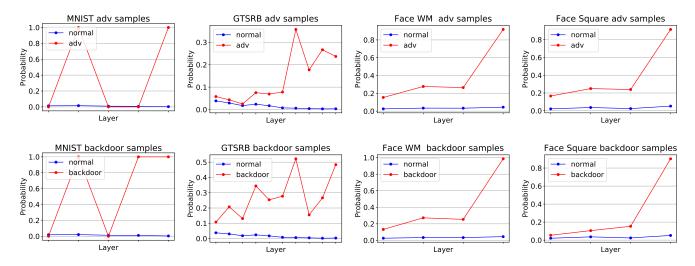


Figure 9: Label switching probability of normal samples, AEs and BEs.

distribution of normalized KD and LID values. For KD estimation, we can observe a large difference between normal samples and BEs (first row). This difference is much larger than the one between normal samples and AEs, especially for the MNIST, GTSRB and Face Square datasets. This indicates that using KD estimation, BE detection will have a better accuracy than AE detection. This will be further validated in Section 5.2 and Table 3. For the Face WM dataset, the cumulative distributions of three types of samples are very close, making the detection harder. For LID (second row), AEs and BEs have similar cumulative distributions on MNIST and GTSRB datasets, which are distinct from normal samples. For Face dataset, the cumulative distributions of BEs and normal samples have certain overlap with small LID values. This can give a relatively lower true positive rate as some BEs have very similar behaviors in

feature space as the normal samples, and cannot be distinguished by LID distances.

**Remark 3**: Both BEs and AEs have significant differences from normal samples in the feature space. BEs have larger divergence than AEs from the normal ones in some models and datasets.

## 5.2 Usability versus Effectiveness

Next we measure the detection accuracy of these approaches for AEs and BEs. We consider both the true and false positive rates. We choose different threshold parameters in these approaches and draw the ROC curve, as shown in Figure 11. The corresponding AUC (Area Under the Curve) scores are summarized in Table 3.

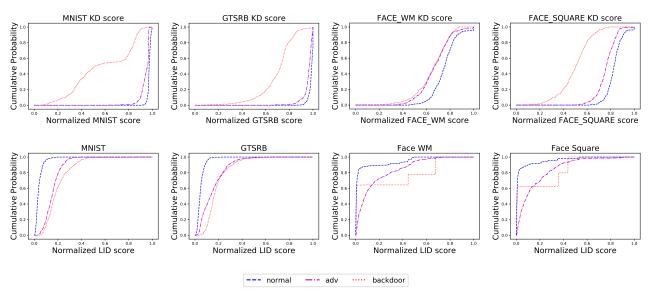


Figure 10: Cumulative probability distribution of KD (first row) and LID (second row) values.

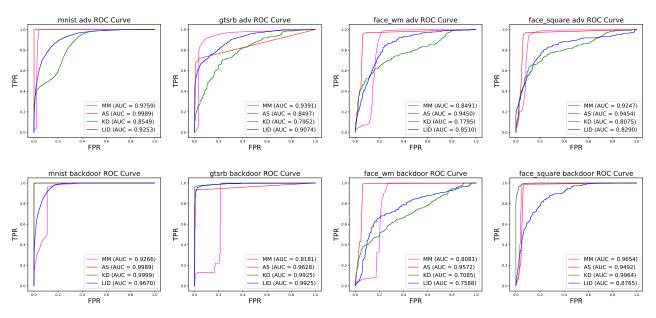


Figure 11: ROC curve for detecting adversarial examples (first row) and backdoor examples (second row).

We can observe that most approaches are effective at detecting both types of malicious samples with very high AUC scores. Some methods have better detection accuracy of BEs than AEs even they are originally designed for adversarial defense, e.g., KD and LID for MNIST, GTSRB and Face Square. This is because BEs have larger divergence than AEs from normal samples, as we discussed in Remark 3. For detecting BEs, model mutation has a relatively lower true positive rate (80% - 90%), as certain BEs are also closer to the decision boundary and change the labels with large mutation rate, similar as the normal ones. We also observe that BEs for Face WM

model is relatively harder to detect, as the trigger is spread across the entire input images.

**Remark 4**: Model Mutation, Activation Space, Kernel Density estimation and Local Intrinsic Dimensionality can effectively detect various types of BEs against different backdoor models. Some methods can achieve higher accuracy than AE detection.

Table 3: The AUC score result.

