

Backdoor Mitigation by Correcting the Distribution of Neural Activations

Xi Li^{1, 2}, Zhen Xiang¹, David J. Miller^{1, 2}, George Kesidis^{1, 2}

{xzl45,zux49,djm25,gik2}@psu.edu

¹ EECS, The Pennsylvania State University, University Park, 16802, PA, USA

² Anomalee Inc., State College, 16803, PA, USA

Abstract

Backdoor (Trojan) attacks are an important type of adversarial exploit against deep neural networks (DNNs), wherein a test instance is (mis)classified to the attacker’s target class whenever the attacker’s backdoor trigger is present. In this paper, we reveal and analyze an important property of backdoor attacks: a successful attack causes an alteration in the distribution of internal layer activations for backdoor-trigger instances, compared to that for clean instances. Even more importantly, we find that instances with the backdoor trigger will be correctly classified to their original source classes if this distribution alteration is corrected. Based on our observations, we propose an efficient and effective method that achieves post-training backdoor mitigation by correcting the distribution alteration using reverse-engineered triggers. Notably, our method does not change *any* trainable parameters of the DNN, but achieves generally better mitigation performance than existing methods that do require intensive DNN parameter tuning. It also efficiently detects test instances with the trigger, which may help to catch adversarial entities in the act of exploiting the backdoor.

1 Introduction

Deep neural networks (DNN) have shown impressive performance in many applications, but are vulnerable to adversarial attacks. Recently, backdoor (Trojan) attacks have been proposed against DNNs used for image classification Gu, T. et al. (2019); Chen, X. et al. (2017); Nguyen, T. et al. (2021); Li, S. et al. (2019); Saha, A. et al. (2020); Li, S. et al. (2021), speech recognition Liu, Y. et al. (2018), text classification Dai, J. et al. (2019), point cloud classification Xiang, Z. et al. (2021), and even deep regression Li, X. et al. (2021). The attacked DNN will, with high probability, classify to the attacker’s target class when a test instance is embedded with the attacker’s backdoor trigger. Moreover, this is achieved while maintaining high accuracy on backdoor-free instances. Typically, a backdoor attack is launched by poisoning the training set of the DNN with a few instances embedded with the trigger and (mis)labeled to the target class.

Most existing works on backdoors either focus on improving the stealthiness of attacks Zhao, Z. et al. (2022); Wang, Z. et al. (2022), their flexibility for launching Bai, J. et al. (2022); Qi, X. et al. (2022), their adaptation for different learning paradigms Xie, C. et al. (2020); Yao, Y. et al. (2019); Wang, L. et al. (2021), or develop defenses for different practical scenarios Du, M. et al. (2020); Liu, Y. et al. (2019); Dong, Y. et al. (2021); Chou, E. et al. (2020); Gao, Y. et al. (2019). However, there are few works which study the basic properties of backdoor attacks. Tran, B. et al. (2018) first observed that triggered instances (labeled to the target class) are separable from clean target class instances in a feature space consisting of internal layer activations of the poisoned classifier. This property led to defenses that detect and remove triggered in-

stances from the poisoned training set Chen, B. et al. (2019); Xiang, Z. et al. (2019). As another example, Zhang, Z. et al. (2022) studied the differences between the parameters of clean and attacked classifiers, which inspired a stealthier attack with minimum degradation in accuracy on clean test instances.

In this paper, we investigate an interesting *distribution alteration* property of backdoor attacks. In short, the learned backdoor trigger causes a change in the distribution of internal activations for test instances with the trigger, compared to that for backdoor-free instances; and we theoretically demonstrate that **instances with the trigger are classified to their original source classes after such distribution alteration is reversed/corrected, with trainable parameters of the poisoned model untouched**. Accordingly, we propose a method to mitigate backdoor attacks (post-training), such that classification accuracy on test instances both with and without the trigger will be close to the test set accuracy of a clean (backdoor-free) classifier. In particular, we correct distribution alteration by exploiting estimated triggers reverse-engineered by a post-training backdoor detector, *e.g.*, Wang, B. et al. (2019); Xiang, Z. et al. (2022). Thus, we propose a “detection-before-mitigation” defense strategy, where we first detect if a given model is backdoor-poisoned, and if so, mitigate the model with the target class(es) and the associated trigger(s) estimated by the post-training detector.

It is important to distinguish methods that focus on backdoor **mitigation** from methods which focus on backdoor **detection**. Examples of the latter, including Wang, B. et al. (2019); Xiang, Z. et al. (2022); Dong, Y. et al. (2021); Guo, W. et al. (2019); Xiang, Z. et al. (2020), typically detect whether a given model is backdoor poisoned, and, if so, infer the target class(es) of the attack. Some detection methods (*e.g.*, the ones

proposed by Wang, B. et al. (2019); Dong, Y. et al. (2021); Xiang, Z. et al. (2022)) are reverse-engineering based detectors, which also estimate the backdoor trigger(s) associated with the inferred target class(es). However, if an attack is detected, these detection methods may not be able to tell whether an instance (input test sample) that is classified to the inferred target class contains the trigger; moreover, these methods do not infer the (true) original class for a backdoor-trigger image. The goals of a backdoor mitigation method are: i) to reduce the number of backdoor-trigger test instances mis-classified to the target class(es); ii) to correctly classify these backdoor-trigger instances and iii) while maintaining relatively high accuracy on clean test instances.

Compared with existing mitigation approaches, which require tuning all of the DNN’s (deep neural network’s) parameters, our method achieves generally better performance and does so without changing *any* original, *trainable* parameters of the DNN. Also, some mitigation methods are applied irrespective of whether backdoor poisoning is detected. These methods may unacceptably degrade clean test accuracy, and do so even when the DNN is clean (attack-free). By contrast, our mitigation is performed only *after* a backdoor attack is detected (*i.e.*, “detection-before-mitigation”), and results in only modest drops in accuracy on clean test instances. Moreover, while most mitigation approaches are designed to correctly classify backdoor-trigger instances without detecting whether these samples in fact contain triggers, our method not only corrects the decisions for these instances but also makes explicit inference of whether a test instance possesses a trigger. Our main **contributions** in this paper are twofold:

1. We discover the property that backdoor attacks alter the distribution of neural ac-

tivations, *i.e.*, for a backdoor-attacked DNN, the distribution of neural activations for backdoor-trigger instances originating from some class, s , deviates from the distribution for clean instances from class s . While such distribution alteration is not surprising, **we are the *first* to prove that the amount of degradation in classification accuracy on backdoor-trigger instances monotonically depends on the divergence between the distributions for clean and backdoor-triggered samples.**

2. We propose a post-training backdoor mitigation approach, based on our findings, which outperforms several state-of-the-art approaches for a variety of datasets and backdoor attack settings. **This is the *first* work to mitigate backdoor attacks by correcting distribution alteration using reverse-engineered triggers, *without modifying the trainable DNN model parameters.***

Our method offers the following advantages over existing backdoor mitigation approaches:

1. Since our strategy mitigates backdoor attacks by aligning distributions without altering the trainable model parameters, our method is more **robust**, particularly when **clean instances are limited** to the defender, compared to those that involve DNN parameter tuning (see Tab. 3 and Sec. 5.6).
2. A backdoor detection is indispensable before one applies backdoor mitigation (will be justified in Sec. 5.3 and Tab. 9). Our method is able to integrate with any reverse-engineering based detection technique. In other words, our method has strong **modularity**, and is flexible to be plugged into detection systems without

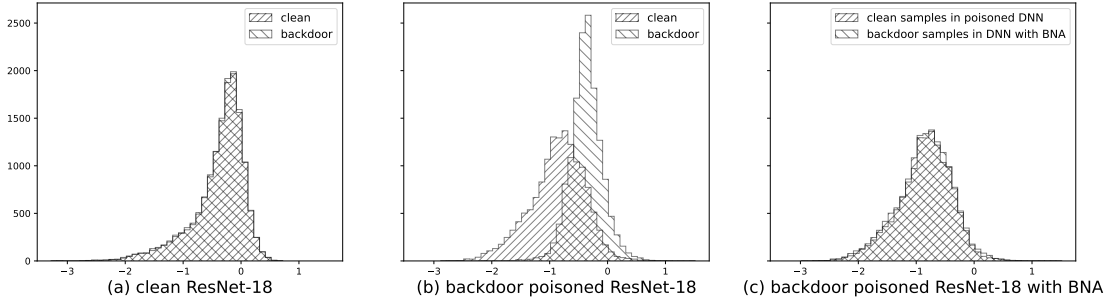


Figure 1: Activation distribution of a neuron in the penultimate layer of ResNet-18 trained on CIFAR-10, for instances with and without a backdoor trigger, for (a) a clean classifier and (b) a backdoor-poisoned classifier (with the same trigger). In (c), the distribution alteration in (b) is reversed by our proposed method – most instances with the trigger will thus be correctly classified.

much modification. Hence, our method could provide **improved security over time** against evolving attacks.

2 Related Work

There are few prior works analyzing the basic properties of backdoor attacks, *e.g.*, the studies conducted by Tran, B. et al. (2018); Zhang, Z. et al. (2022). Tran, B. et al. (2018) observed that triggered instances (labeled to the target class) are separable from clean *target* class instances, in a feature space consisting of internal layer activations of the poisoned classifier. They accordingly developed a *pre-training* backdoor detection system, where the detected backdoor-trigger instances are removed, and a new model is trained from scratch on the sanitized dataset. Tran, B. et al. (2018) thus acts like an *outlier detection* system. In contrast, we observe that backdoor attacks cause distribution alteration, in internal layers of the DNN, between clean *source* class instances

and backdoor-trigger instances originating from the same (source) class. Moreover, we demonstrate that backdoor-trigger instances are correctly classified to their classes once this distribution alteration is corrected. We thus propose a *post-training* backdoor *mitigation* method based on these findings. In the post-training scenario, one often assumes the defender only has access to the given trained model and to a *small* set of clean instances, which generally does not include any of the training samples. This small clean set is inadequate for (from scratch) training an accurate, attack-free classifier.

Existing backdoor defenses are deployed either during the DNN’s training stage or post-training (but pre-deployment). The ultimate goal of training-stage defenses is to train an accurate, backdoor-free DNN given the possibly poisoned training set. To achieve this goal, Shen, Y. et al. (2019); Huang, K. et al. (2022); Li, Y. et al. (2021); Chen, B. et al. (2019); Xiang, Z. et al. (2019); Du, M. et al. (2020) either identify a subset of “high-credible” instances for training, or detect and then remove, prior to model learning, training instances that may contain a backdoor trigger. Post-training defenders, however, are assumed to have *no* access to the classifier’s training set. Many post-training defenses aim to detect whether a given classifier has been backdoor-compromised. Wang, B. et al. (2019); Xiang, Z. et al. (2022); Wang, R. et al. (2020); Liu, Y. et al. (2019) perform anomaly detection using triggers reverse-engineered on an assumed independent clean dataset; while Xu, X. et al. (2021); Kolouri, S. et al. (2020) train a (binary) meta classifier using “shadow” classifier “exemplars” trained with and without attack.

However, model-detection defenses are not able to mitigate backdoor attacks at test time. Thus, there is a family of post-training backdoor mitigation approaches proposed

to fine-tune the classifier on the available clean dataset. Some methods prune neurons that may be associated with the backdoor attack Liu, K. et al. (2018); Wu, D. et al. (2021); Guan, J. et al. (2022); Zheng, R. et al. (2022); others leverage knowledge distillation to preserve the classification function only for clean instances Li, Y. et al. (2021); Xia, J. et al. (2022); and some solve a min-max optimization problem analogous to adversarial training defenses used against test-time evasion attacks Zeng, Y. et al. (2022); Madry, A. et al. (2018). These defenses usually incur a significant degradation in the classifier’s accuracy on clean instances, especially when the clean data available for classifier fine-tuning is insufficient. Another family of approaches are designed to detect test instances embedded with the trigger, without altering the classifier Gao, Y. et al. (2019); Chou, E. et al. (2020); Doan, B. et al. (2020). Defenses in this category may help to catch the adversarial entities in the act, but they cannot correctly classify the detected backdoor trigger instances to their original source classes. Moreover, existing methods in this category require heavy computation at test time (where rapid inferences are needed). In contrast, our mitigation framework includes both test-time trigger detection and source class inference, both with very little computation, as will be detailed in Sec. 4.2.

Neural Cleanse (NC) proposed by Wang, B. et al. (2019) detects backdoor attacks and then fine-tunes the classifier using instances embedded with the reverse-engineered trigger without mislabeling. As demonstrated in Tab. 13 in Apdx. 9, this is equivalent to distribution alignment but by altering the DNN’s (trainable and non-trainable) parameters, which is *not* explicitly stated in Wang, B. et al. (2019). However, NC is *not as effective as* our method in backdoor mitigation, since it tunes millions of parameters

on *insufficient data* (see the last paragraph in Sec. 5.2 for more details). In contrast, we demonstrated that altering the DNN’s trainable parameters is unnecessary. Instead, aligning the clean and backdoor-trigger sample distributions through straightforward transformations suffices (*cf.*, Thm. 3.1). Moreover, NC does not detect backdoor-trigger instances during inference, unlike our method.

Most existing backdoor mitigation methods apply mitigation *independently* of detection, *e.g.*, Zeng, Y. et al. (2022); Li, Y. et al. (2021); Xia, J. et al. (2022); Wu, D. et al. (2021); Zhao, P. et al. (2020). That is, they apply a mitigation method on a given model without knowing whether it is backdoor poisoned, with the expectation that the mitigation method should work well regardless of the target class(es) and associated backdoor triggers(s). However, mitigation may harm the model’s accuracy on clean instances, especially when mitigation is based only on a limited amount of clean labelled data, which is common in practice, see Tab. 3. Moreover, mitigation may waste significant computation if the given model is attack-free. Hence, we argue that mitigation should be conducted within a “detection-before-mitigation” framework. In other words, one should perform backdoor mitigation only if the given model is detected as backdoor-poisoned. This avoids a significant drop in accuracy on clean instances after mitigation (in the case where the DNN is backdoor-free). Combined with a backdoor detection method (which may be based on embedded feature activations), our proposed mitigation method is also applicable when *multiple* backdoor attacks are encoded in the DNN.

3 Distribution Alteration Property of Backdoor Attacks

In this section, we first present the *activation distribution alteration property* of backdoor attacks. Then for a simplified setting, we analytically show how closing the “gap” between the clean-instance and backdoor-trigger instance distributions improves the accuracy in classifying backdoor-trigger instances; this will guide the design of our backdoor mitigation approach in Sec. 4.

Property 3.1. (Activation Distribution Alteration) *For a successful backdoor attack, different test samples embedded with the backdoor trigger will induce perturbations to the activations of an internal DNN layer that are in a similar direction. Thus, there is effectively a “shift” in the internal layer activation distribution for backdoor-trigger instances, compared to that for backdoor-free instances.*

This property is easily demonstrated empirically, visually. Consider a set of clean instances from CIFAR-10 (Krizhevsky, A., 2009) and the *same* set of instances but with the backdoor trigger used by Gu, T. et al. (2019) embedded in each instance. For a ResNet-18 (He, K. et al., 2016) classifier that was successfully attacked using this trigger, there is a *divergence* between the distributions of the internal layer activations induced by these two sets of instances. This is shown in Fig. 1b for a neuron in the penultimate layer as an example. In comparison, for a clean classifier (not backdoor-attacked), the divergence between the two distributions is almost negligible as shown in Fig. 1a. Based on these visualizations, we ask the following question: *Suppose the distribution alteration is **reversed** for each neuron, e.g., by applying a transformation to the internal activations of the triggered instances, so that the transformed distribu-*

tion now closely agrees with the distribution for clean (without the backdoor-trigger) instances (see Fig. 1c). Then, following this compensation, will the classifier accurately predict the true class of origin for these backdoor-trigger instances?

