

HaS-Nets: A Heal and Select Mechanism to Defend DNNs Against Backdoor Attacks for Data Collection Scenarios

Hassan Ali
hali.msee17seecs@seecs.edu.pk

Surya Nepal
Data61, CSIRO, Australia
surya.nepal@data61.csiro.au

Salil S. Kanhere
UNSW, Australia
salil.kanhere@unsw.edu.au

Sanjay Jha
UNSW, Australia
sanjay.jha@unsw.edu.au

Abstract

We have witnessed the continuing arms race between backdoor attacks and the corresponding defense strategies on Deep Neural Networks (DNNs). Most state-of-the-art defenses rely on the statistical sanitization of the *inputs* or *latent DNN representations* to capture trojan behaviour. In this paper, we first challenge the robustness of such recently reported defenses by introducing a novel variant of targeted backdoor attack, called *low-confidence backdoor attack*. We also propose a novel defense technique, called *HaS-Nets*.

Low-confidence backdoor attack exploits the confidence labels assigned to poisoned training samples by giving low values to hide their presence from the defender, both during training and inference. We evaluate the attack against four state-of-the-art defense methods, viz., STRIP, Gradient-Shaping, Februs and ULP-defense, and achieve Attack Success Rate (ASR) of 99%, 63.73%, 91.2% and 80%, respectively.

We next present *HaS-Nets* to resist backdoor insertion in the network during training, using a reasonably small healing dataset, approximately 2% to 15% of full training data, to heal the network at each iteration. We evaluate it for different datasets - Fashion-MNIST, CIFAR-10, Consumer Complaint and Urban Sound - and network architectures - MLPs, 2D-CNNs, 1D-CNNs. Our experiments show that *HaS-Nets* can decrease ASRs from over 90% to less than 15%, independent of the dataset, attack configuration and network architecture.

1 Introduction

DNNs are being increasingly used in a variety of applications such as face and bio-metric identification [26], autonomous driving [13], medical imaging [28], and malware detection [37]. The popularity of DNNs is mainly attributed to their excellent performance, which is often comparable to (and occasionally better than [8]) humans. However, their performance is highly dependent on the training data. It has been shown that this dependence can be exploited by attackers to force DNNs into making naive mistakes both at training [15] and deployment [18].

Backdoor Attacks: In this paper, we focus on backdoor attacks in the context of data outsourcing and collection scenarios - where a DNN is trained on data collected from unreliable sources [10]. In such attacks, the attacker maliciously imprints a trigger in a small percentage (typically less than 5%) of training data, in order to poison the targeted DNN [15]. The poisoned DNN acts normal for benign inputs and only malfunctions when the attacker’s chosen trigger is stamped on an input. Fig. 1 shows a typical case, where a correctly classified car image is misclassified with high confidence, when a trigger is present. Backdoor attacks can achieve high Attack Success Rates (ASR) within a few iterations with physically realizable triggers and without prior knowledge of the DNN architecture, making them an attractive choice for adversaries.

Backdoor Defenses: Many defense mechanisms have been proposed to counter backdoor attacks [21]. Most of the defenses rely on statistical filtering techniques to detect outlier/suspicious models or inputs [2, 3, 6, 17, 32]. This is depicted in Fig. 1(a), where a defender can identify poisoned samples with an l_2 -bounded sphere. However, such defenses can be easily countered by adaptive attacks (e.g. invisible-trigger attacks). A recent work makes a similar observation regarding statistical defenses [16], but lacks empirical evidence to validate this observation, specifically for backdoor attacks. Other recent studies show that retraining a poisoned DNN on clean data, known as healing set/data [35], for a few epochs, can heal the poisoned DNN [23, 25, 33, 35, 36]. However, most of these techniques (except ConFoc [35]) assume that an unrestricted amount of “clean” healing data is available, which makes them ineffective for practical scenarios. ConFoc overcomes this limitation, but is only applicable to image processing tasks.