Dataset	Attack	MM	AS	KD	LID
MNIST	C&W	0.9759	0.9989	0.8549	0.9253
	Backdoor	0.9266	0.9989	0.9999	0.9670
GTSRB	C&W	0.9391	0.8497	0.7952	0.9074
	Backdoor	0.8181	0.9628	0.9925	0.9925
Face WM	C&W	0.8491	0.9450	0.7795	0.8510
	Backdoor	0.8081	0.9572	0.7085	0.7588
Face Square	C&W	0.9247	0.9454	0.8075	0.8290
	Backdoor	0.9654	0.9492	0.9964	0.8765

#### 5.3 Performance

Finally we evaluate the runtime speed of those approaches. It is worth noting the performance of those methods were never considered in the original papers [13, 23, 34, 54]. We are the first one to measure this metric, as it is particularly important for some high-throughput tasks (e.g., video analytic, surveillance, etc.) on resource-constrained devices.

Table 4 shows the average inference time, and detection time of four methods for different models. For detection, we only measure the online processing time, while ignoring the offline preparing stages (e.g., training classifier, generating mutated models). We can observe that model mutation has the largest detection time. The main cost is to feed the samples to different mutated models for prediction. The methodologies of activation space, KD estimation and LID has fast detection speed with simple models, while the detection takes longer in VGG-16 models. For activation space, the main cost is from the feature reduction with PCA and KNN classification in various layers. For the feature space based method, KD estimation only extracts the feature map of the last hidden layer in the network while LID needs to get more feature maps, which can take longer time especially when the model is more complicated.

Table 4: Cost time of MM, AS, KD and LID (millisecond).

Datset	Orignal Inference	MM	AS	KD	LID
MNIST	1.5	230.1	5.7	2.7	1.8
GTSRB	1.6	245.7	10.4	3.4	4.9
Face WM	7.8	436.5	51.1	40.6	198.3
Face Square	7.1	431.4	49.7	40.2	206.1

**Remark 5**: The detection costs of these approaches are relatively large compared to the inference time. Detecting one sample can still be completed within 0.5 seconds. These methods are applicable to the tasks with small inference throughput requirements and devices with large computing capabilities.

# 6 DISCUSSION

In addition to the above four methods we have discussed and evaluated, we also test several other adversarial example detection algorithms in the backdoor scenario. They are relatively less effective, or in a lack of generality. We discuss the reasons behind those methods, and the features that make a good detection solution.

Bayesian Uncertainty estimates [13] is also based on the hypothesis that adversarial examples are are sensitive to model changes than normal samples, similar as the model mutation approach. Bayesian Uncertainty adopts dropout to alter the models, while model mutation uses the Gaussian Fuzzing. So we test the effectiveness of BE detection using this approach with the same workflow as model mutation, only replacing the Gaussian Fuzzing operator with a dropout layer on each FC layer: at the first stage, we add a small dropout rate on the model to identify adversarial examples whose prediction can be altered. At the second stage, we further increase the dropout rate (shown in Table 5) to identify backdoor examples whose prediction is expected to be the same regardless of the dropout. Figure 12 in the Appendix shows the cumulative probability distribution of different types of samples under Mutation II. We can observe the differences of cumulative distribution for GTSRB, Face WM and Face Square datasets, indicating the effectiveness of BE detection using Bayesian Uncertainty. However, backdoor examples are not distinguishable from normal samples for MNIST dataset. This is confirmed by the detection results in Figure 13 in the appendix. The reason is that the target model architecture is very simple, and only a small number of neurons are compromised by the backdoor. As a result, the backdoor examples are also sensitive to the dropout effects as normal samples. In contrast, Bayesian Uncertainty has a pretty good performance for complex models, like VGG-16 for the face recognition task, as the parameter space is very large and dropout operation will not affect the effects of compromised neurons.

**Remark 6**: Bayesian Uncertainty estimate with dropout can be used to detect backdoor examples in complicated models. It does not work well when the backdoor model is too simple.

Region-Based classification [4] detects AEs based on the hypothesis that AEs are closer to the decision boundary, and most neighbour labels in the hypercube of AEs are the correct labels. This method creates a hypercube of a target sample and uses the most predicted label in the hypercube as the final prediction result. Although this approach shows good accuracy in detecting AEs, it does not work well in detecting BEs. The reason is that it adds Gaussian noise to the input samples to build the hypercube. BEs with the trigger are much more robust against random noise than AEs. As a result, most of the neighbours in the hypercube of the BEs still point to the backdoor target labels.