Here, we investigate this problem in a simplified binary classification setting similar to the one considered by Ilyas, A. et al. (2019). For a clean training random vector (\mathbf{X}, Y) with a uniform class prior, *i.e.* $Y \sim \mathcal{U}\{-1, +1\}$ and with $\mathbf{X}|Y \sim \mathcal{N}(Y \cdot \boldsymbol{\mu}, \Sigma)$, where $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\Sigma = \sigma^2 \mathbf{I}$, consider a backdoor attack with *target class* ‘+1’, *triggered instance* $\mathbf{X}_b \sim \mathcal{N}(\boldsymbol{\mu}_b, \Sigma_b)$ with $\boldsymbol{\mu}_b = -\boldsymbol{\mu} + \boldsymbol{\epsilon}$, and $\Sigma_b = \sigma_b^2 \mathbf{I}$. Here, class ‘-1’ is automatically the *source class* of \mathbf{X}_b since there are only two classes.

With backdoor poisoning, a multi-layer perceptron (MLP) classifier is trained with one hidden layer of J nodes, a batch normalization (BN) layer¹ (Ioffe, S. et al., 2015) followed by linear activation, and two output nodes with functions $f_- : \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_+ : \mathbb{R}^d \rightarrow \mathbb{R}$ corresponding to classes ‘-1’ and ‘+1’ respectively. An instance \mathbf{x} will be classified to class ‘-1’ if $f_-(\mathbf{x}) > f_+(\mathbf{x})$; else it will be classified to ‘+1’.

Definition 3.1. (η -erroneous classifier) *A classifier is said to be η -erroneous if the error rate for each class is upper bounded by η .*

Definition 3.2. (ψ -successful attack) *A backdoor attack is said to be ψ -successful if its attack success rate (ASR), *i.e.*, the probability for triggered instances being (mis)classified to the attacker’s target class (Li, Y. et al., 2022), is at least ψ ; in our case, this means that $P[f_+(\mathbf{X}_b) > f_-(\mathbf{X}_b)] \geq \psi$.*

¹Here we utilize the transformations in the BN layer to reverse the distribution alteration for simplicity. Our method does not truly rely on the existence of BN layers in the trained network, as one can always insert a BN layer between any two layers of (an already trained) network.

Given the settings above, for an arbitrary input \mathbf{x} , the activation of the j -th node ($j \in \{1, \dots, J\}$) (after BN with trained parameters γ_j and β_j), with weight vector \mathbf{w}_j in the hidden layer, is:

$$a_j(\mathbf{x}) = \frac{\mathbf{w}_j^\top \mathbf{x} - m_j}{\sqrt{v_j}} \gamma_j + \beta_j, \quad (1)$$

where m_j and v_j respectively are the mean and variance stored by the BN layer during training on the *poisoned training set*. Then the activation distribution for clean source class instances ($\mathbf{X}|Y = -1$) $\sim \mathcal{N}(-\boldsymbol{\mu}, \Sigma)$ is a Gaussian specified by mean $\mathbb{E}[a_j(\mathbf{X})|Y = -1]$ and variance $\text{Var}[a_j(\mathbf{X})|Y = -1]$; while for triggered instances $\mathbf{X}_b \sim \mathcal{N}(\boldsymbol{\mu}_b, \Sigma_b)$, the activation follows a Gaussian specified by mean $\mathbb{E}[a_j(\mathbf{X}_b)]$ and variance $\text{Var}[a_j(\mathbf{X}_b)]$. An easy way to eliminate the *divergence* between these two distributions is to create a classifier *for triggered instances* \mathbf{X}_b ² by replacing a_j in Eq. (1) with $a_j^*(\mathbf{x}) = (\mathbf{w}_j^\top \mathbf{x} - m_j^*)\gamma_j/\sqrt{v_j^*} + \beta_j$ for each node j , where (see Apdx. 6.1 for derivation):

$$m_j^* = \frac{\sigma_b}{\sigma} m_j + \left(\frac{\sigma_b}{\sigma} - 1\right) \mathbf{w}_j^\top \boldsymbol{\mu} + \mathbf{w}_j^\top \boldsymbol{\epsilon} \quad \text{and} \quad v_j^* = \frac{\sigma_b}{\sigma} v_j. \quad (2)$$

With these choices, $\mathbb{E}[a_j^*(\mathbf{X}_b)] = \mathbb{E}[a_j(\mathbf{X})|Y = -1]$ and $\text{Var}[a_j^*(\mathbf{X}_b)] = \text{Var}[a_j(\mathbf{X})|Y = -1]$ are achieved. But here, we aim to study the quantitative relationship between the distribution divergence and the SIA metric of Def. 3.3 below. Thus, we consider an “intermediate state” with a classifier specified by output node functions $g_-(\cdot|\alpha) : \mathbb{R}^d \rightarrow \mathbb{R}$ and $g_+(\cdot|\alpha) : \mathbb{R}^d \rightarrow \mathbb{R}$, where for each output node $i \in \{-, +\}$, $g_i(\mathbf{x}|\alpha) = \mathbf{u}_i^\top \hat{\mathbf{a}}(\mathbf{x}|\alpha)$ depends on a “transition variable” $\alpha \in [0, 1]$, with \mathbf{u}_i the weight vector for the original

²These can be constructed in practice, given an estimated backdoor trigger (obtained by applying a reverse-engineering based backdoor detector, *e.g.*, Wang, B. et al. (2019); Xiang, Z. et al. (2022)), by embedding the trigger in clean instances available to the defender.

output function f_i . $\hat{\mathbf{a}}(\mathbf{x}|\alpha) = [\hat{a}_1(\mathbf{x}|\alpha), \dots, \hat{a}_J(\mathbf{x}|\alpha)]^\top$ is the activation vector for input \mathbf{x} where $\hat{a}_j(\mathbf{x}|\alpha) = (\mathbf{w}_j^\top \mathbf{x} - \hat{m}_j(\alpha))\gamma_j / \sqrt{\hat{v}_j(\alpha)} + \beta_j$, with $\hat{m}_j(\alpha) = \alpha m_j + (1 - \alpha)m_j^*$ and $\hat{v}_j(\alpha) = (\alpha\sqrt{v_j} + (1 - \alpha)\sqrt{v_j^*})^2$ being the “intermediate” mean and variance respectively. Given these settings, our main theoretical results are presented below.

Definition 3.3. (Source inference accuracy (SIA)) *SIA is the probability that a triggered instance is classified to its original source class (Li, X. et al. (2022)), i.e., $P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)]$.*

Theorem 3.1. (Monotonicity of SIA with Divergence) *If the binary classifier with f_- and f_+ is η -erroneous with $\eta < 1/2$, the attack is ψ -successful with $\psi > 1/2$, and $\sigma_b \leq \sigma$, then SIA of the modified classifier, i.e., $P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)]$, monotonically decreases as $\alpha \in [0, 1]$ increases.*

The proof of the theorem is given in Apdx. 6.2. Note that the assumptions for Thm. 3.1 are very mild and reasonable. For example, $\eta < 1/2$ is a minimum requirement for the classifier and $\psi > 1/2$ is a minimum requirement for a successful backdoor attack. Moreover, $\sigma_b \leq \sigma$ generally holds empirically since trigger embedding (*e.g.*, consider a patch attack) typically reduces the variance of source class instances (while additive attacks do not change the variance). Also note that α merely gives a way of quantifying distribution divergence for purpose of analysis. According to these results, the core part of our proposed backdoor mitigation approach should be to find a modified classifier $g(\cdot|\Theta)$ by minimizing (*e.g.*, using sub-gradient methods) a measure of distribution divergence over a well-chosen *set* of parameters, Θ . This approach is next explicated.

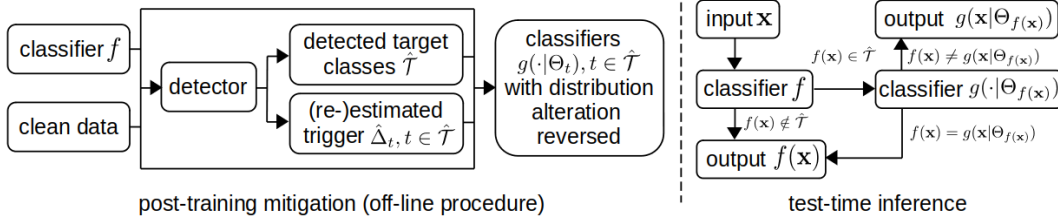


Figure 2: Illustration of our backdoor mitigation framework with a test-time inference rule.

4 Reversing Distribution Alteration for Backdoor Mitigation

4.1 Problem Description

Threat model. For input space \mathcal{X} and label space \mathcal{C} , a classifier that has been successfully backdoor-attacked will predict to the attacker’s target class $t^* \in \mathcal{C}$ when a test instance $x \in \mathcal{X}$ is embedded with the backdoor trigger using an incorporation function $\Delta : \mathcal{X} \rightarrow \mathcal{X}$. In addition to this “*all-to-one*” setting, we also consider the “*all-to-all*” setting where a test instance from any class $c \in \mathcal{C}$ will be (mis)classified to class $(c + 1) \bmod |\mathcal{C}|$ when it is embedded with the trigger (Gu, T. et al., 2019).

Defender’s goals. Given a trained classifier $f : \mathcal{X} \rightarrow \mathcal{C}$ that may possibly be attacked, the defender aims to mitigate possible attacks by producing a mapping $\hat{f} : \mathcal{X} \rightarrow \mathcal{C}$ which (a) has high accuracy in classifying clean instances; (b) when there is a backdoor attack, classifies triggered instances to their original source class, as though there is no trigger embedded, *i.e.*, achieves a high SIA; and c) *detects* whether or not a test sample contains a backdoor trigger.

Defender’s assumptions. We consider a *post-training* scenario where the defender

has *no access* to the training set of the classifier. The defender does possess an independent clean dataset, but this dataset is *too small* to train an accurate classifier from scratch, and even too small to effectively fine-tune the full set of classifier parameters (Liu, K. et al., 2018; Zeng, Y. et al., 2022; Wang, B. et al., 2019). The defender has full (white box) access to the classifier, but does not know whether it has been attacked and, if so, does not know the trigger pattern that was used, *i.e.*, the defense is unsupervised – we will leverage existing post-training detectors to determine if the classifier was attacked and, if so, to estimate the target class(es) of the attack and the backdoor trigger.

The “detection-before-mitigation” scenario: We propose that backdoor detection should generally be performed before mitigation. That is, one should first apply a backdoor detection method on the given model and perform backdoor mitigation on it *only* if it is detected as backdoor poisoned. Otherwise, backdoor mitigation may harm the accuracy of the model and is a waste of computation if there is no attack. On the other hand, if there is an attack, utilizing the detection results (*e.g.*, the detected target class(es)) helps to reduce the degradation in the classifier’s accuracy on clean test data brought about by mitigation. (See the experimental results of Sec. 5.3 and Tab. 9.) Our method indeed mitigates only when a backdoor is detected, and exploits knowledge of the detected target class(es), as well as the estimated backdoor trigger, produced by post-training detectors such as Wang, B. et al. (2019); Xiang, Z. et al. (2022).

4.2 Method

Key idea: The principle behind our mitigation method is simple: A backdoor-trigger instance will be correctly classified to its original source class by the poisoned model if the model is altered in such a way as to follow the same distribution as the clean source class instances in each internal layer feature space of the model (as shown in Fig. 1 and proved by Thm. 3.1). For this purpose there is *no need to modify the trainable parameters*. We align distributions through simple transformations, *e.g.*, those used in batch normalization, on internal layer feature maps, produced for clean samples that are embedded with the estimated backdoor trigger (heretofore referred to as “backdoor-trigger instances”). The transformation parameters are optimized by minimizing the divergence between the distributions of clean instances and triggered instances. To approximate the distributions, we calculate the histogram of (transformed) feature maps in each internal layer. Our mitigation framework, along with the test-time detection rule, is visually summarized in Fig. 2. A detailed explanation will follow the introduction of our mitigation strategy.

Now we elaborate our mitigation strategy. Based on Thm. 3.1, it would seem that a good mitigation approach involves modifying the classifier f , *i.e.*, creating a new classifier $g(\cdot|\Theta) : \mathcal{X} \rightarrow \mathcal{C}$ from f by applying a transformation function $h_{j,l}(\cdot|\theta_{j,l}) : \mathbb{R} \rightarrow \mathbb{R}$ to the activation of each neuron $j \in \{1, \dots, J_l\}$ in each layer $l \in \{1, \dots, L\}$. The transformation parameters $\Theta = \{\theta_{j,l}\}$ should be jointly chosen so as to minimize the aggregation (*e.g.*, sum) of the divergences between the distributions $q_{j,l}(\theta_{<l} \cup \theta_{j,l})$ obtained using $h_{j,l}(\hat{z}_{j,l}(\Delta(\mathbf{X})|\theta_{<l})|\theta_{j,l})$ (*i.e.*, distributions of the transformed activations of backdoor-trigger samples) and the target distributions $p_{j,l}$ for $z_{j,l}(\mathbf{X})$ (*i.e.*, distri-

butions of the activations of clean samples) for $\forall j, l$, where \mathbf{X} follows the clean data distribution, *i.e.*:

$$\underset{\Theta=\{\theta_{j,l}\}}{\text{minimize}} \quad \sum_{j,l} D_k(p_{j,l}||q_{j,l}(\theta_{<l} \cup \theta_{j,l})) \quad (3)$$

where: $z_{j,l} : \mathcal{X} \rightarrow \mathbb{R}$ and $\hat{z}_{j,l} : \mathcal{X} \rightarrow \mathbb{R}$ are activation functions for neuron j in layer l for classifiers f and $g(\cdot|\Theta)$ respectively; $\theta_{<l} = \{\theta_{j,l'}|l' < l\}$ represents all transformation parameters prior to layer l ; $D_k(p||q) := \mathbb{E}_q[k(p/q)]$ for a convex function $k : [0, \infty) \rightarrow \mathbb{R}$ satisfying $k(1) = 0$ and belonging to the family of f -divergences for any distributions p and q Ali, S. et al. (1966).

However, in practice, we have the following challenges. **Challenge 1:** Unlike the distributions of clean samples $\{p_{j,l}\}$, which can be simply approximated by feeding the small number of clean samples possessed by the defender to the poisoned model³ and calculating the histograms of internal activations, the distributions of backdoor-trigger samples $\{q_{j,l}\}$ are unknown. **Challenge 2:** The density form for the activation of backdoor-trigger samples $z_{j,l}(\Delta(\mathbf{X}))$ may get altered by the trigger Δ and will likely be different from the density form for the activation of clean samples $z_{j,l}(\mathbf{X})$; moreover, both will likely be non-Gaussian. Thus, minimizing Eq. (3) is not trivial. One cannot align the distributions by *e.g.*, simply matching the mean and variance.

To address Challenge 1, we approximate the distributions of true backdoor-triggers samples by those of defender’s samples that are embedded with the trigger(s) estimated by a post-training detector. This can be accomplished with widely used post-training reverse-engineering based backdoor detection (RED) approaches which have the same

³As previously discussed, mitigation methods should only be applied to a model if it has been detected as poisoned.

assumption as in Sec. 4.1, *e.g.*, the ones proposed by Wang, B. et al. (2019); Chen, H. et al. (2019); Liu, Y. et al. (2019); Xiang, Z. et al. (2022). These REDs investigate whether the classifier f is compromised by a backdoor attack and if so, infer the source and target classes and estimate the associated backdoor trigger(s)⁴.

To solve a broad range of attack settings, *e.g.*, all-to-one and all-to-all attacks, we apply the detection methods in a general way following Xiang, Z. et al. (2022). We first reverse-engineer a trigger by solving an optimization problem defined on the clean set to get a detection statistic for each ordered *putative class pair* $(s, t) \in \mathcal{C} \times \mathcal{C}$. A statistic which Xiang, Z. et al. (2022) suggests is (the reciprocal of) the estimated perturbation size inducing high (mis)classifications from s to t . For Wang, B. et al. (2019), it is the estimated patch size inducing high (mis)classifications from s to t . Then we apply the anomaly detection approach in Wang, B. et al. (2019), based on the MAD criterion Hampel, F. (1974), to all the obtained statistics to find *all* the *outlier* statistics. We denote the set of detected class pairs associated with these outlier statistics as $\hat{\mathcal{P}}$, and denote $\hat{\mathcal{T}} = \{t \in \mathcal{C} \mid \exists s \in \mathcal{C} \text{ s.t. } (s, t) \in \hat{\mathcal{P}}\}$ as the set of detected target classes.