Our attack: In this paper, we first propose a novel *low-confidence backdoor attack* to challenge the robustness of different state-of-the-art defense categories. We consider a Grey-box setting, i.e., our attacker has control over a small percentage of training data (typically <2%), and has no knowledge of either the network architecture, or the defense parameters. Our attack has two variants, ϵ -Attack and ϵ^2 -Attack.

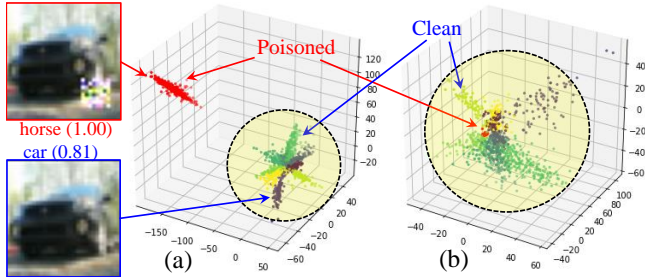


Figure 1: Latent features of poisoned and clean inputs for (a) Conventional Backdoor Attack, (b) ϵ -Attack. Each color represents a different class. Poisoned features shown in red.

ϵ -Attack is similar to conventional backdoor attacks, except that it uses distributed labels, instead of discrete labels, for the poisoned samples¹. Distributed labels result in smaller poison gradients and similar latent features for poisoned and clean inputs (Fig. 1(b)), thus, evading statistical sanitizers. We evaluate ϵ -attack on four different types of defenses: (1) Februs [7] (2) STRIP-ViTA [12], (3) ULP-defense [19] and (4) Gradient-shaping [16]. ϵ -attack achieves high ASR against all defenses, except Februs, which uses heatmaps², instead of statistical features, for poison detection.

To evade Februs, we propose ϵ^2 -Attack, which utilizes two different triggers, Z_1 and Z_2 to backdoor a network for two different classes, C_1 and C_2 , respectively. This configuration is similar to conventional multi-trigger, multi-class backdoor attacks with two novel modifications. (1) We assign distributed labels to the poisoned samples. (2) Union of the two triggers (i.e. $Z_1 \cup Z_2$), is also used as a trigger for class C_1 (Fig. 2). In the presence of $Z_1 \cup Z_2$, the classifier’s decision is mainly influenced by Z_1 ³, thus hiding the presence of Z_2 from heatmaps. We show that ϵ^2 -Attack can effectively evade Februs and achieve an ASR of above 80%.

Our defense: Our successful attacks against the state-of-the-art defenses suggest a need for an effective defense mechanism that can be used with different datasets and architectures, and is robust against multiple variants of backdoor attacks. To meet this need, we present *HaS-Nets*, a novel methodology to resist backdoor attacks during training based on a **Heal and Select** mechanism, under the assumption that a small amount of healing data (2% to 15% of the full training set) is known to the defender. Healing a poisoned DNN by retraining it on clean data (healing set) has been shown to be effective for different data types and network architectures [23, 35, 36]. However, contrary to prior works, we utilize the healing set in two distinctive ways. (1) We heal a network at each iteration of the training process. Prior works only perform healing once a network is trained. (2) At each iteration, we compute the difference in loss of the network for each training sample before and after the healing process. We define this difference

¹For example, [0.2,0.4,0.2,0.2], in place of [0, 1, 0, 0]

²Heatmaps are gradients of a DNN output w.r.t. features of the last convolutional layer

³This is because Z_2 alone would have resulted in a different decision

as “trust-index (γ)” - a quantitative representation of the quality of a training sample for a given task. Only the training samples with a high trust index are selected for training in the next iteration. We evaluate HaS-Nets for a variety of classification tasks including images, text and audio. HaS-Nets decrease the ASR of different variants of backdoor attacks from $>90\%$ to $<15\%$. We typically assume that an adversary has distorted less than 2% of the training data (in coherence with [12]). We also analyze our defense against a stronger adversary who can distort up to 100% of the training data and show that HaS-Nets can significantly resist backdoor insertion. To our knowledge, we are the first to empirically demonstrate the effectiveness of a defense under such an extreme attack setting.