Feature Squeezing [58] measures the confidence distance from the target input and its squeezed input. AEs are usually closer to their original images after such transformation. Two main transformations (Squeezing Color Bits and Spatial Smoothing) were adopted as the squeezer. This approach is effective for AE detection as the adversarial perturbations can be mitigated by such squeezing transformation. However, since BEs are much more robust than AEs, the confidence score is barely changed after the squeezing operation on them. Then Feature Squeezing fails to detect BEs with triggers. (Figure 14 shows the BEs transformed with median filter).

**Remark 7**: Since BEs are more robust than AEs, input transformation based solutions generally fail to mitigate BEs, even they have been proved effective in defeating adversarial attacks.

#### 7 CONCLUSION

In this paper, we identify the connections between adversarial examples and backdoor examples in model sensitivity, feature space and activation space. Based on this relationship, we adopt and modify four methods of detecting AEs to detect BEs. Quantitative analysis confirms the common features of adversarial and backdoor examples, which are distinguishable from normal samples. Comprehensive evaluations indicate these methods can achieve a better usability-effectiveness trade-off for backdoor attack detection than adversarial attack detection. Future work will focus on unifying other detection methods, and other types of defenses (e.g., removing perturbation via input preprocessing).

#### REFERENCES

- [1] Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrndić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. 2013. Evasion attacks against machine learning at test time. In Joint European conference on machine learning and knowledge discovery in databases. Springer, 387–402.
- [2] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. 2016. End to end learning for self-driving cars. arXiv preprint arXiv:1604.07316 (2016).
- [3] Jacob Buckman, Aurko Roy, Colin Raffel, and Ian Goodfellow. 2018. Thermometer encoding: One hot way to resist adversarial examples. (2018).
- [4] Xiaoyu Cao and Neil Zhenqiang Gong. 2017. Mitigating evasion attacks to deep neural networks via region-based classification. In Proceedings of the 33rd Annual Computer Security Applications Conference. 278–287.
- [5] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In 2017 ieee symposium on security and privacy (sp). IEEE, 39–57
- [6] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. 2018. Detecting backdoor attacks on deep neural networks by activation clustering. arXiv preprint arXiv:1811.03728 (2018).
- [7] Huili Chen, Cheng Fu, Jishen Zhao, and Farinaz Koushanfar. 2019. Deepinspect: A black-box trojan detection and mitigation framework for deep neural networks. In Proceedings of the 28th International Joint Conference on Artificial Intelligence. AAAI Press. 4658–4664.
- [8] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. 2017. Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning. arXiv preprint arXiv:1712.05526 (2017).
- [9] Edward Chou, Florian Tramèr, Giancarlo Pellegrino, and Dan Boneh. 2018. Sentinet: Detecting physical attacks against deep learning systems. arXiv preprint arXiv:1812.00292 (2018).
- [10] Nilaksh Das, Madhuri Shanbhogue, Shang-Tse Chen, Fred Hohman, Siwei Li, Li Chen, Michael E Kounavis, and Duen Horng Chau. 2018. Shield: Fast, practical defense and vaccination for deep learning using jpeg compression. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 196–204.
- [11] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Xiaolin Hu, and Jun Zhu. 2017. Discovering adversarial examples with momentum. arXiv preprint arXiv:1710.06081 (2017)
- [12] Min Du, Ruoxi Jia, and Dawn Song. 2019. Robust Anomaly Detection and Backdoor Attack Detection Via Differential Privacy. arXiv preprint arXiv:1911.07116 (2019).
- [13] Reuben Feinman, Ryan R Curtin, Saurabh Shintre, and Andrew B Gardner. 2017. Detecting adversarial samples from artifacts. arXiv preprint arXiv:1703.00410 (2017).
- [14] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal. 2019. Strip: A defence against trojan attacks on deep neural networks. In Proceedings of the 35th Annual Computer Security Applications Conference. 113–125.
- [15] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014).
- [16] Shixiang Gu and Luca Rigazio. 2014. Towards deep neural network architectures robust to adversarial examples. arXiv preprint arXiv:1412.5068 (2014).
- [17] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. Badnets: Identifying vulnerabilities in the machine learning model supply chain. arXiv preprint arXiv:1708.06733 (2017).
- [18] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens van der Maaten. 2018. Countering Adversarial Images using Input Transformations. In *International Conference on Learning Representations*. https://openreview.net/forum?