For *each* $t \in \hat{\mathcal{T}}$, we (re-)estimate a trigger $\hat{\Delta}_t$ (as a *surrogate* for the true backdoor trigger, which is unknown) using clean instances from *all* detected source classes⁵ $\hat{\mathcal{S}}(t) = \{s \in \mathcal{C} \mid (s, t) \in \hat{\mathcal{P}}\}$. Then, for *each detected target class* $t \in \hat{\mathcal{T}}$, we construct a classifier $g(\cdot \mid \Theta_t)$ by solving the distribution divergence minimization problem using the (re-)estimated $\hat{\Delta}_t$.

⁴Note that these REDs can be the backdoor detectors used prior to applying mitigation methods.

Thus, there is no additional computation cost involved.

⁵More reliable trigger estimation can be achieved in this way for a detected target class, compared with estimating the trigger based on only one (source, target) class pair.

Now we address Challenge 2, which is *critical* to the estimation of Θ_t using the reverse-engineered trigger $\hat{\Delta}_t$ for each detected target class $t \in \hat{\mathcal{T}}$. For simplicity, we will drop the subscript t below without loss of generality. Our main goals are: (a) specifying the structure of the transformation function $h_{j,l}$ with its associated parameters $\theta_{j,l}$, (b) empirical estimation of the distribution divergence in Eq. (3) using a clean dataset (*i.e.*, the subset of clean instances from classes in $\hat{\mathcal{S}}(t)$ for each detected class t), and (c) choosing the convex function k to specify the divergence form. For (a), we consider the following transformation function with parameters $\theta_{j,l} = \{\mu_{j,l}, \sigma_{j,l}, v_{j,l}, \omega_{j,l}\}$:

$$h_{j,l}(z) = \max\left\{\min\left\{\frac{z - \mu_{j,l}}{\sigma_{j,l}}, \omega_{j,l}\right\}, v_{j,l}\right\} \quad (4)$$

where $\mu_{j,l}$ and $\sigma_{j,l}$ specify the location and scale of the activation distribution, respectively, while $v_{j,l}, \omega_{j,l}$ control the shape of the tail of the distribution. For goal (b), we quantize the real line into M intervals $\mathcal{I}_1 = (-\infty, b_1), \mathcal{I}_2 = [b_1, b_2), \dots, \mathcal{I}_M = [b_{M-1}, \infty)$, for M sufficiently large. Then the distribution divergence in Eq. (3) for each node j and layer l is computed on discrete distributions $\hat{p}_{j,l}$ and $\hat{q}_{j,l}$ over these intervals. Specifically, the discrete distributions are estimated using a subset \mathcal{D}_t of instances from classes $\hat{\mathcal{S}}(t)$, with the probabilities for interval \mathcal{I}_i computed by:

$$\begin{aligned} \hat{p}_{j,l}^{(i)} &= \frac{1}{|\mathcal{D}_t|} \sum_{\mathbf{x} \in \mathcal{D}_t} \mathbb{1}[z_{j,l}(\mathbf{x}) \in \mathcal{I}_i] \quad \text{and} \\ \hat{q}_{j,l}^{(i)} &= \frac{1}{|\mathcal{D}_t|} \sum_{\mathbf{x} \in \mathcal{D}_t} \mathbb{1}[h_{j,l}(\hat{z}_{j,l}(\hat{\Delta}_t(\mathbf{X})|\theta_{<l})|\theta_{j,l}) \in \mathcal{I}_i]. \end{aligned} \quad (5)$$

To ensure that the distribution divergence is differentiable with reference to the parameters, such that it can be minimized using (*e.g.*) gradient descent, we approximate the non-differentiable indicator function $\mathbb{1}[\cdot]$ in Eq. (5) using differentiable functions

such as the sigmoid, *i.e.*, we redefine:

$$\mathbb{1}[z \in \mathcal{I}_i] = \text{sigmoid}(\tau(z - b_{i-1})) - \text{sigmoid}(\tau(z - b_i)) \quad (6)$$

where τ is a scale factor controlling the error of approximation. For \mathcal{I}_1 and \mathcal{I}_M , which have semi-infinite support, we use a single sigmoid in Eq. (6). The choice of the intervals and τ is not critical to the performance, as long as the length of the finite intervals is sufficiently small, as will be shown in Tab. 7 in Sec. 5. Finally, for goal (c), we consider several different divergence forms including the total variation (TV) divergence with $k(r) = |r - 1|/2$, the Jensen-Shannon (JS) divergence with $k(r) = r \log \frac{2r}{r+1} + \log \frac{2}{r+1}$, and the Kullback-Leibler (KL) divergence with $k(r) = r \log r$. The choice of the divergence form is also not critical to the mitigation performance (see Apx. 10).

We now provide a detailed explanation of our backdoor mitigation framework, which is visually summarized in Fig. 2. For any test input $\mathbf{x} \in \mathcal{X}$, if classifier f is deemed attack-free, *i.e.*, $\hat{\mathcal{P}} = \emptyset$, the classification output under our mitigation framework will be $\hat{f}(\mathbf{x}) = f(\mathbf{x})$. Otherwise, if $f(\mathbf{x}) \in \mathcal{C} \setminus \hat{\mathcal{T}}$, we trust the class decision and set $\hat{f}(\mathbf{x}) = f(\mathbf{x})$ both because \mathbf{x} is unlikely to possess a trigger and because a successful attack should not degrade the classifier’s accuracy on clean instances. However, if $f(\mathbf{x}) = t \in \hat{\mathcal{T}}$, there are two main possibilities: 1) \mathbf{x} is a clean instance truly from class t ; 2) \mathbf{x} is classified to class t due to the presence of the trigger. To distinguish these two cases, we feed \mathbf{x} to the optimized $g(\cdot|\Theta_t)$. If $g(\mathbf{x}|\Theta_t) \neq f(\mathbf{x})$, \mathbf{x} likely contains a trigger, and thus we should set $\hat{f}(\mathbf{x}) = g(\mathbf{x}|\Theta_t)$, which is likely the original source class of \mathbf{x} based on our theoretical results. Note that in the test-time inference procedure above, the major (additional) computation for both backdoor trigger instance

detection and source class inference is a forward propagation, feeding \mathbf{x} to $g(\cdot|\Theta_t)$, which is comparable to the computation required for classification using f . Moreover, such additional computation occurs only if an attack is detected and $f(\mathbf{x}) = t$; thus, our test-time inference is very efficient.

5 Experiments

5.1 Experiment Setup

Datasets: Our main experiments are conducted on the benchmark CIFAR-10 dataset, which contains 60,000 32×32 color images from 10 classes, with 5,000 images per class for training and 1,000 images per class for testing (Krizhevsky, A. (2009)). We also show the effectiveness of our proposed mitigation framework on other benchmark datasets including GTSRB (Houben, S. et al. (2013)), CIFAR-100 (Krizhevsky, A. (2009)), ImageNette (Howard, J. (2020)), TinyImageNet Le, Y. et al. (2015), and VGGFace2 Cao, Q. et al. (2018). Details of these datasets can be found in Apdx. 7.1. Data allocation in our experiments strictly follows the assumptions in Sec. 4.1. For each dataset, we randomly sample 10% of the test set to form the small, clean dataset $\mathcal{D}_{\text{Defense}}$ assumed for the defender. The remaining test instances, denoted by $\mathcal{D}_{\text{Test}}$, are reserved for performance evaluation.

Attack settings: In this paper, we consider standard backdoor attacks launched by poisoning the training set of the classifier (Gu, T. et al., 2019; Chen, X. et al., 2017). In particular, we consider both the *all-to-one* (**A2O**) attacks and the *all-to-all* (**A2A**) attacks in our main experiments on CIFAR-10. For A2O attacks on CIFAR-10, we

arbitrarily choose class 9 as the target class; while for A2A attacks, as described in Sec. 4.1, triggered instances from any class $c \in \mathcal{C}$ are supposed to be (mis)classified to class $(c + 1) \bmod |\mathcal{C}|$. For each attack setting, we consider the following triggers: 1) a 3×3 random patch (**BadNet**) with a randomly selected location (fixed for all triggered images for each attack) used in Gu, T. et al. (2019); 2) an additive perturbation (with size $2/255$) resembling a chessboard (**CB**) used in Xiang, Z. et al. (2022); 3) a single pixel (**SP**) perturbed by $75/255$ with a randomly selected location (fixed for all triggered images for each attack) used by Tran, B. et al. (2018); 4) invisible triggers generated with l_0 and l_2 norm constraints (l_0 **inv** and l_2 **inv** respectively) proposed by Li, S. et al. (2021); 5) a warping-based trigger (**WaNet**) proposed by Nguyen, T. et al. (2021); 6) a Hello Kitty blending trigger (**Blend**) used by Chen, X. et al. (2017); 7) a trigger generated by the horizontal sinusoidal function (**SIG**) defined in Barni, M. et al. (2019). Details for generating these triggers are deferred to Apdx. 7.2. We randomly created 5 attacks for each attack setting *e.g.* by randomly locating the trigger. We also evaluated against the “label-consistent” (**CL**) backdoor attack proposed by Turner, A. et al. (2019) on the CIFAR-10 dataset, which only embeds the backdoor trigger into the target class training samples. Details are given in Apdx. 7.2. For experiments on the other five datasets, we only consider A2O attacks for a subset of triggers where sufficiently high success rate can be achieved. For each dataset, we create one attack for each trigger being considered. A2A attacks are not considered for these datasets since there is insufficient data per class for them to achieve successful attacks. More details about the attacks, including the number of backdoor-trigger images used for poisoning and the target class selected to create A2O attacks for the five datasets other than

CIFAR-10, are shown in Apdx. 7.2.

Performance evaluation metrics: 1) The attack success rate (**ASR**) is the fraction of clean instances in $\mathcal{D}_{\text{Test}}$ (mis)classified to the designated target class when the backdoor trigger is embedded. 2) The clean test accuracy (**ACC**) is the DNN’s accuracy on $\mathcal{D}_{\text{Test}}$ without trigger embedding. 3) The **SIA** (Def 3.3) is the fraction of clean instances in $\mathcal{D}_{\text{Test}}$ classified to the original source class when the trigger is embedded. For a successful backdoor attack, ASR and ACC should be high, while SIA should be low. For a successful mitigation approach, the resulting ASR should be low, while ACC and SIA should be high.

Training settings: We train one classifier for each attack to evaluate our mitigation approach against existing ones. Training configurations, including the DNN architecture, batch size, number of epochs, etc., are detailed in Tab. 11 in Apdx. 7.3. Data augmentation choices, including random cropping and horizontal flipping, are applied to each training instance. As shown in Tab. 1, the defenseless “vanilla” classifiers being attacked achieve high ACC but suffer high ASR and low SIA (averaged over the five attacks we created) for all trigger types and for both A2O and A2A settings, *i.e.*, the attacks are all successful and hence adequate for performance evaluation.

Hyper-parameter settings: We compare our mitigation approach (named ‘Batch Normalization Alteration’ (**BNA**) in the sequel) with six well-known and/or state-of-the-art methods, including NC(Wang, B. et al. (2019)), NAD(Li, Y. et al. (2021)), I-BAU(Zeng, Y. et al. (2022)), ANP(Wu, D. et al. (2021)), ARGD(Xia, J. et al. (2022)), and MCRZhao, P. et al. (2020). For MCR, in their original paper, the defender is assumed to have access to two poisoned models, which may be impractical. Thus, we fine-tune the given

model on the defender’s dataset and use it as the second model (which is also suggested in their paper). For all these other methods, we used their officially posted code for implementation. For BNA, following Sec. 4.2, we first perform detection by reverse-engineering a backdoor trigger for each class pair using objective functions from Wang, B. et al. (2019); Xiang, Z. et al. (2022) and then feed the statistics obtained based on the estimated trigger to an anomaly detector. Our anomaly detector is based on MAD, which is a classical approach also used by Wang, B. et al. (2019); Chen, H. et al. (2019); Wang, R. et al. (2020). Here, we set the detection threshold at “7-MAD” which easily catches all the backdoor class pairs. More details, including pattern estimation and detection statistics are shown in Apdx. 8. Then, for each detected target class, we solve the divergence minimization problem to optimize the transformation functions using learning rate 0.01 for 10 epochs. Since our mitigation method applies simple transformations which are also used in BN, we consider model structures that contain BN layers (which is very common) for simplicity. But note that the proof of monotonicity of SIA with distribution divergence (Thm. 3.1) and our method (Sec. 4.2) do not truly rely on the presence of BN layers – one can always insert a BN layer between any two given layers of the trained network. If a neuron is followed by a BN, instead of applying an additional transformation function $h_{j,l}$, we treat the mean and standard deviation of BN as the parameters $\mu_{j,l}$ and $\sigma_{j,l}$ associated with $h_{j,l}$ respectively. We optimize the mean and standard deviation by minimizing distribution divergence for all the BN layers. In Sec. 5.2, we only show results for BNA with the total variation divergence. Results for KL-divergence and JS-divergence are deferred to Apdx. 10. To compute the divergence, we use the “interval trick” (Eq. (5)) to obtain the discrete empirical distribution.

For simplicity, we let all finite intervals, $\mathcal{I}_i = [b_{i-1}, b_i)$, $i = 1, \dots, M$, have the same length $\Delta b = 0.1$. For each neuron, we set b_{\min} and b_{\max} as the minimum and maximum activations, respectively, when feeding in clean instances from $\mathcal{D}_{\text{Defense}}$ to the poisoned classifier f . Then, the number of intervals is $M = \lceil \frac{b_{\max} - b_{\min}}{\Delta b} \rceil$; and all intervals can be specified by $b_0 = b_{\min}$ and $b_i = b_{i-1} + \Delta b$. Finally, the scale factor in Eq. (6) is set to $\tau = 150$, which is obtained by line search to minimize the total variation between the “soft” distribution and the empirical one on $\mathcal{D}_{\text{Defense}}$. In fact, the choices for Δb and τ (over reasonable ranges) have little impact on our mitigation performance, as shown in Tab. 7.