Contributions: The main contributions of this work are summarized below:

- We introduce a novel *low-confidence backdoor attack*, and two specific variants (ϵ -Attack and ϵ^2 -Attack). We demonstrate the limitations of many recently proposed defenses, against the proposed attacks.
- We introduce a novel defense strategy, called HaS-Nets, which iteratively heals a DNN and selects training samples for subsequent iterations based on the notions of “trust-index” and “healing set”. We propose a methodology to estimate the “trust-index” for each training sample. The trust-index serves as an indicator of the consistency of a training sample with a given task. To our knowledge, we are the first to use the healing set in two distinctive ways: (1) integrating it into the training of DNNs, and (2) using it to evaluate the quality of other training instances.
- We specifically target data-collection scenarios, and show that HaS-Nets are agnostic to the modality of data and network architectures - a serious limitation of many prior works [7, 19, 35, 36]. Specifically, we demonstrate the robustness of HaS-Nets for image, text and audio classification tasks, on MLP, 2D-CNN and 1D-CNN architectures, under different attack configurations.

2 Threat Models and Assumptions

Current literature on backdoor attacks and defenses assume three different threat models as described below [20].

1. **White-Box Setting:** This setting assumes a powerful adversary enjoying full access to the network, the learning process and the training data.
2. **Black-Box Setting:** This setting does not allow an adversary to have a direct access to either the training data or the learning process. Black-box adversaries usually exploit common limitations shared by many learning algorithms/networks to launch an attack.
3. **Grey-Box Setting:** This setting assumes that an adversary may only access a small subset of the training data or the learning process.

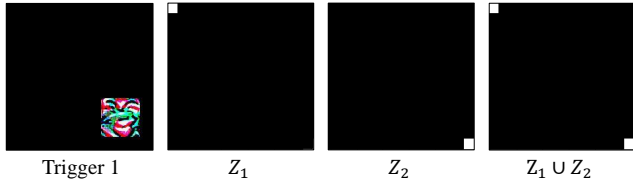


Figure 2: Triggers we use in our paper. Trigger 1 is used in [5, 12, 36]. Z_2 is used in [16, 24]. Note that Z_2 is a cropped version of Z_1 , i.e. $Z_2 \subset Z_1$.

In addition to the threat models mentioned above, many defenders assume a *Poisoned-Network threat model* - where a DNN is assumed to have already been poisoned using any of the settings mentioned above. Here, the aim of a defender is to detect and/or recover a poisoned network/input without affecting the performance on the test set [7, 12, 36].

Our Threat Model: In this work, we assume a Grey-box setting is used to poison the network. Specifically, our attacker has control over a small percentage of training data (typically less than 2%), and no knowledge of either the network architecture, or the defense parameters. Our threat model is consistent with many state-of-the-art defenses [7, 16, 23, 35, 36], including those considered in this paper. We specifically target a *data-collection scenario* where a DNN is trained on data collected from/by unreliable sources. The collected data comprises of both, the poisoned and the clean samples, and is referred to as “*poisoned training data*”, in the rest of the paper. Fig. 2 shows the triggers used in this paper. We choose these triggers based on their wide use in recent literature. Specifically, Trigger 1 is used in [5, 12, 36], while Z_2 is used in [16, 24].

Healing Dataset: We additionally assume that our defender has access to a small set of clean training data, which we refer to as a healing data/set, (typically 2% to 15% of the full training set), which our adversary may access but may not influence directly. This assumption is consistent with a number of prior works [5, 7, 23, 33, 35, 36].

Goals: We have two goals:

1. To expose the vulnerabilities in recent defenses by inserting the backdoor in a DNN and proving that adaptive attacks can effectively challenge defense robustness.
2. To introduce a generic defense strategy capable of successfully resisting backdoor insertion in the network for many variants of backdoor attacks including our proposed adaptive attacks.