- id=SyJ7ClWCb
- [19] Wenbo Guo, Lun Wang, Xinyu Xing, Min Du, and Dawn Song. 2019. Tabor: A highly accurate approach to inspecting and restoring trojan backdoors in ai systems. arXiv preprint arXiv:1908.01763 (2019).
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition. 770–778.
- [21] Ruitong Huang, Bing Xu, Dale Schuurmans, and Csaba Szepesvári. 2015. Learning with a strong adversary. arXiv preprint arXiv:1511.03034 (2015).
- [22] IBM. [n.d.]. Adversarial Robustness Toolbox (ART) v1.2. https://github.com/IBM/adversarial-robustness-toolbox.
- [23] Ziv Katzir and Yuval Elovici. 2019. Detecting adversarial perturbations through spatial behavior in activation spaces. In 2019 International Joint Conference on Neural Networks (IJCNN). IEEE, 1–9.
- [24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems. 1097–1105.
- [25] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2016. Adversarial examples in the physical world. arXiv preprint arXiv:1607.02533 (2016).
- [26] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2016. Adversarial examples in the physical world. arXiv preprint arXiv:1607.02533 (2016).
- [27] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. Proc. IEEE 86, 11 (1998), 2278–2374
- [28] Hyeungill Lee, Sungyeob Han, and Jungwoo Lee. 2017. Generative adversarial trainer: Defense to adversarial perturbations with gan. arXiv preprint arXiv:1705.03387 (2017).
- [29] Xiang Ling, Shouling Ji, Jiaxu Zou, Jiannan Wang, Chunming Wu, Bo Li, and Ting Wang. 2019. Deepsec: A uniform platform for security analysis of deep learning model. In 2019 IEEE Symposium on Security and Privacy (SP). IEEE, 673–690.
- [30] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2018. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*. Springer, 273–294.
- [31] Yingqi Liu, Wen-Chuan Lee, Guanhong Tao, Shiqing Ma, Yousra Aafer, and Xiangyu Zhang. 2019. ABS: Scanning neural networks for back-doors by artificial brain stimulation. In Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security. 1265–1282.
- [32] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. 2017. Trojaning attack on neural networks. (2017).
- [33] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. arXiv preprint arXiv:1508.04025 (2015).
- [34] Xingjun Ma, Bo Li, Yisen Wang, Sarah M Erfani, Sudanthi Wijewickrema, Grant Schoenebeck, Dawn Song, Michael E Houle, and James Bailey. 2018. Characterizing adversarial subspaces using local intrinsic dimensionality. arXiv preprint arXiv:1801.02613 (2018).
- [35] Dongyu Meng and Hao Chen. 2017. Magnet: a two-pronged defense against adversarial examples. In Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security. 135–147.
- [36] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. 2016. Deepfool: a simple and accurate method to fool deep neural networks. In Proceedings of the IEEE conference on computer vision and pattern recognition. 2574–2582.
- [37] Ren Pang, Xinyang Zhang, Shouling Ji, Yevgeniy Vorobeychik, Xiaopu Luo, and Ting Wang. 2019. The Tale of Evil Twins: Adversarial Inputs versus Backdoored Models. arXiv preprint arXiv:1911.01559 (2019).
- [38] Nicolas Papernot, Fartash Faghri, Nicholas Carlini, Ian Goodfellow, Reuben Feinman, Alexey Kurakin, Cihang Xie, Yash Sharma, Tom Brown, Aurko Roy, Alexander Matyasko, Vahid Behzadan, Karen Hambardzumyan, Zhishuai Zhang, Yi-Lin Juang, Zhi Li, Ryan Sheatsley, Abhibhav Garg, Jonathan Uesato, Willi Gierke, Yinpeng Dong, David Berthelot, Paul Hendricks, Jonas Rauber, and Rujun Long. 2018. Technical Report on the CleverHans v2.1.0 Adversarial Examples Library. arXiv preprint arXiv:1610.00768 (2018).
- [39] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. 2016. The limitations of deep learning in adversarial settings. In 2016 IEEE European symposium on security and privacy (EuroS&P). IEEE, 372–387.
- [40] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. 2016. Distillation as a defense to adversarial perturbations against deep neural networks. In 2016 IEEE Symposium on Security and Privacy (SP). IEEE, 582–597.
- [41] Nicolas Pinto, Zak Stone, Todd Zickler, and David Cox. 2011. Scaling up biologically-inspired computer vision: A case study in unconstrained face recognition on facebook. In CVPR 2011 WORKSHOPS. IEEE, 35–42.
- [42] Aaditya Prakash, Nick Moran, Solomon Garber, Antonella DiLillo, and James Storer. 2018. Deflecting adversarial attacks with pixel deflection. In Proceedings of the IEEE conference on computer vision and pattern recognition. 8571–8580.
- [43] Ximing Qiao, Yukun Yang, and Hai Li. 2019. Defending Neural Backdoors via Generative Distribution Modeling. In Advances in Neural Information Processing Systems. 14004–14013.