5.2 Backdoor Mitigation Results

In Tab. 1, we show the ASR, ACC, and SIA for BNA compared with the other six methods (which are all DNN tuning-based) for attacks on CIFAR-10. Each metric is averaged over the five attacks created for each trigger type and attack setting, with the highest ACC and SIA, and the lowest ASR in bold. We found that these tuning-based methods are sensitive to the choices of hyper-parameters, such as the learning rate. Hence, for these methods, we optimize the hyper-parameter values to show the *best* results for these methods in Tab. 1. Although these tuning-based methods (except for MCR) can effectively deactivate backdoor attacks (*i.e.*, significantly reduce ASRs), there is a clear drop (3%-20%) in both ACC and SIA, compared with those for the vanilla DNN (the first row of Tab. 1). This is possibly due to tuning many DNN parameters using very limited data. (Note that BNA mitigation uses much less clean labelled

Trigger type	BadNet		CB		l_0 inv		l_2 inv		SP		WaNet		
	A2O	A2A	A2O	A2A	A2O	A2A	A2O	A2A	A2O	A2A	A2O	A2A	
Vanilla	ACC	0.9122	0.9121	0.9135	0.9098	0.9135	0.9131	0.9130	0.9126	0.9138	0.9060	0.9032	0.8994
	ASR	0.9573	0.8658	0.9685	0.8692	0.9989	0.8973	0.9889	0.8620	0.8912	0.8550	0.9153	0.8216
	SIA	0.0397	0.0432	0.0293	0.0257	0.0010	0.0151	0.0107	0.0194	0.1016	0.0593	0.0772	0.0714
NC	ACC	0.8797	0.8762	0.8735	0.8776	0.8835	0.8767	0.8750	0.8690	0.8854	0.8614	0.8748	0.8756
	ASR	0.0130	0.0154	0.0064	0.0155	0.0120	0.0150	0.0080	0.0179	0.0335	0.0188	0.0144	0.1381
	SIA	0.8532	0.8614	0.8312	0.8597	0.8654	0.8650	0.7932	0.8254	0.8362	0.8477	0.8231	0.7183
I-BAU	ACC	0.8500	0.8758	0.8812	0.8719	0.8452	0.8800	0.8825	0.8726	0.8666	0.8745	0.8777	0.8700
	ASR	0.0094	0.0164	0.1973	0.0811	0.0091	0.0133	0.2600	0.3353	0.0172	0.0154	0.1339	0.1253
	SIA	0.8301	0.8583	0.6399	0.7756	0.8277	0.8673	0.5549	0.4928	0.8479	0.8609	0.7059	0.7269
ANP	ACC	0.8644	0.8492	0.8241	0.8577	0.8455	0.8648	0.8345	0.8421	0.8195	0.8411	0.8298	0.8607
	ASR	0.0474	0.1199	0.3351	0.0927	0.0836	0.1326	0.4703	0.2648	0.1229	0.0495	0.0263	0.0835
	SIA	0.8184	0.7205	0.4587	0.7168	0.7697	0.7324	0.3351	0.4976	0.7060	0.7942	0.7368	0.7451
NAD	ACC	0.8814	0.8819	0.8800	0.8908	0.8958	0.9047	0.8991	0.8781	0.8813	0.8761	0.8592	0.8963
	ASR	0.0193	0.7132	0.0871	0.0681	0.0356	0.0457	0.0254	0.0191	0.0667	0.0647	0.0571	0.1056
	SIA	0.8498	0.1520	0.7711	0.8084	0.8504	0.8534	0.8221	0.8337	0.8123	0.8082	0.7710	0.7773
ARGD	ACC	0.8689	0.8482	0.8800	0.8774	0.8880	0.8885	0.8669	0.8583	0.8899	0.8728	0.8739	0.8755
	ASR	0.0368	0.0839	0.0099	0.0117	0.0079	0.0122	0.0125	0.0179	0.0955	0.0452	0.0111	0.0362
	SIA	0.8217	0.7544	0.8657	0.8690	0.8725	0.8786	0.8168	0.8297	0.7934	0.8295	0.8283	0.8241
MCR	ACC	0.8751	0.8840	0.8126	0.8618	0.8534	0.8834	0.8688	0.8481	0.8886	0.8613	0.8808	0.8661
	ASR	0.1447	0.1001	0.6015	0.1004	0.2900	0.0000	0.9769	0.0971	0.0136	0.1166	0.0268	0.0974
	SIA	0.7572	0.0811	0.3151	0.1163	0.0927	0.1000	0.0210	0.3524	0.8570	0.2239	0.7999	0.7789
BNA (ours)	ACC	0.9032	0.8951	0.9072	0.8615	0.9068	0.8944	0.9005	0.8638	0.9058	0.8921	0.8945	0.8792
	ASR	0.0139	0.0189	0.0127	0.0202	0.0033	0.0111	0.0042	0.0168	0.0104	0.0225	0.0041	0.0191
	SIA	0.8835	0.8841	0.8787	0.8820	0.8924	0.8942	0.8383	0.8522	0.8863	0.8811	0.8530	0.8607

Table 1: Average ACC, ASR, and SIA for BNA, compared with NC, NAD, I-BAU, ANP, and ARGD, against all the created attacks applied to ResNet-18 trained on the CIFAR-10 dataset. Best performances are indicated in bold.

		Vanilla	NC	I-BAU	ANP	NAD	ARGD	MCR	BNA
	ACC	0.9062	0.9061	0.1735	0.8998	0.3354	0.2362	0.8829	0.8967
CL	ASR	0.9304	0.1444	0.4940	0.4632	0.0636	0.0581	0.7756	0.0594
	SIA	0.0667	0.7558	0.0896	0.4977	0.2985	0.2349	0.1932	0.8057

Table 2: ACC, ASR, and SIA for BNA, compared with those of NC, NAD, I-BAU, ANP, and ARGD, against the ResNet-18 trained on the CIFAR-10 dataset poisoned by the label-consistent (CL) backdoor attack.

	VGGFace2				CIFAR-10							
	Vanilla	NC	I-BAU	BNA	Vanilla	NC	I-BAU	ANP	NAD	ARGD	MCR	BNA
ACC	0.8989	0.8967	0.6828	0.8917	0.9122	0.8848	0.8539	0.6463	0.8731	0.4946	0.8650	0.9026
ASR	0.9771	0.9693	0.9737	0.0046	0.9573	0.2531	0.4175	0.0117	0.8880	0.0227	0.8078	0.0185
SIA	0.0216	0.0294	0.0196	0.8889	0.0397	0.6842	0.5148	0.6324	0.1007	0.4731	0.1684	0.8807

Table 3: ACC, ASR, and SIA for BNA, compared with those of NC, I-BAU, ANP, NAD, ARGD, and MCR, with limited amount of clean data on VGGFace2 and CIFAR-10. Both datasets are poisoned by the BadNet attack.

data than what was reported for these other methods in their original papers.) Though MCR can effectively deactivate most of the backdoor attacks, excluding the global pattern CB and l_2 inv, it fails to infer the true source classes for the backdoor-triggered instances. For ANP with neuron pruning, the performance is acceptable only for A2O with the BadNet trigger. One possible reason is that invisible, perturbation-based triggers affect most neurons only moderately (which is also discussed in Wang, H. et al. (2022)); thus, pruning a small number of neurons cannot mitigate the attack. In contrast, our method successfully mitigates all these backdoor attacks (with generally the best ACC and ASR compared with the others) regardless of the trigger type and attack setting. Notably, *since the purpose of BNA’s divergence minimization is to maximize the*

SIA, it unsurprisingly achieves the best *SIA* with a clear margin over all other methods, in all cases (the corresponding distribution divergences are shown in Tab. 13 in Apdx. 9). We also tune the poisoning ratio and perturbation size used in A2O CB attacks, and the performance for BNA slightly declines as the attack is strengthened, as shown in Tab. 4. However, it still outperforms the other methods (see Tab. 15 in Apdx. 11).

	Number of poisoned instances per class					Perturbation size (*255)				
	50	100	150	200	250	2	3	4	5	6
ACC	0.9112	0.9094	0.9098	0.9102	0.9015	0.9094	0.9041	0.9079	0.8992	0.8912
ASR	0.0095	0.0141	0.0121	0.0170	0.0090	0.0141	0.0395	0.0222	0.0109	0.0388
SIA	0.8851	0.8837	0.8728	0.8840	0.8662	0.8851	0.8783	0.8711	0.8814	0.8435

Table 4: ACC, ASR, and SIA for BNA as a function of (1) the number of poisoned instances injected into the training set; (2) the perturbation size under all-to-one CB attack.

For the CL attack, we poison half (2500) of the target-class training samples to achieve an effective attack (which is stronger than in the original paper (Turner, A. et al., 2019)), as shown in Tab. 2. Although NAD and ARGD effectively deactivate the attack, both ACC and ASR drop significantly. For NC, ANP, and MCR, the ACC after mitigation is almost the same as the ACC before mitigation, but the ASR is still high. For ANP, nearly half of the backdoor-trigger images are unimpeded by the mitigation system. The attack is still effective after MCR is deployed. I-BAU does not perform well in mitigating the CL attack – the mitigated model fails to correctly classify most of the clean test images, but still recognizes half of the backdoor-trigger images to the target class. By contrast, BNA decreases the ACC by only a small amount, reduces ASR to around 6%, and correctly classifies 80% of the backdoor-trigger images.

Results of BNA on other datasets are shown in Tab. 3 and 5. We first train a DNN on the VGGFace2 dataset poisoned by the BadNet attack. As shown in Tab. 3, the BadNet attack is effective, with a high ASR and a nearly unchanged ACC. (The ACC for the DNN trained on the clean VGGFace2 dataset is 0.9211.) We then apply BNA, NC, and I-BAU on the poisoned DNN.⁶ For all mitigation methods, we only preserve 10 clean images per class since there is a severely limited number of samples for VGGFace2. BNA effectively reduces the ASR and yields a high SIA, outperforming the other mitigation methods. We will thoroughly discuss the impact of the number of clean images possessed by the defender in Sec. 5.6. The ACC for DNNs trained without attack for GTSRB, CIFAR-100, ImageNette, and TinyImageNet are 0.9567, 0.6926, 0.8726, and 0.5224, respectively; while ACC, ASR, and SIA for attacked DNNs are shown in the row “Vanilla” in Tab. 5, which demonstrate that all the attacks are effective. We apply BNA on the poisoned DNNs, with the same settings as for CIFAR-10, which significantly reduces ASR (to less than 1.3% in all cases), with uniformly high SIA and ACC.

We also evaluated the performance of our BNA against all-to-one backdoor attacks that utilize more complex global backdoor patterns, such as the blended backdoor attack (Chen, X. et al., 2017) and Sinusoidal Signal backdoor attack (SIG) (Barni, M. et al., 2019). Details of the attack configurations can be found in Apdx. 7.2. The performance of our BNA as well as other mitigation methods are shown in Tab. 6. The results

⁶We did not evaluate the performance of ANP, NAD, and ARGD on VGGFace2, since these references do not provide the architecture of VGG-16 that fits their respective mitigation system. Although MCR provides the architecture of VGG-16, it is different from the one provided by PyTorch. Therefore we cannot load the pre-trained weights and make a fair comparison.

demonstrate the effectiveness of our BNA mitigation method even when dealing with complicated backdoor patterns. Compared with the other methods⁷, our BNA significantly reduces the ASRs and produces relatively satisfactory SIAs, while maintaining ACCs that are competitive with pre-mitigation figures.

Trigger	GTSRB					CIFAR-100				TinyImageNet	ImageNette	
	BadNet	CB	l_0 inv	l_2 inv	WaNet	BadNet	CB	l_0 inv	l_2 inv	BadNet	BadNet	
Vanilla	ACC	0.9517	0.9556	0.9531	0.9521	0.9408	0.6796	0.6917	0.6863	0.6804	0.5192	0.8626
	ASR	1.0000	1.0000	1.0000	0.9794	0.9000	0.9037	0.9169	0.9935	0.9097	0.8058	0.9144
	SIA	0.0000	0.0000	0.0000	0.0169	0.0905	0.0781	0.0646	0.0063	0.0707	0.1134	0.0771
BNA	ACC	0.9491	0.9548	0.9505	0.9500	0.9404	0.6770	0.6863	0.6858	0.6787	0.5178	0.7941
	ASR	0.0000	0.0000	0.0122	0.0001	0.0041	0.0002	0.0524	0.0062	0.0016	0.0043	0.0016
	SIA	0.9312	0.9454	0.9330	0.8945	0.9338	0.6526	0.5880	0.6169	0.5224	0.4965	0.7940

Table 5: ACC, ASR, and SIA for BNA against all-to-one attacks on CIFAR-100, GTSRB, ImageNette, and TinyImageNet datasets.

		Vanilla	NC	I-BAU	ANP	NAD	ARGD	BNA
	ACC	0.9264	0.8132	0.7932	0.8667	0.8221	0.4856	0.8942
Blend	ASR	0.9731	0.0488	0.7419	0.5119	0.0521	0.0782	0.1283
	SIA	0.0254	0.5593	0.1427	0.3369	0.6003	0.4181	0.6252
	ACC	0.9266	0.8234	0.5696	0.8251	0.7794	0.4233	0.8716
SIG	ASR	0.9991	0.1414	0.2594	0.9451	0.3223	0.0988	0.0158
	SIA	0.0008	0.2980	0.1319	0.0383	0.2503	0.3271	0.3357

Table 6: ACC, ASR, and SIA for BNA, compared with those of NC, NAD, I-BAU, ANP, and ARGD, against the ResNet-18 trained on the CIFAR-10 dataset poisoned by the all-to-one Blend and SIG backdoor attacks.

⁷We weren’t able to reproduce the results reported in the published papers describing these methods due to different defense settings – in our experiments, the defender possesses far fewer clean samples.

τ ($\Delta b=0.1$)	10	100	200	300	400	500	600	700	800	900	1000
ACC	0.9024	0.9025	0.9022	0.9019	0.9024	0.9019	0.9022	0.9018	0.9017	0.9015	0.9021
ASR	0.0257	0.022	0.0206	0.0207	0.0214	0.0202	0.0212	0.0209	0.0207	0.0201	0.0204
SIA	0.8744	0.8758	0.8774	0.8768	0.8768	0.8775	0.8773	0.8768	0.8772	0.8778	0.8777
Δb ($\tau=150$)	0.1	0.11	0.12	0.13	0.14	0.15	0.16	0.17	0.18	0.19	0.2
ACC	0.9028	0.9018	0.9019	0.9023	0.9019	0.9023	0.9019	0.9019	0.9022	0.9022	0.9022
ASR	0.0197	0.0202	0.0199	0.0204	0.0199	0.0204	0.0198	0.0206	0.0206	0.0207	0.0212
SIA	0.8799	0.8775	0.8779	0.8773	0.8772	0.8779	0.8773	0.8767	0.8779	0.8769	0.8764

Table 7: ACC, ASR, and SIA for BNA as a function of scale factor and bin size on ResNet-18 trained on CIFAR-10 poisoned by all-to-one BadNet attack.

5.3 The “detection-before-mitigation” scenario

As discussed in Sec. 4.1, BNA performs backdoor mitigation only after the model has been detected as backdoor-poisoned. To justify the “detection-before-mitigation scenario”, we first apply the mitigation methods that do not involve a detection system (*i.e.*, I-BAU, ANP, NAD, ARGD, and MCR) on a ResNet-18 trained on the attack-free CIFAR-10 dataset. The resulting (absolute) drop in ACC is shown in column “clean” in Tab. 9. ANP has the largest impact on ACC – the ACC drops by 0.1877 after mitigation. I-BAU and ARGD respectively decrease the ACC by 0.0684 and 0.0343. NAD and MCR keep the ACC almost as high as that of the vanilla model, but they are not sufficiently effective in terms of SIA when the model is poisoned. NC and BNA detect the backdoor attack before mitigation; thus there is no impact on the ACC for clean classifiers (for which no attack is detected). We also found reduction in ACCs when applying all mitigation methods for a ResNet-18 model trained on CIFAR-10 poisoned by an all-to-one BadNet attack in column “BadNet”. All the other methods decrease

ACC by more than 0.03, while our method has little impact on ACC.

Trigger	BadNet		CB		l_0 inv		l_2 inv		SP		WaNet		
	A2O	A2A	A2O	A2A	A2O	A2A	A2O	A2A	A2O	A2A	A2O	A2A	
BNA	FPR	0.1390	0.0606	0.1144	0.1092	0.1413	0.0600	0.1976	0.1027	0.1323	0.0656	0.1406	0.0865
	TPR	0.9872	0.9508	0.9872	0.9682	0.9967	0.9873	0.9958	0.9793	0.9894	0.9294	0.9959	0.9248
STRIP	FPR	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15	0.15
	TPR	0.9638	0.5147	0.6802	0.2089	0.9995	0.1053	0.9924	0.5272	0.8522	0.3123	0.0202	0.0411

Table 8: TPR and FPR for BNA, compared with STRIP, against all attacks created on CIFAR-10.

I-BAU		ANP		NAD		ARGD		MCR		NC		BNA	
Clean	BadNet	Clean	BadNet	Clean	BadNet	Clean	BadNet	Clean	BadNet	Clean	BadNet	Clean	BadNet
0.0684	0.0622	0.1877	0.0478	0.0064	0.0308	0.0343	0.0433	0.0038	0.0371	NA	0.0325	NA	0.0090

Table 9: Drop in ACC when applying I-BAU, ANP, NAD, ARGD, NC and BNA on ResNet-18 trained on the clean (attack-free) CIFAR-10 dataset and trained on CIFAR-10 poisoned by the all-to-one BadNet attack.