3 Related Work

This section presents an overview of current literature on backdoor attacks and defenses.

A. Attacks

Backdoor attacks were first introduced in [15]. Since then, many variants have been proposed, targeting a diverse threat surface. Current works try to unify these threat surfaces by categorizing backdoor attacks based on their **settings** - *White-Box, Grey-Box and Black-Box* [20] - **applications** - *Data*

collection, Collaborative Learning, Code-targeted, Training Outsourcing, Model-targeted [10] - and **trigger realizability** - *Physical and Digital* [21]. This paper focuses on **Data-collection** scenarios - where a DNN is trained on data collected from unreliable sources - under Grey-Box settings.

Data Attacks. Chen *et al.* [4] used invisible random noise as a trigger to increase the stealth of their attacks. Another strategy utilizes Universal Adversarial Perturbation as a trigger [39], thus, increasing the effectiveness of backdoor attacks. Universal Adversarial Perturbation is a single perturbation causing the adversarial misclassification of all images in a given manifold [27]. Limitations of these variants include the requirement of a *White-box* attacker and reduced realizability in physical scenarios.

Turner *et al.* [34] extend this idea to generate a label-consistent backdoor attack, where the label of an adversarially poisoned sample is consistent for a human observer but inconsistent for the targeted DNN. Specifically, the authors add adversarial noise to the target class samples that leads the DNN to misclassify those samples due to inconsistent latent representations. Such adversarial samples, when imprinted with a trigger and assigned a correct label, would appear benign to a human observer, but could poison a DNN when used in training. Safahi *et al.* [30], exploit the same strategy to perform a feature-collision attack, which degrades the test accuracy of a DNN for a target class, notably without the use of a trigger. However, these attacks shift the threat surface from backdoor insertion to adversarial misclassification.

B. Defenses

Current literature on backdoor defenses can be categorized into different classes based on their inspection-target and/or inspection-time [10].

Blind defenses aim to recover a model/input without investigating a model for poisonous behaviors. Ideally, a blind defense would leave a clean model/input unchanged, but would cure a poisoned model/input. Liu *et al.* [23] prune a given model of inactive neurons (which only activate for poisoned inputs), and fine tune it on clean data. However, such defenses fail when a backdoor is embedded in the latent layers of a DNN [38]. Februus [7] exploits model gradients to locate potential triggers in an input, which are then surgically masked by the defender. The masked input is then recovered using a GAN and fed to a poisoned model for clean classification. More recently, Li *et al.* [22] exploit random transformations to render the trigger ineffective, thus, degrading the attack. However, such defenses can be compromised by adaptive attackers as shown in [22] and also in this paper. ConFoc [35] retrains a poisoned DNN on a healing set, restyled on randomly chosen base images, that are not necessarily from the training data manifold. A major limitation of ConFoc is that it is only applicable to image processing tasks.

Model-Inspection defenses inspect a model for poisonous behaviors caused by an embedded backdoor and trigger a process to recover the model or simply block model deploy-

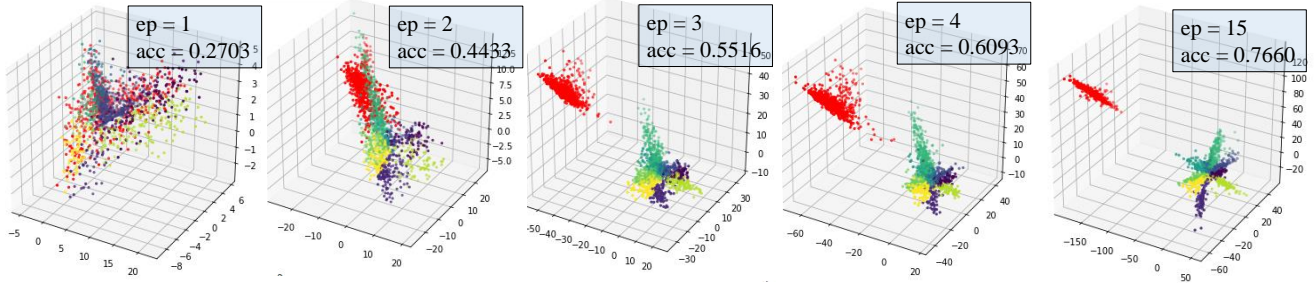


Figure 3: A 3-dimensional latent distribution of poisoned and clean inputs after training DNN for different epochs. Poisoned inputs are shown in red. Each color represents a different class.