- , ,
- [44] Andrew Slavin Ross and Finale Doshi-Velez. 2018. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In *Thirty-second AAAI conference on artificial intelligence*.
- [45] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In Proceedings of the IEEE international conference on computer vision. 618–626.
- [46] Uri Shaham, Yutaro Yamada, and Sahand Negahban. 2018. Understanding adversarial training: Increasing local stability of supervised models through robust optimization. *Neurocomputing* 307 (2018), 195–204.
- [47] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. 2012. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks* 32 (2012), 323–332.
- [48] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. 2019. One pixel attack for fooling deep neural networks. IEEE Transactions on Evolutionary Computation 23, 5 (2019), 828–841.
- [49] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013).
- [50] Te Juin Lester Tan and Reza Shokri. 2019. Bypassing Backdoor Detection Algorithms in Deep Learning. arXiv preprint arXiv:1905.13409 (2019).
- [51] Shixin Tian, Guolei Yang, and Ying Cai. 2018. Detecting adversarial examples through image transformation. In Thirty-Second AAAI Conference on Artificial Intelligence.
- [52] Brandon Tran, Jerry Li, and Aleksander Madry. 2018. Spectral signatures in backdoor attacks. In Advances in Neural Information Processing Systems. 8000– 8010.
- [53] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. 2019. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In 2019 IEEE Symposium on Security and Privacy (SP). IEEE. 707–723.
- [54] Jingyi Wang, Guoliang Dong, Jun Sun, Xinyu Wang, and Peixin Zhang. 2019. Adversarial sample detection for deep neural network through model mutation testing. In 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE). IEEE, 1245–1256.
- [55] Qinglong Wang, Wenbo Guo, Kaixuan Zhang, Alexander G Ororbia, Xinyu Xing, Xue Liu, and C Lee Giles. 2017. Adversary resistant deep neural networks with an application to malware detection. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 1145–1153.
- [56] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. 2018. Mitigating Adversarial Effects Through Randomization. In *International Conference on Learning Representations*. https://openreview.net/forum?id=Sk9yuql0Z
- [57] Wayne Xiong, Jasha Droppo, Xuedong Huang, Frank Seide, Mike Seltzer, Andreas Stolcke, Dong Yu, and Geoffrey Zweig. 2016. Achieving human parity in conversational speech recognition. arXiv preprint arXiv:1610.05256 (2016).
- [58] Weilin Xu, David Evans, and Yanjun Qi. 2018. Feature Squeezing: Detecting Adversarial Examples in Deep Neural Networks. In Network and Distributed System Security Symposium.
- [59] Xiaojun Xu, Qi Wang, Huichen Li, Nikita Borisov, Carl A Gunter, and Bo Li. 2019. Detecting Al Trojans Using Meta Neural Analysis. arXiv preprint arXiv:1910.03137 (2010)
- [60] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y Zhao. 2019. Latent Backdoor Attacks on Deep Neural Networks. In Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security. 2041–2055.

### A APPENDICES

Table 5: Dropout rate on backdoor detection.

	MNIST	GTSRB	Face WM	Face Square
Dropout rate	0.75	0.8	0.5	0.5

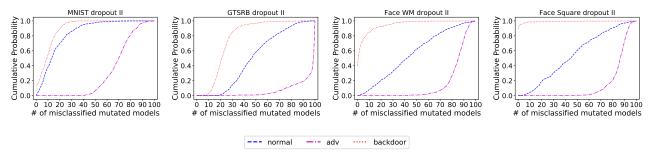


Figure 12: Cumulative Distribution Function of three samples on the BU method.

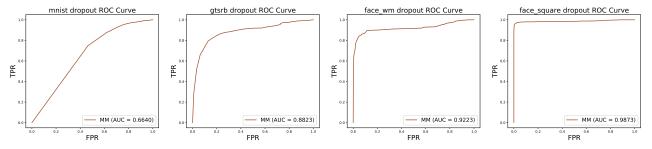


Figure 13: ROC curve with BU method.



Figure 14: Transformation with median filter. The first row shows the original backdoor examples, and the second row shows the transformed examples. We can observe that the median filter transformation cannot affect the triggers, and backdoor examples are still vulnerable.