5.4 Test-Time backdoor-trigger instance detection

Different from other tuning-based backdoor mitigation approaches, our BNA can also detect backdoor-trigger instances at test-time, as described in Sec. 4.2 and shown in Fig. 2. Here, we evaluate accuracy of our test-time detector compared with a state-of-the-art detector named STRIP (Gao, Y. et al., 2019). For any input image during inference, STRIP blends it with clean images possessed by the defender. The blended image is fed into the poisoned DNN, with an entropy calculated on the output posteriors. If the entropy is lower than a prescribed detection threshold, the input is deemed to

be embedded with the trigger. Here, we set the detection threshold to achieve 15% FPR for STRIP, a choice which achieves a generally good trade-off between TPR and FPR. In contrast, BNA does not require setting a detection threshold. In Tab. 8, we show the True Positive Rate (TPR, *i.e.*, the fraction of backdoor-trigger images correctly detected) and the False Positive Rate (FPR, *i.e.*, the fraction of clean test images from the backdoor target class(es) that are falsely detected) for both methods. Although STRIP performs well on A2O attacks for some trigger types, *e.g.*, BadNet, l_0 inv, and l_2 inv, its TPR drops drastically on attacks using human-imperceptible triggers, especially the WaNet attacks. Moreover, it does not perform well on all A2A attacks, with a largest TPR of only 0.5272. By contrast, BNA is effective for all these attacks – it detects almost all the backdoor-trigger images, with FPRs comparable to STRIP.

5.5 Mitigation performance against adaptive attacks

A recent backdoor attack proposed by Doan, K. et al. (2021) minimizes a metric similar to that used by our BNA defense, in order to achieve better stealthiness. This attack can be viewed as an adaptive attack against our mitigation defense since the trained classifier will be more sensitive to even a smaller distribution divergence than for ordinary backdoor attacks. Nevertheless, our method successfully mitigates this attack. In our experiment on CIFAR-10, the average distribution total variation divergence over all neurons is reduced from 8067 to 2789. Accordingly, the ACC/ASR before and after mitigation are 0.9162/0.9978 and 0.8906/0.0072, respectively, with an SIA of 0.8496.

5.6 Mitigation with a limited amount of clean data

Why tuning-based methods like NC cannot achieve SIAs as high as BNA (which does not alter the DNN’s parameters) ? Note that NC tunes the classifier using instances embedded with the estimated trigger but without label flipping. This is equivalent to minimizing the divergence between internal activation distributions for clean and triggered instances (see Tab. 13 in Apdx. 9), but by altering the DNN’s parameters. Even for an optimal (zero) divergence, the best achievable SIA of NC is still upper-bounded by the ACC of the classifier after tuning, which usually *drops* due to the **data insufficiency**. By contrast, the reference distribution for BNA’s divergence minimization is obtained by feeding clean instances to the poisoned classifier *without* changing its parameters; thus, it is a “better” reference with a higher upper-bound ACC.

For the main experiments on CIFAR-10, we preserve 100 clean test images for all mitigation methods. In other words, all mitigation methods, excluding our BNA, tune the (around 11 million) trainable parameters of the poisoned ResNet-18 based only on 1000 clean labelled images (2% of the training set) for a few epochs. Our method is light-weight, since it only updates the (less than 10 thousand) non-trainable parameters (*i.e.*, mean and standard deviation). Note that all the mitigation methods use more clean images in their original papers than in the experiments reported herein. For example, NC and NAD respectively chose 10% and 5% of the clean training samples for mitigation. For all mitigation methods, excluding our BNA, the **insufficiency** of clean labelled data reduces the ACC by at most 10%. The SIA is upper-bounded by the ACC after mitigation and is also affected by data insufficiency. This is also verified in Li, Y. et al. (2021), where they varied the number of clean images from 0% to 20% of

the clean training samples, with the performance of NAD significantly degraded as the number of clean samples decreases. Our BNA only aligns internal layer distributions without affecting the trainable parameters, and thus is more robust when the amount of clean samples is limited.

To further demonstrate the impact of the number of clean samples on mitigation performance, for the ResNet-18 trained on CIFAR-10 poisoned by the BadNet attack, we reduce the number of clean images used by the defender to just 10 images per class. The corresponding performance of all mitigation methods is shown in Tab. 3. Although ANP and ARGD effectively de-activate the backdoor attack, both ACC and SIA dramatically decrease. For NC, I-BAU, NAD, and MCR, the ACC changes a little, but the attack is still effective, especially for NAD and MCR. However, our BNA is still effective, with the ASR less than 2% and both ACC and SIA comparable to the ACC before mitigation.

Data insufficiency is a common phenomenon in real-world applications. For example, we use a subset of VGGFace2, which consists of 18 identities, each of which has 450 training face images and 100 test face images. So, VGGFace2 is much smaller than other benchmark datasets such as CIFAR-10. We only assign 10 clean images per class for the defender. The results are shown in Tab. 3. On this high-resolution and insufficient dataset, both NC and I-BAU fail to de-activate the BadNet attack. In contrast, our BNA successfully reduces the ASR to 0.46% and has ACC and SIA about 89%.

Conclusion

In this paper, we revealed an activation distribution alteration property for backdoor attacks. We theoretically proved that by correcting such alteration, backdoor trigger instances will be correctly classified to their original source classes. Accordingly, we proposed a post-training backdoor mitigation approach to align distributions of clean and backdoor-trigger samples through simple transformations, *without changing millions trainable parameters* of the classifier, which outperformed methods that use DNN fine-tuning. The proposed method is *robust* especially when there is *limited amount* of clean data available to the defender, compared with parameter-tuning based methods. Besides, the proposed method is *flexible* to be integrated with existing detection systems. Moreover, our method can detect instances with the trigger during inference.

Ethics Statement

The main purpose of this research is to understand the behavior of deep learning systems facing malicious activities, and enhance their safety level. The backdoor attack considered in this paper is well-known, with open-sourced implementation code. Thus, publication of this paper will be beneficial to the community in defending against backdoor attacks. The code of our defense will be released if the paper is accepted.

Acknowledgments

This research was supported in part by NSF SBIR grant 2132294.

Appendix

6 Proof of Theorems in the Main Paper

6.1 Derivation of Eq. (2)

Here, we provide the derivation showing that m_j^* and v_j^* in Eq. (2) are the solutions to:

$$\mathbb{E}[a_j^*(\mathbf{X}_b)] = \mathbb{E}[a_j(\mathbf{X})|Y = -1] \quad (7)$$

$$\text{Var}[a_j^*(\mathbf{X}_b)] = \text{Var}[a_j(\mathbf{X})|Y = -1] \quad (8)$$

Based on Eq. (1), the above equations can be expanded as follows:

$$\mathbb{E}\left[\frac{\mathbf{w}_j^\top \mathbf{X}_b - m_j^*}{\sqrt{v_j^*}} \gamma_j + \beta_j\right] = \mathbb{E}\left[\frac{\mathbf{w}_j^\top \mathbf{X} - m_j}{\sqrt{v_j}} \gamma_j + \beta_j | Y = -1\right] \quad (9)$$

$$\text{Var}\left[\frac{\mathbf{w}_j^\top \mathbf{X}_b - m_j^*}{\sqrt{v_j^*}} \gamma_j + \beta_j\right] = \text{Var}\left[\frac{\mathbf{w}_j^\top \mathbf{X} - m_j}{\sqrt{v_j}} \gamma_j + \beta_j | Y = -1\right] \quad (10)$$

We first solve Eq. (10) for $(\mathbf{X}|Y = -1) \sim \mathcal{N}(-\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ and $\mathbf{X}_b \sim \mathcal{N}(\boldsymbol{\mu}_b, \sigma_b^2 \mathbf{I})$, which leads to:

$$v_j^* = \frac{\sigma_b}{\sigma} v_j \quad (11)$$

By substituting Eq. (11) into Eq. (9), and since $\boldsymbol{\mu}_b = -\boldsymbol{\mu} + \boldsymbol{\epsilon}$, we get the following:

$$\begin{aligned} m_j^* &= \sqrt{\frac{v_j^*}{v_j}} (m_j - \mathbf{w}_j^\top \boldsymbol{\mu}) + \mathbf{w}_j^\top \boldsymbol{\mu}_b \\ &= \frac{\sigma_b}{\sigma} m_j + \left(\frac{\sigma_b}{\sigma} - 1\right) \mathbf{w}_j^\top \boldsymbol{\mu} + \mathbf{w}_j^\top \boldsymbol{\epsilon}. \end{aligned}$$

6.2 Proof of Theorem 3.1

Proof. First, let's specify the following vector/matrix representations that will be used throughout this proof:

$$\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_J]^\top \in \mathbb{R}^{J \times d}$$

$$\mathbf{V} = \begin{bmatrix} v_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & v_J \end{bmatrix} \in \mathbb{R}^{J \times J} \quad \hat{\mathbf{V}}(\alpha) = \begin{bmatrix} \hat{v}_1(\alpha) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{v}_J(\alpha) \end{bmatrix} \in \mathbb{R}^{J \times J}$$

$$\mathbf{m} = \begin{bmatrix} m_1 \\ \vdots \\ m_J \end{bmatrix} \in \mathbb{R}^J \quad \hat{\mathbf{m}}(\alpha) = \begin{bmatrix} \hat{m}_1(\alpha) \\ \vdots \\ \hat{m}_J(\alpha) \end{bmatrix} \in \mathbb{R}^J \quad \mathbf{\Gamma} = \begin{bmatrix} \gamma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \gamma_J \end{bmatrix} \in \mathbb{R}^{J \times J} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_J \end{bmatrix} \in \mathbb{R}^J$$

$$\mathbf{a}(\cdot) = \begin{bmatrix} a_1(\cdot) \\ \vdots \\ a_J(\cdot) \end{bmatrix} \in \mathbb{R}^J \quad \hat{\mathbf{a}}(\cdot|\alpha) = \begin{bmatrix} \hat{a}_1(\cdot|\alpha) \\ \vdots \\ \hat{a}_J(\cdot|\alpha) \end{bmatrix} \in \mathbb{R}^J \quad \mathbf{a}^*(\cdot) = \begin{bmatrix} a_1^*(\cdot) \\ \vdots \\ a_J^*(\cdot) \end{bmatrix} \in \mathbb{R}^J$$

Let $\mathbf{X}_- = (\mathbf{X}|Y = -1) \sim \mathcal{N}(-\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ denote a random instance from source class ‘-1’ for simplicity. Let $\mathbf{u} = \mathbf{u}_- - \mathbf{u}_+$ with \mathbf{u}_- and \mathbf{u}_+ being the weight vectors associated with the node for class ‘-1’ and the node for class ‘+1’ respectively. Then, it is easy to see that:

$$\hat{\mathbf{a}}(\mathbf{X}_b|\alpha)|_{\alpha=1} = \mathbf{a}(\mathbf{X}_b) \quad \text{and} \quad \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)|_{\alpha=0} = \mathbf{a}^*(\mathbf{X}_b),$$

and taking one step further by setting $\alpha = 1$, we have the following:

$$\begin{aligned}
P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)|\alpha = 1] &= P[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha) > 0|\alpha = 1] & (12) \\
&= P[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_b) > 0] \\
&= P[f_-(\mathbf{X}_b) > f_+(\mathbf{X}_b)] \\
&\leq 1 - \psi. & (13)
\end{aligned}$$

This is, when $\alpha = 1$, the classifier is not modified at all; thus the SIA will be no larger than $1 - \psi$ since the attack is ψ -successful (see Definition 3.2). On the other hand, by setting $\alpha = 0$, we will have the following:

$$\begin{aligned}
P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)|\alpha = 0] &= P[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha) > 0|\alpha = 0] \\
&= P[\mathbf{u}^\top \mathbf{a}^*(\mathbf{X}_b) > 0] \\
&= P[f_-(\mathbf{X}_-) > f_+(\mathbf{X}_-)] & (14) \\
&\geq 1 - \eta. & (15)
\end{aligned}$$

That is, when $\alpha = 0$, the distribution shift will be fully recovered, such that SIA is equally high as the accuracy of the source class. Recall that Eq. (14) is due to Eq. (7) and Eq. (8). The inequality (15) is because the classifier specified by f_- and f_+ is assumed η -erroneous (see Definition 3.1). Here, we prove the theorem by showing that the partial derivative of $P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)]$ over α is strictly *negative* when $\sigma_b \leq \sigma$ (*i.e.*, triggered instances have smaller standard deviation than clean instances, which is generally true). To achieve this, we notice that

$$\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha) = \mathbf{u}^\top \hat{\mathbf{V}}(\alpha)^{-\frac{1}{2}} \Gamma(\mathbf{W}\mathbf{X}_b - \hat{\mathbf{m}}(\alpha)) + \mathbf{u}^\top \boldsymbol{\beta}$$

follows a Gaussian distribution with

$$\mathbb{E}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)] = \mathbf{u}^\top \hat{\mathbf{V}}(\alpha)^{-\frac{1}{2}} \Gamma(-\mathbf{W}\boldsymbol{\mu} + \mathbf{W}\boldsymbol{\epsilon} - \hat{\mathbf{m}}(\alpha)) + \mathbf{u}^\top \boldsymbol{\beta} \quad (16)$$

$$\text{Var}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)] = \sigma_b^2 \|\mathbf{W}^\top \Gamma \hat{\mathbf{V}}(\alpha)^{-\frac{1}{2}} \mathbf{u}\|_2^2 \quad (17)$$

We also notice that for source class instances \mathbf{X}_-

$$\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-) = \mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma(\mathbf{W}\mathbf{X}_- - \mathbf{m}) + \mathbf{u}^\top \boldsymbol{\beta}$$

follows a Gaussian distribution with

$$\mathbb{E}[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-)] = \mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma(-\mathbf{W}\boldsymbol{\mu} - \mathbf{m}) + \mathbf{u}^\top \boldsymbol{\beta} \quad (18)$$

$$\text{Var}[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-)] = \sigma^2 \|\mathbf{W}^\top \Gamma \mathbf{V}^{-\frac{1}{2}} \mathbf{u}\|_2^2 \quad (19)$$

Then we have

$$P[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha) > 0] = 1 - \Phi\left(-\frac{\mathbb{E}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}{\sqrt{\text{Var}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}}\right) \quad (20)$$

$$P[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-) > 0] = 1 - \Phi\left(-\frac{\mathbb{E}[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-)]}{\sqrt{\text{Var}[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-)]}}\right) \quad (21)$$

where Φ is the cumulative distribution function of standard Gaussian. Now let's consider Eq. (21) first. Since $\eta < \frac{1}{2}$ as we have reasonably assumed (otherwise the classifier may be worse than a random guess), and also according to Eq. (15), we have

$$P[\mathbf{u}^\top \mathbf{a}(\mathbf{X}_-) > 0] = P[f_-(\mathbf{X}_-) > f_+(\mathbf{X}_-)] > \frac{1}{2}$$

Thus, based on Eq. (21) and Eq. (18), we get

$$\mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma(-\mathbf{W}\boldsymbol{\mu} - \mathbf{m}) + \mathbf{u}^\top \boldsymbol{\beta} > 0 \quad (22)$$

Next, let us focus on Eq. (20). Again, we set $\alpha = 1$. Based on (12)-(13) and the reasonable assumption that $\psi > \frac{1}{2}$ (otherwise the attack is not deemed successful since

the success rate will be even lower than the accuracy on clean instances), we have

$$P[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha) > 0 | \alpha = 1] = P[f_-(\mathbf{X}_b) > f_+(\mathbf{X}_b)] < \frac{1}{2}$$

Then, based on Eq. (20) and Eq. (16), we get

$$\mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma (-\mathbf{W}\boldsymbol{\mu} + \mathbf{W}\boldsymbol{\epsilon} - \mathbf{m}) + \mathbf{u}^\top \boldsymbol{\beta} < 0 \quad (23)$$

Subtracting Eq. (23) from Eq. (22) we get:

$$-\mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma \mathbf{W}\boldsymbol{\epsilon} > 0 \quad (24)$$

Based on Eq. (20), we also have

$$\frac{\partial P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)]}{\partial \alpha} = \frac{\partial \frac{\mathbb{E}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}{\sqrt{\text{Var}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}}{\partial \alpha} \cdot \phi\left(-\frac{\mathbb{E}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}{\sqrt{\text{Var}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}}\right) \quad (25)$$

where ϕ is the probability density function (PDF) for standard normal distribution.