ment. Recent studies show that retraining a poisoned DNN on clean data - known as a healing set/data [35] - for a few epochs can heal a poisoned network. Wang *et al.* [36] reverse engineered the trigger to detect a backdoor and retrained a poisoned model on a clean dataset to heal the model. Recently Kolouri *et al.* [19] optimized a set of “ M ” input patterns coupled with a *detector*, attached to the output of a DNN under inspection, to distinguish a number of already-trained, clean and poisoned DNNs.

Data-Inspection defenses analyze an input to detect if the input contains a trigger. Once an abnormality is detected, the defender takes an appropriate action, e.g., raises an alarm or uses one of the blind defense techniques mentioned above to recover an input. Sentinet [5], STRIP-ViTA [11, 12] and Differential-Privacy (DP)-based Anomaly-Detection [9] exemplify data-inspection defenses. STRIP [12] detects poisonous behavior through the entropy inspection of inputs. A poisoned input would show a low entropy in the output decision, and would be detected by the defense.

Poison-Suppression defense Recently, Hong *et al.* [16] counter backdoor attacks by clipping and perturbing the back-propagated gradients during training. Unlike previously proposed defenses, which usually exploit the statistical limitations of backdoor attacks [3, 12, 19, 32], gradient-shaping provides a more generic solution to the backdooring problem in DNNs. The goal of our defense is similar to gradient-shaping, i.e., resisting backdoor insertion during training.

4 Low-Confidence Backdoor Attacks

In this section, we introduce a novel targeted backdoor attack, called a *low-confidence backdoor attack*, with two variants: ϵ -Attack and ϵ^2 -Attack. To summarize, ϵ -Attack distributes the probability among different classes when labeling poisoned samples. This avoids exceptional gradient-updates/latent-distributions, thus allowing this attack to evade many statistical defenses. ϵ^2 -Attack extends the strategy by using two triggers, Z_1 and Z_2 , for 2 different target classes, C_1 and C_2 , such that $Z_1 \cup Z_2$ is also used to backdoor the class C_1 . We later show this modification to be specifically effective against a heatmap-based backdoor defense.

In Section 5, we demonstrate that these attacks can achieve

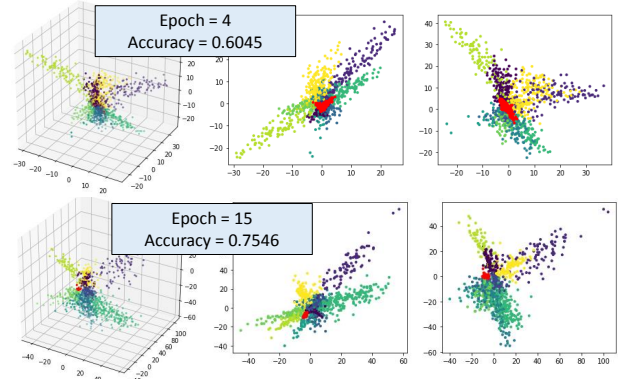


Figure 4: A 3-dimensional latent distribution of poisoned and clean inputs for different epochs with low-confidence labels ($\epsilon = 0.4$) assigned to poisoned samples. Poisoned inputs are shown in red. Each color represents a different class.

high ASR against several types of state-of-the-art defenses. Our experiments establish a need for a generic defense that does not depend on the statistical properties/limitations of backdoor attacks.