Based on Eq. (16), Eq. (17), and Eq. (2), we have

$$\begin{aligned} & \frac{\mathbb{E}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}{\sqrt{\text{Var}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}} \\ &= \frac{\mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma [-\mathbf{W}\boldsymbol{\mu} + \mathbf{W}\boldsymbol{\epsilon} - (\alpha \mathbf{m} + (1-\alpha)\mathbf{m}^*)] + (\alpha + (1-\alpha)\frac{\sigma_b}{\sigma}) \mathbf{u}^\top \boldsymbol{\beta}}{\sigma_b \|\mathbf{W}^\top \Gamma \mathbf{V}^{-\frac{1}{2}} \mathbf{u}\|_2} \\ &= \frac{\alpha \cdot \mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma [(\frac{\sigma_b}{\sigma} - 1)\mathbf{m} + (\frac{\sigma_b}{\sigma} - 1)\mathbf{W}\boldsymbol{\mu} + \mathbf{W}\boldsymbol{\epsilon}] - \alpha \cdot (\frac{\sigma_b}{\sigma} - 1) \mathbf{u}^\top \boldsymbol{\beta}}{\sigma_b \|\mathbf{W}^\top \Gamma \mathbf{V}^{-\frac{1}{2}} \mathbf{u}\|_2} + \text{constant} \end{aligned}$$

and thus, based on Eq. (23) and Eq. (24)

$$\begin{aligned} & \frac{\partial \frac{\mathbb{E}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}{\sqrt{\text{Var}[\mathbf{u}^\top \hat{\mathbf{a}}(\mathbf{X}_b|\alpha)]}}{\partial \alpha} \\ &= \frac{(\frac{\sigma_b}{\sigma} - 1)[\mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma (\mathbf{m} + \mathbf{W}\boldsymbol{\mu} - \mathbf{W}\boldsymbol{\epsilon}) - \mathbf{u}^\top \boldsymbol{\beta}] + \frac{\sigma_b}{\sigma} \mathbf{u}^\top \mathbf{V}^{-\frac{1}{2}} \Gamma \mathbf{W}\boldsymbol{\epsilon}}{\sigma_b \|\mathbf{W}^\top \Gamma \mathbf{V}^{-\frac{1}{2}} \mathbf{u}\|_2} \end{aligned}$$

< 0

when $\sigma_b \leq \sigma$. Substitute it into Eq. (25) and given the Gaussian PDF being strictly positive, we have

$$\frac{\partial P[g_-(\mathbf{X}_b|\alpha) > g_+(\mathbf{X}_b|\alpha)]}{\partial \alpha} < 0$$

□

7 Datasets, Training settings, and attack settings

7.1 Datasets

In the experiments, we show the effectiveness of our proposed backdoor mitigation method on several benchmark datasets including CIFAR-10 (Krizhevsky, A. (2009)), GTSRB (Houben, S. et al. (2013)), CIFAR-100 (Krizhevsky, A. (2009)), ImageNette (Howard, J. (2020)), and TinyImageNet. CIFAR-10 dataset contains 60,000 32×32 color images from 10 classes, with 5,000 images per class for training and 1,000 images per class for testing. GTSRB has more than 50,000 traffic sign images with different sizes from 43 classes. Here, we resize all images in GTSRB to 32×32 . CIFAR-100 contains 60,000 32×32 color images evenly from 100 classes, where 500 images per class are used for training, while the others are used for testing. ImageNette is a subset of 10 easily classified classes from Imagenet⁸, with image size of 256×256 . For each class, there are around 900 images for training and 400 images for testing. TinyImageNet is a subset of ImageNet (Russakovsky, O. et al. (2015)). It contains 100,000 64×64 color images evenly distributed in 200 classes (500 training images

⁸The 10 classes are tench, English springer, cassette player, chain saw, church, French horn, garbage truck, gas pump, golf ball, and parachute.

and 50 test images for each class). VGGFace2 is a large-scale face recognition dataset. It contains more than 3.3 million face images from more than 9000 identities. In our experiments, we use a subset of VGGFace2, which consists of 18 identities, each of which has 550 face images. For each class (identity), 450 images are used for training and 100 images are used for testing. All images are resized to 224×224 pixels.

7.2 Attack Settings

On CIFAR-10, we consider the following triggers: 1) a 3×3 random patch (**BadNet**) with a randomly selected location (fixed for all triggered images for each attack) used in Gu, T. et al. (2019), as visualized in Fig. 3b; 2) an additive perturbation (with size $2/255$) resembling a chessboard (**CB**) used in Xiang, Z. et al. (2022), as visualized in Fig. 3c; 3) a single pixel (**SP**) perturbed by $75/255$ with a randomly selected location (fixed for all triggered images for each attack) used by Tran, B. et al. (2018), as visualized in Fig. 3d; 4) invisible triggers generated with l_0 and l_2 norm constraints (l_0 **inv** and l_2 **inv** respectively) proposed by Li, S. et al. (2021), as visualized in Fig. 3e and 3f; 5) a warping-based trigger (**WaNet**) proposed by Nguyen, T. et al. (2021), as visualized in Fig. 3g; 6) a Hello Kitty trigger with a blend ratio of $\alpha = 0.15$ used by Chen, X. et al. (2017), as is visualized in Fig. 4b; 7) a sinusoidal signal trigger generated by the horizontal sinusoidal function defined in Barni, M. et al. (2019) with $\Delta = 20$ and $f = 6$, as visualized in Fig. 4c.

Attack settings for CIFAR-10 are summarized in Tab. 10. For all-to-one attacks, we arbitrarily choose class 9 as the target class, and embed the backdoor triggers in 100 randomly chosen training samples per class (excluding the target class). To achieve

similarly effective attacks as for other triggers, we poison 900 images per source class in the all-to-one attack using WaNet. For all-to-all attacks, we embed the backdoor triggers into 300 images for each class. For effective attacks, we poison 800 training images and 1500 training images in the all-to-all attacks using SP and WaNet, respectively.

For CIFAR-10, we also launched the label-consistent (CL) backdoor attack proposed in Turner, A. et al. (2019), as visualized in Fig. 5b. We followed the same settings as reported in their paper. Specifically, we used Projected Gradient Descent (PGD) to generate adversarial perturbations bounded by a L_{inf} maximum perturbation of $\epsilon = 16$. A 3×3 white square is embedded at the four corners of the perturbed images. Half (2500) of the target class training images are poisoned by the CL attack. All non-target class test images are also poisoned in the same way for performance evaluation.

Attack settings for other datasets are summarized in Tab. 12. Due to the insufficiency of data, we only conduct all-to-one attacks on these datasets for effective attacks. We arbitrarily choose class 0 as the target class for CIFAR-100, GTSRB, TinyImageNet, and VGGface2, and class 9 for ImageNet. All the classes other than the target class are source classes of the attack. For CIFAR-100, we use the same BadNet, l_0 inv, and l_2 inv triggers as for CIFAR-10. We increase the perturbation size to $6/255$ for the CB pattern for an effective backdoor attack. For each attack, we poison 10 images per source class using the above triggers. Triggers SP and WaNet are not considered since we cannot launch a successful backdoor attack using these triggers on CIFAR-100. For GTSRB, in addition to the same triggers as for CIFAR-100, we also use the warping-based trigger (WaNet). We poison 2% of the training images per source class using the BadNet and l_2 inv triggers, and 5% of the training images per source class with the CB and l_0 inv trig-

gers. To achieve similarly effective attacks, we embed the WaNet trigger into 24% of the training images per source class. For TinyImageNet, ImageNette, and VGGFace2, we only consider BadNet as the trigger, as the DNN cannot learn the backdoor mapping using the other (relatively simple and small) triggers in datasets that are much more complicated than CIFAR-10. To successfully plant backdoors, we increase the size the the BadNet patch to 6×6 for TinyImageNet and to 21×21 for ImageNette and VGGFace2. We embed the trigger in 10 training images per source class in TinyImageNet and VGGFace2, and in 5% of the training images per source class for ImageNette.

7.3 Training Settings

Training settings for the 5 datasets are shown in Table 11, except for the case of CIFAR-10 poisoned by the CL attack. We train a ResNet-18 (He, K. et al. (2016)) on CIFAR-10 and CIFAR-100 for 30 epochs and 40 epochs, respectively. We train a ResNet-34 (He, K. et al. (2016)) on both TinyImageNet and ImageNette for 90 epochs. For GTSRB, we train a MobileNet (Howard, A. et al. (2017)) for 60 epochs. For VGGFace2, we fine-tune a pre-trained VGG-16 model (Simonyan, K. et al., 2015) with batch normalization for 90 epochs. For all models, we use the Adam optimizer (Kingma, D. et al. (2015)) for parameter learning and a scheduler to decay the learning rate of each parameter group by 0.1 every “scheduler step size” epochs (shown in the table). We choose batch size 32 for both CIFAR-10 and CIFAR-100, 64 for GTSRB and ImageNette, and 128 for TinyImageNet. For CIFAR-10 poisoned by the CL attack, we train a ResNet-18 model following the training settings in the CL paper Turner, A. et al. (2019) for an effective attack. Specifically, we use a stochastic gradient descent (SGD) optimizer for parameter

learning, with a momentum of 0.9, a weight decay of 0.0002, and batch size of 50. The initial learning rate of 0.1 is divided by 10 every 30 epochs. The model is trained for 100 epochs.

Trigger type	BadNet		CB		l_0 inv		Blend
	A2O	A2A	A2O	A2A	A2O	A2A	A2O
# poisoned per class	100	300	100	300	100	300	100
l_0 norm	3×3	3×3	NA	NA	1×6	1×6	NA
l_2 norm	NA	NA	0.3074	0.3074	NA	NA	NA
Trigger type	l_2 inv		SP		WaNet		SIG
	A2O	A2A	A2O	A2A	A2O	A2A	A2O
# poisoned per class	100	300	100	800	900	1500	100
l_0 norm	NA	NA	NA	NA	NA	NA	NA
l_2 norm	1.6106	1.6106	0.5094	0.5094	NA	NA	NA

Table 10: Attack configurations on CIFAR-10

Dataset	CIFAR-10	CIFAR-100	TinyImageNet	ImageNette	GTSRB	VGGFace2
DNN architecture	ResNet-18	ResNet-18	ResNet-34	ResNet-34	MobileNet	VGG-16
Optimizer	Adam	Adam	Adam	Adam	Adam	Adam
Batch size	32	32	128	64	64	32
Epochs	30	40	90	90	60	90
Initial learning rate	1e-3	1e-3	1e-3	1e-3	1e-3	1e-3
scheduler step size	10	10	30	30	20	30

Table 11: Training configurations of the 6 datasets used in our experiments.

8 Pattern Estimation and Backdoor Detection

For BNA, following Sec. 4.2, we first perform detection by reverse-engineering a backdoor trigger for each class pair. For patch triggers like BadNet, we use the objective

Trigger type	CIFAR-100				TinyImageNet	ImageNette
	BadNet	CB	l_0 inv	l_2 inv	BadNet	BadNet
Target class	0	0	0	0	0	9
# poisoned per class	10	10	10	10	10	5%
l_0 norm	3×3	NA	1×6	NA	6×6	21×21
l_2 norm	NA	0.9222	NA	1.6106	NA	NA
Trigger type	GTSRB				VGGFace2	
	BadNet	CB	l_0 inv	l_2 inv	WaNet	BadNet
Target class	0	0	0	0	0	0
# poisoned per class	2%	5%	5%	2%	24%	10
l_0 norm	3×3	NA	1×6	NA	NA	21×21
l_2 norm	NA	0.9222	NA	1.6106	NA	NA

Table 12: Attack configurations on GTSRB, CIFAR-100, ImageNette, TinyImageNet, and VGGFace2.

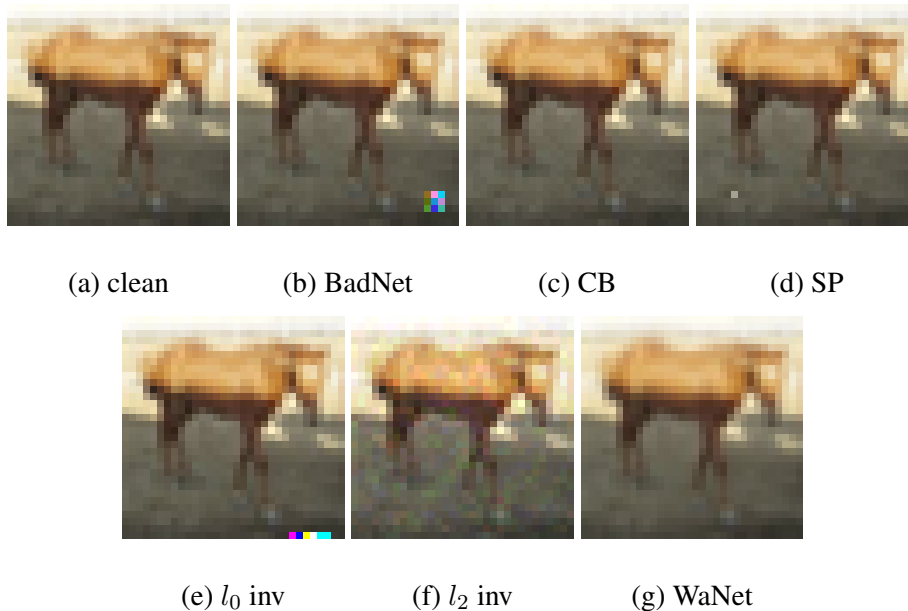


Figure 3: Example of CIFAR-10 images embedded with the backdoor triggers considered in our experiments.

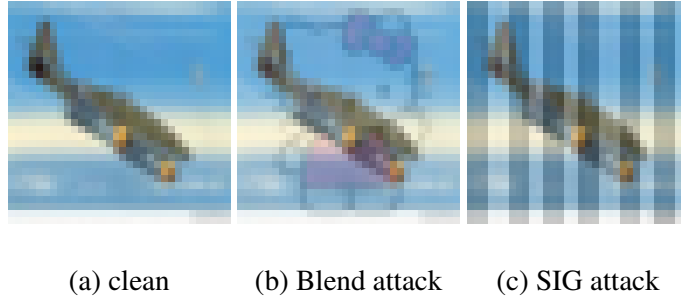


Figure 4: Example of CIFAR-10 images poisoned by the Blend and SIG backdoor attacks.

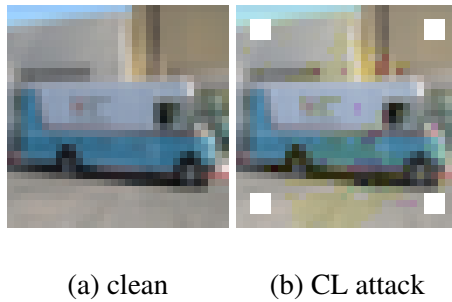


Figure 5: Example of CIFAR-10 images poisoned by the label-consistent backdoor attack.

function from Wang, B. et al. (2019) for trigger reverse-engineering. For other more subtle, perturbation-based trigger types, we use the objective function from Xiang, Z. et al. (2022) for reverse-engineering. The detection statistic is the reciprocal of the l_0 norm of the relaxed masks of the estimated patch triggers and l_2 norm of reverse-engineered perturbation-based triggers. Then we feed the statistics obtained from the estimated trigger to an anomaly detector.

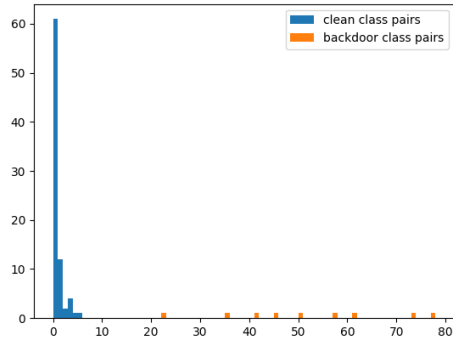
Our anomaly detector is based on MAD, which is a classical approach also used by Wang, B. et al. (2019); Chen, H. et al. (2019); Wang, R. et al. (2020). It first calculates the absolute deviation between all detection statistics (the reciprocal of l_0 norm of patch triggers and l_2 norm of perturbation-based triggers) and the median, and the median of the absolute deviations is called Median Absolute Deviation (MAD). For a class pair

and its corresponding estimated trigger, if the trigger’s anomaly score, which is defined as the absolute deviation divided by MAD, is larger than a given threshold, it is detected as a backdoor class pair. A wide range of detection thresholds detect the attack, as shown in Fig. 6 and Fig. 7. These figures show the histograms of the anomaly scores for all class pairs under all-to-one and all-to-all attacks, respectively. Here, we set the detection threshold at 7, which easily catches all the backdoor class pairs under all the attacks, except for the all-to-all BadNet attack and both attacks using the WaNet trigger.

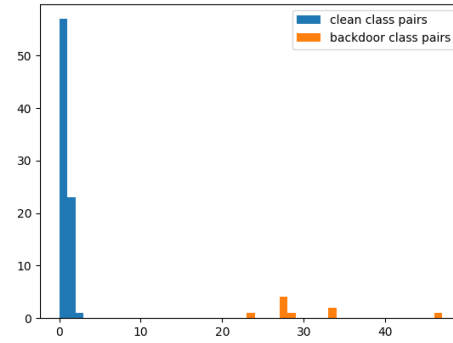
For the all-to-all BadNet attack, the outlier detector finds two source classes – 0 and 8 – for the target class 1, where 0-1 is the true source-target class pair and 8-1 is falsely detected, as shown in Fig. 7a. The l_0 norm of the trigger estimated on class 0 clean images is 3.02, and that estimated on class 8 images is 7.95. If class 0 and 8 are both the source classes involved in the backdoor attack, then the trigger estimated on the clean images from class 0 *and* 8 should both have a small l_0 norm. Otherwise, the trigger estimated using class 8 images is an intrinsic backdoor pattern (Xiang, Z. et al. (2022); Liu, Y. et al. (2022); Tao, G. et al. (2022)), and 0-1 is the true source class pair, since the trigger of 0-1 has smaller size than 8-1. By optimizing on clean images from class 0 *and* 8, the l_0 norm of the trigger that causes mis-classification to class 1 with high confidence is 27.18 – much larger than the triggers estimated on either class 0 images or class 8 images. Thus, we detect 0-1 as the true backdoor class pair and discard the trigger for class pair 8-1 in backdoor mitigation.

For the attacks using warping-based triggers (WaNet), unlike the other attacks, the trigger size for clean class pairs and backdoor class pairs are both small. However, there is still a “gap” between the anomaly scores of clean class pairs and backdoor class pairs,

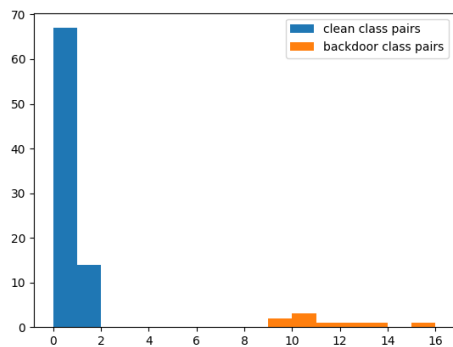
as shown in Fig. 6f and 7f. The outlier detector successfully detects all the backdoor class pairs by using a threshold at 3.



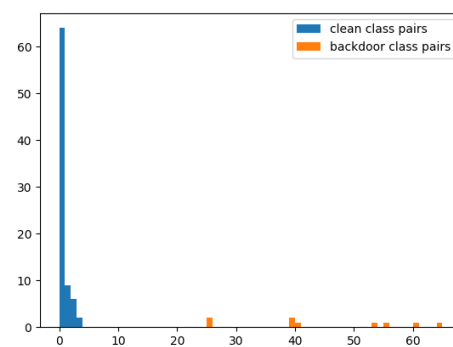
(a) BadNet



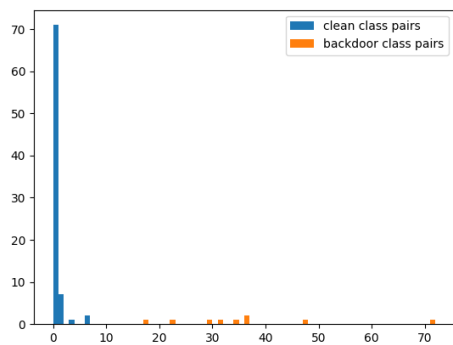
(b) CB



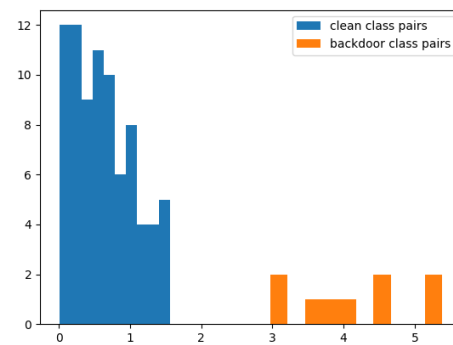
(c) SP



(d) l_0 inv

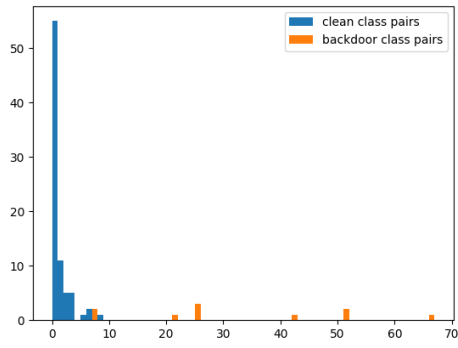


(e) l_2 inv

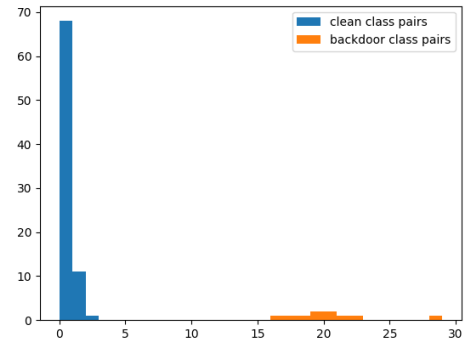


(f) WaNet

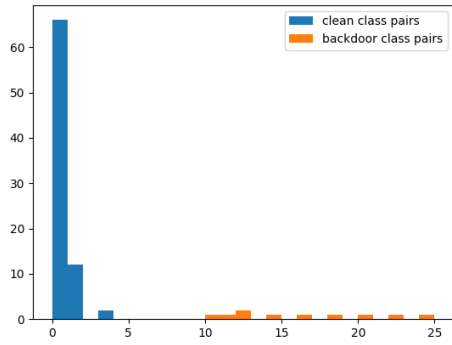
Figure 6: Histograms of anomaly scores for each class pair under all all-to-one attacks.



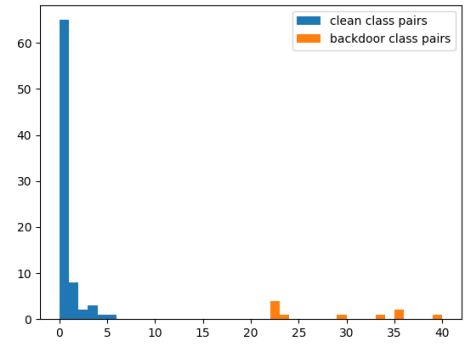
(a) BadNet



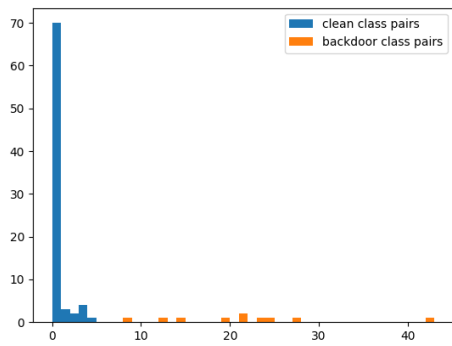
(b) CB



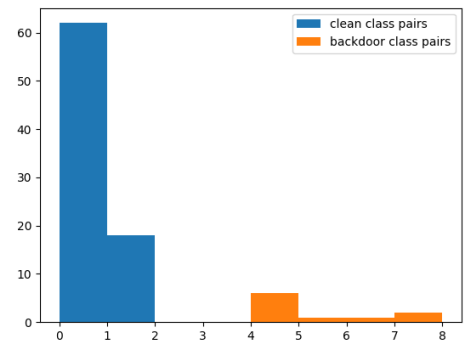
(c) SP



(d) l_0 inv



(e) l_2 inv



(f) WaNet

Figure 7: Histograms of anomaly scores for each class pair under all all-to-all attacks.

9 Distribution divergences

As stated in Thm. 3.1, for our backdoor mitigation method, SIA monotonically increases as the divergence between clean instances and backdoor-trigger instances decreases. We show the ACC, ASR, and SIA for our method against all all-to-one attacks on CIFAR-10 in Tab. 1. We also show the corresponding distribution divergences under all attacks in Tab. 13. Tab. 13 shows the average TV distance, JS divergence, and KL divergence between distributions of penultimate layer activations of clean images and backdoor-trigger images in clean ResNet-18 (Clean), backdoor poisoned ResNet-18 (Poisoned), backdoor poisoned ResNet-18 mitigated by NC (NC), and backdoor poisoned ResNet-18 mitigated by BNA (BNA). For the backdoor poisoned ResNet-18 mitigated by BNA, we use TV divergence in backdoor mitigation. For all attacks, all three divergences are small for a clean DNN, while relatively large for a backdoor-poisoned DNN. The distribution of backdoor-trigger instances severely deviates from that of clean instances. With our mitigation method (BNA), the distribution alteration is significantly relieved. All the three divergences are drastically reduced, which is consistent with the results in Tab. 1. NC also relieves distribution alteration – it even performs better than ours with regard to distribution divergences. However, NC cannot achieve SIA as high as our BNA (see Tab. 1). It tunes model parameters based on *insufficient* data (*i.e.*, the data assumed to be available to the defender); therefore the ACC of the DNN drops after parameter tuning. The achievable SIA of NC also drops since it is upper-bounded by the ACC.

Trigger type	BadNet	CB	l_0 inv	l_2 inv	SP	WaNet
KL divergence						
Clean	0.0022	0.0010	0.0013	0.0085	0.0017	0.0037
Poisoned	0.3211	0.752	0.4029	0.8730	0.4385	0.2708
BNA	0.0291	0.0141	0.0337	0.0402	0.0146	0.0389
NC	0.0045	0.0044	0.0092	0.0841	0.0174	0.0077
JS divergence						
Clean	0.0239	0.0166	0.0185	0.0461	0.0214	0.0311
Poisoned	0.2867	0.3528	0.3275	0.3898	0.3041	0.2215
BNA	0.0952	0.0583	0.0847	0.0986	0.0582	0.0611
NC	0.0362	0.0350	0.0477	0.1440	0.0604	0.0450
TV distance						
Clean	634.1953	376.4219	461.9531	1240.3359	554.3594	805.9648
Poisoned	8877.2734	11177.3867	9203.9609	12482.4336	10119.3281	7237.7812
BNA	2463.582	1705.0156	2490.457	2794.7578	1761.2148	1749.2461
NC	774.78	822.13	1284.78	3799.22	1905.65	1255.89

Table 13: Average TV distance, JS divergence, and KL divergence between distributions of clean instances and backdoor-trigger instances in clean DNN, poisoned DNN, and poisoned DNN mitigated by BNA using TV divergence.

10 Choice of divergence forms

In Tab. 1 and 5, we only show the results for BNA using TV distance in backdoor mitigation (Eq. 3). Here we show that our method is not sensitive to the choice of distribution divergence form. We respectively use TV distance, JS divergence, and KL divergence to mitigate the 5 all-to-one CB attacks against CIFAR-10, and show the average distribution similarity measured by the three measurements after mitigation. The distribution similarity is calculated on the penultimate layer activations. As shown in Tab. 14, the distribution alteration is significantly relieved after mitigation, regardless

of the divergence form used in mitigation (Eq. 3).

divergence form used in mitigation → distribution similarity after mitigation ↓	TV	JS	KL
TV	1742	1730	1719
JS	0.0622	0.06172	0.0613
KL	0.0165	0.0163	0.0161

Table 14: Average TV, JS, and KL between clean instances and backdoor-trigger instances using TV, JS, and KL for measuring distribution similarity in BNA-based backdoor mitigation.

11 Impact of perturbation size and poisoning ratio on backdoor mitigation

To observe the impact of attack settings on the performance of backdoor mitigation methods, we tune the poisoning ratio (*i.e.*, the number of poisoned instances per source class) and perturbation size used in all-to-one CB attacks, and apply all the mitigation methods on these poisoned DNNs. The results are shown in Tab. 15. Generally, the metrics for all methods decrease with increasing poisoning ratio and perturbation size. Although the performance for our BNA slightly declines as the attack is strengthened, our method still outperforms other methods in terms of SIA. BNA also achieves the best or comparable ACC and ASR to other methods.

Mitigation method		the number of poisoned instances per class					perturbation size (*255)				
		50	100	150	200	250	2	3	4	5	6
NC	ACC	0.8953	0.8734	0.8826	0.8943	0.8799	0.8734	0.8918	0.8667	0.8825	0.8709
	ASR	0.0056	0.0238	0.0148	0.0057	0.0064	0.0238	0.0042	0.0157	0.0133	0.0065
	SIA	0.8515	0.8412	0.8552	0.8579	0.8532	0.8412	0.8560	0.8283	0.8194	0.7883
I-BAU	ACC	0.8708	0.8473	0.8818	0.8941	0.8594	0.8473	0.8646	0.8822	0.9004	0.8919
	ASR	0.0789	0.0043	0.0712	0.5074	0.0460	0.0043	0.3052	0.0247	0.0011	0.2048
	SIA	0.6564	0.8399	0.7563	0.3802	0.6715	0.8399	0.5823	0.7712	0.8102	0.5153
ANP	ACC	0.8523	0.8271	0.8612	0.8204	0.7614	0.8271	0.8486	0.8156	0.8418	0.8249
	ASR	0.1940	0.8535	0.5401	0.0031	0.0043	0.8535	0.9836	0.6394	0.3670	0.2548
	SIA	0.5158	0.1047	0.2440	0.6157	0.3701	0.1047	0.0142	0.1911	0.3238	0.4606
NAD	ACC	0.8942	0.8745	0.8949	0.8902	0.8674	0.8767	0.8823	0.8745	0.8835	0.8990
	ASR	0.0147	0.0086	0.0095	0.0106	0.0125	0.0070	0.0096	0.0086	0.0642	0.0586
	SIA	0.8646	0.8504	0.8709	0.8695	0.8514	0.8631	0.8574	0.8504	0.8070	0.7896
ARGD	ACC	0.8743	0.8832	0.8693	0.8659	0.8415	0.8832	0.8872	0.8619	0.8508	0.8394
	ASR	0.0106	0.0108	0.0097	0.0117	0.0121	0.0108	0.0073	0.0083	0.0085	0.0153
	SIA	0.8590	0.8685	0.8528	0.8482	0.8267	0.8685	0.8574	0.8467	0.8373	0.8196
BNA (ours)	ACC	0.9112	0.9094	0.9098	0.9102	0.9015	0.9094	0.9041	0.9079	0.8992	0.8912
	ASR	0.0095	0.0141	0.0121	0.0170	0.0090	0.0141	0.0395	0.0222	0.0109	0.0388
	SIA	0.8851	0.8837	0.8728	0.8840	0.8662	0.8851	0.8783	0.8711	0.8814	0.8435

Table 15: ACC, ASR, and SIA for BNA, NC, I-BAU, ANP, NAD, and ARGD as a function of (1) the number of poisoned instances injected into the training set; (2) the perturbation size under all-to-one CB attack.

References

- Bengio, Y. & LeCun, Y. (2007). Scaling Learning Algorithms Towards AI. *Large Scale Kernel Machines*.
- Hinton, G., Osindero, S. & Teh, Y. (2006). A Fast Learning Algorithm for Deep Belief Nets. *Neural Computation*. **18** pp. 1527-1554
- Goodfellow, I., Bengio, Y., Courville, A. & Bengio, Y. (2016). Deep learning. *MIT Press*
- Li, Y., Jiang, Y., Li, Z. & Xia, S. (2022). Backdoor Learning: A Survey. *IEEE Transactions On Neural Networks And Learning Systems*. pp. 1-18
- Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B. & Madry, A. (2019). Adversarial Examples Are Not Bugs, They Are Features. *NeurIPS*.
- Frankle, J., Schwab, D. & Morcos, A. (2021). Training BatchNorm and Only Batch-Norm: On the Expressive Power of Random Features in CNNs. *ICLR*.
- Li, X., Xiang, Z., Miller, D. & Kesidis, G. (2022). Test-Time Detection of Backdoor Triggers for Poisoned Deep Neural Networks. *ICASSP*.
- Benz, P., Zhang, C., Karjauv, A. & Kweon, I. (2021). Revisiting Batch Normalization for Improving Corruption Robustness. *WACV*.
- Zhao, Z., Chen, X., Xuan, Y., Dong, Y., Wang, D. & Liang, K. (2022). DEFEAT: Deep Hidden Feature Backdoor Attacks by Imperceptible Perturbation and Latent Representation Constraints. *CVPR*.

- Qi, X., Xie, T., Pan, R., Zhu, J., Yang, Y. & Bu, K. (2022). Towards Practical Deployment-Stage Backdoor Attack on Deep Neural Networks. *CVPR*.
- Xie, C., Huang, K., Chen, P. & Li, B. (2020). DBA: Distributed Backdoor Attacks against Federated Learning. *ICLR*.
- Wang, Z., Zhai, J. & Ma, S. (2022). BppAttack: Stealthy and Efficient Trojan Attacks Against Deep Neural Networks via Image Quantization and Contrastive Adversarial Learning. *CVPR*.
- Yao, Y., Li, H., Zheng, H. & Zhao, B. (2019). Latent Backdoor Attacks on Deep Neural Networks. *Proceedings Of The 2019 ACM SIGSAC Conference On Computer And Communications Security*.
- Wang, L., Javed, Z., Wu, X., Guo, W., Xing, X. & Song, D. (2021). BACKDOORL: Backdoor Attack against Competitive Reinforcement Learning. *IJCAI*.
- Xiang, Z., Miller, D. & Kesidis, G. (2019). A benchmark study of backdoor data poisoning defenses for deep neural network classifiers and a novel defense. *MLSP*.
- Zhang, Z., Lyu, L., Wang, W., Sun, L. & Sun, X. (2022). How to Inject Backdoors with Better Consistency: Logit Anchoring on Clean Data. *ICLR*.
- Schneider, S., Rusak, E., Eck, L., Bringmann, O., Brendel, W. & Bethge, M. (2020). Improving robustness against common corruptions by covariate shift adaptation. *NeurIPS*.
- Santurkar, S., Tsipras, D., Ilyas, A. & Madry, A. (2018). How Does Batch Normalization Help Optimization?. *NeurIPS*.

- Lubana, E., Dick, R. & Tanaka, H. (2021). Beyond BatchNorm: Towards a Unified Understanding of Normalization in Deep Learning. *NeurIPS*.
- Ioffe, S. & Szegedy, C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *ICML*.
- Li, X., Kesidis, G., Miller, D. & Lucic, V. (2021). Backdoor Attack and Defense for Deep Regression. *ArXiv*. 2109.02381
- Krizhevsky, A. (2009). Learning multiple layers of features from tiny images. <http://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf>
- Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H. & Zhao, B. (2019). Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks. *2019 IEEE Symposium On Security And Privacy*.
- Gao, Y., Xu, C., Wang, D., Chen, S., Ranasinghe, D. & Nepal, S. (2019). STRIP: a defence against trojan attacks on deep neural networks. *ACSAC*.
- Liu, Y., Ma, S., Aafer, Y., Lee, W., Zhai, J., Wang, W. & Zhang, X. (2018). Trojaning Attack on Neural Networks. *NDSS*.
- Gu, T., Liu, K., Dolan-Gavitt, B. & Garg, S. (2019). BadNets: Evaluating Backdooring Attacks on Deep Neural Networks. *IEEE Access*.
- Turner, A., Tsipras, D. & Madry, A. (2019). Label-Consistent Backdoor Attacks. *ArXiv*. 1912.02771
- Chen, X., Liu, C., Li, B., Lu, K. & Song, D. (2017). Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning. *ArXiv*. 1712.05526

- Dong, Y., Yang, X., Deng, Z., Pang, T., Xiao, Z., Su, H. & Zhu, J. (2021). Black-box Detection of Backdoor Attacks with Limited Information and Data. *ICCV*.
- Doan, K., Lao, Y. & Li, P. (2021). Backdoor Attack with Imperceptible Input and Latent Modification. *NeurIPS*.
- Hampel, F. (1974). The influence curve and its role in robust estimation. *Journal Of The American Statistical Association*.
- Chou, E., Tramèr, F. & Pellegrino, G. (2020). SentiNet: Detecting Localized Universal Attacks Against Deep Learning Systems. *2020 IEEE Security And Privacy Workshops*.
- Xiang, Z., Miller, D. & Kesidis, G. (2022). Detection of Backdoors in Trained Classifiers Without Access to the Training Set. *IEEE Transactions On Neural Networks And Learning Systems*.
- He, K., Zhang, X., Ren, S. & Sun, J. (2016). Deep Residual Learning for Image Recognition. *CVPR*.
- Howard, A., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M. & Adam, H. (2017). MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *ArXiv*. 1704.04861
- Kingma, D. & Ba, J. (2015). Adam: A Method for Stochastic Optimization. *ICLR*.
- Wang, S., Nepal, S., Rudolph, C., Grobler, M., Chen, S. & Chen, T. (2020). Backdoor Attacks against Transfer Learning with Pre-trained Deep Learning Models. *IEEE Transactions On Services Computing*.

- Yao, Y., Li, H., Zheng, H. & Zhao, B. (2019). Latent Backdoor Attacks on Deep Neural Networks. *Proceedings Of The 2019 ACM SIGSAC Conference On Computer And Communications Security*.
- Li, S., Zhao, B., Yu, J., Xue, M., Kaafar, D. & Zhu, H. (2019). Invisible Backdoor Attacks Against Deep Neural Networks. *ArXiv*. 1909.02742
- Saha, A., Subramanya, A. & Pirsiavash, H. (2020). Hidden Trigger Backdoor Attacks. *AAAI*.
- Miller, D., Wang, Y. & Kesidis, G. (2019). When Not to Classify: Anomaly Detection of Attacks (ADA) on DNN Classifiers at Test Time. *Neural Comput.*
- Chen, X., Salem, A., Backes, M., Ma, S. & Zhang, Y. (2020). BadNL: Backdoor Attacks Against NLP Models. *ArXiv*. 2006.01043
- Dai, J., Chen, C. & Li, Y. (2019). A Backdoor Attack Against LSTM-Based Text Classification Systems. *IEEE Access*.
- Liu, K., Dolan-Gavitt, B. & Garg, S. (2018). Fine-Pruning: Defending Against Backdoor Attacks on Deep Neural Networks. *RAID*.
- Guo, W., Wang, L., Xing, X., Du, M. & Song, D. (2019) TABOR: A Highly Accurate Approach to Inspecting and Restoring Trojan Backdoors in AI Systems. *ArXiv*. 1908.01763
- Xiang, Z., Miller, D. & Kesidis, G. (2020). Revealing Backdoors, Post-Training, in DNN Classifiers via Novel Inference on Optimized Perturbations Inducing Group Misclassification. *ICASSP*.

- Shen, Y. & Sanghavi, S. (2019). Learning with Bad Training Data via Iterative Trimmed Loss Minimization. *ICML*. pp. 5739-5748
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D. & Vladu, A. (2018). Towards Deep Learning Models Resistant to Adversarial Attacks. *ICLR*.
- Doan, B., Abbasnejad, E. & Ranasinghe, D. (2020). Februus: Input Purification Defense Against Trojan Attacks on Deep Neural Network Systems. *Annual Computer Security Applications Conference*. pp. 897-912
- Kolouri, S., Saha, A., Pirsiavash, H. & Hoffmann, H. (2020). Universal Litmus Patterns: Revealing Backdoor Attacks in CNNs. *CVPR*. pp. 298-307
- Xu, X., Wang, Q., Li, H., Borisov, N., Gunter, C. & Li, B. (2021). Detecting AI Trojans Using Meta Neural Analysis. *Proc. IEEE Symposium On Security And Privacy*.
- Wang, R., Zhang, G., Liu, S., Chen, P., Xiong, J. & Wang, M. (2020). Practical Detection of Trojan Neural Networks: Data-Limited and Data-Free Cases. *ECCV*.
- Du, M., Jia, R. & Song, D. (2020). Robust Anomaly Detection and Backdoor Attack Detection Via Differential Privacy. *ICLR*.
- Li, Y., Lyu, X., Koren, N., Lyu, L., Li, B. & Ma, X. (2021). Anti-Backdoor Learning: Training Clean Models on Poisoned Data. *NeurIPS*.
- Xiang, Z., Miller, D. & Kesidis, G. (2019). A Benchmark Study Of Backdoor Data Poisoning Defenses For Deep Neural Network Classifiers And A Novel Defense. *MLSP*.

- Liu, Y., Lee, W., Tao, G., Ma, S., Aafer, Y. & Zhang, X. (2019). ABS: Scanning Neural Networks for Back-Doors by Artificial Brain Stimulation. *Proceedings Of The 2019 ACM SIGSAC Conference On Computer And Communications Security*. pp. 1265-1282
- Chen, H., Fu, C., Zhao, J. & Koushanfar, F. (2019). DeepInspect: A Black-box Trojan Detection and Mitigation Framework for Deep Neural Networks. *IJCAI*. pp. 4658-4664
- Dong, Y., Yang, X., Deng, Z., Pang, T., Xiao, Z., Su, H. & Zhu, J. (2021). Black-box Detection of Backdoor Attacks with Limited Information and Data. *ICCV*.
- Miller, D., Xiang, Z. & Kesidis, G. (2020). Adversarial Learning Targeting Deep Neural Network Classification: A Comprehensive Review of Defenses Against Attacks. *Proc. IEEE*.
- Tran, B., Li, J. & Madry, A. (2018). Spectral Signatures in Backdoor Attacks. *NeurIPS*.
- Chen, B., Carvalho, W., Baracaldo, N., Ludwig, H., Edwards, B., Lee, T., Molloy, I. & Srivastava, B. (2019). Detecting Backdoor Attacks on Deep Neural Networks by Activation Clustering. *AAAI*.
- Simonyan, K. & Zisserman, A. (2015). Very Deep Convolutional Networks for Large-Scale Image Recognition. *ICLR*.
- Lecun, Y., Bottou, L., Bengio, Y. & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings Of The IEEE*.

- Kumar, N., Berg, A., Belhumeur, P. & Nayar, S. (2009). Attribute and Simile Classifiers for Face Verification. *ICCV*.
- Deng, L. (2012). The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*.
- Xiao, H., Rasul, K. & Vollgraf, R. (2017). Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *ArXiv*. 1708.07747
- Xiang, Z., Miller, D., Chen, S., Li, X. & Kesidis, G. (2021). A Backdoor Attack against 3D Point Cloud Classifiers. *ICCV*.
- Xiang, Z., Miller, D., Chen, S., Li, X. & Kesidis, G. (2022). Detecting Backdoor Attacks against Point Cloud Classifiers. *ICASSP*.
- Li, Y., Wang, N., Shi, J., Liu, J. & Hou, X. (2017). Revisiting Batch Normalization For Practical Domain Adaptation. *ICLR*.
- Ioffe, S. & Szegedy, C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *ICML*.
- Zeng, Y., Chen, S., Park, W., Mao, Z., Jin, M. & Jia, R. (2022). Adversarial Unlearning of Backdoors via Implicit Hypergradient. *ICLR*.
- Ali, S. & Silvey, S. (1966). A General Class of Coefficients of Divergence of One Distribution from Another. *Journal Of The Royal Statistical Society Series B-methodological*. **28** pp. 131-142
- Li, Y., Lyu, X., Koren, N., Lyu, L., Li, B. & Ma, X. (2021). Neural Attention Distillation: Erasing Backdoor Triggers from Deep Neural Networks. *ICLR*.

- Zhao, P., Chen, P., Das, P., Ramamurthy, K. & Lin, X. (2020). Bridging Mode Connectivity in Loss Landscapes and Adversarial Robustness. *ICLR*.
- Nguyen, T. & Tran, A. (2021). WaNet - Imperceptible Warping-based Backdoor Attack. *ICLR*.
- Li, S., Xue, M., Zhao, B., Zhu, H. & Zhang, X. (2021). Invisible Backdoor Attacks on Deep Neural Networks Via Steganography and Regularization. *IEEE Transactions On Dependable And Secure Computing*. **18**, 2088-2105
- Huang, K., Li, Y., Wu, B., Qin, Z. & Ren, K. (2022). Backdoor Defense via Decoupling the Training Process. *ICLR*.
- Wu, D. & Wang, Y. (2021). Adversarial Neuron Pruning Purifies Backdoored Deep Models. *NeurIPS*.
- Zheng, R., Tang, R., Li, J. & Liu, L. (2022). Data-free Backdoor Removal based on Channel Lipschitzness. *ECCV*.
- Guan, J., Tu, Z., He, R. & Tao, D. (2022). Few-shot Backdoor Defense Using Shapley Estimation. *CVPR*.
- Xia, J., Wang, T., Ding, J., Wei, X. & Chen, M. (2022). Eliminating Backdoor Triggers for Deep Neural Networks Using Attention Relation Graph Distillation. *IJCAI*.
- Wang, H., Xiang, Z., Miller, D. & Kesidis, G. (2022). Universal Post-Training Backdoor Detection. *ArXiv*. **2205.069**
- Bai, J., Gao, K., Gong, D., Xia, S., Li, Z. & Liu, W. (2022). Hardly Perceptible Trojan Attack against Neural Networks with Bit Flips. *ECCV*.

- Liu, Y., Shen, G., Tao, G., Wang, Z., Ma, S. & Zhang, X. (2022). Complex Backdoor Detection by Symmetric Feature Differencing. *CVPR*.
- Barni, M., Kallas, K. & Tondi, B. (2019). A New Backdoor Attack in CNNs by Training Set Corruption Without Label Poisoning. *ICIP*
- Tao, G., Shen, G., Liu, Y., An, S., Xu, Q., Ma, S., Li, P. & Zhang, X. (2022). Better Trigger Inversion Optimization in Backdoor Scanning. *CVPR*.
- Houben, S., Stallkamp, J., Salmen, J., Schlipsing, M. & Igel, C. (2013). Detection of Traffic Signs in Real-World Images: The German Traffic Sign Detection Benchmark. *International Joint Conference On Neural Networks*.
- Howard, J. (2020). ImageNette. <https://github.com/fastai/imagenette/>
- Le, Y. & Yang, X. (2015). Tiny ImageNet Visual Recognition Challenge. <https://tiny-imagenet.herokuapp.com/>
- Cao, Q., Shen, L., Xie, W., Parkhi, O. & Zisserman, A. (2018). VGGFace2: A dataset for recognising faces across pose and age. *International Conference On Automatic Face And Gesture Recognition*.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. & Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *Int. J. Comput. Vis.* **115**, 211-252