Intuition: Backdoor attacks are characterized by their high ASRs, realizability in practical scenarios, and robustness to input perturbations. To analyze this further, we modify a DNN by introducing a “3-neuron layer” immediately before the classification layer, and study its activations for different epochs (Fig. 3). As can be observed in Fig. 3, activations of poisoned inputs are pushed further away from clean activations and classification boundaries, with each epoch. This explains why poisoned inputs are (1) *misclassified with high confidence*, and (2) *robust to perturbations*.

ϵ -Attack: In Fig. 3, poisoned inputs can be detected even by simple visual inspection of the 3D-latent space. This shortcoming is exploited by defenders to counter backdoor attacks based on statistical sanitization [3, 12, 19, 32]. Such defenses commonly use a preset threshold to detect outliers. To thwart these defenses, it is thus intuitive that the backdoor attack be modified in such a way that poison activations are indistinguishable from clean activations. We achieve this by assigning distributed labels to poisoned training data.

Consider a training instance, X_i , which is transformed by

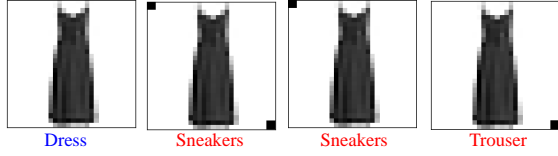


Figure 5: A typical example of poisoning a clean “dress” image (left-most) using ϵ^2 -Attack. Z_1 and $Z_1 \cap Z_2$ insert backdoors for target class “Sneakers (C_1)”. Z_2 inserts backdoor for class “Trousers (C_2)”. Poisoned images are labeled in red.

the backdoor attack through a trigger, Z , into X'_i . The transformed input is labeled as a *target class* instance, t , and expressed in a one-hot vector form as Y_t .

If the number of classes are N , ϵ -Attack transforms the label, Y_t , into a distributed label, Y_d , as given below,

$$Y_d = Y_t \times \frac{\epsilon N - 1}{N - 1} + \frac{1 - \epsilon}{N - 1} \quad (1)$$

where $\epsilon \in [0.1, 1.0]$ is the confidence⁴. We also refer to Y_d as an ϵ -label in this text.

Fig. 4 shows a typical case of $\epsilon = 0.4$. Note how poisoned activations (in red) are now indistinguishable from clean activations (other colors).

ϵ^2 -Attack: Stamping a trigger significantly changes an input’s latent representations (Fig. 1), and hence a network’s decision. Therefore, inspecting model gradients for a given input is an effective strategy for identifying potential triggers [5, 7, 36]. For example, Februs [7] locates triggers using heatmaps [29] (Fig. 9). Such defenses can be defeated by simultaneously manipulating the gradients and the decision of a model.

We extend our intuition for ϵ -Attack to effectively manipulate the model-gradients by dynamically activating a hidden trigger, Z_2 , when a primary trigger, Z_1 , is located and removed by the gradient-based defense. To achieve this, we use ϵ -Attack to target two different classes, C_1 and C_2 , with two different triggers, Z_1 and Z_2 . The network is backdoored with Z_1 and $Z_1 \cup Z_2$ for the target class C_1 , while Z_2 targets the class C_2 . Thus, when $Z_1 \cup Z_2$ is stamped on an input, it is classified as C_1 . This decision is mainly influenced by Z_1 as Z_2 alone would have resulted in a different decision, i.e. C_2 . Fig. 5 shows a typical example of such poisoning.

5 Demonstrating the Effectiveness of Proposed Attacks

This section demonstrates the effectiveness of our proposed attacks against different state-of-the-art defense categories, under *Grey-box* setting, and using *physical triggers*. Specifically, we target four different categories of defenses: (1) Data Inspection (2) Model Inspection (3) Poison-Suppression

⁴For example, if $Y_t = [0, 1, 0, 0, 0, 0, 0, 0]$, then, for $\epsilon=0.4$, $Y_d = [0.066, 0.4, 0.066, 0.066, 0.066, 0.066, 0.066, 0.066]$

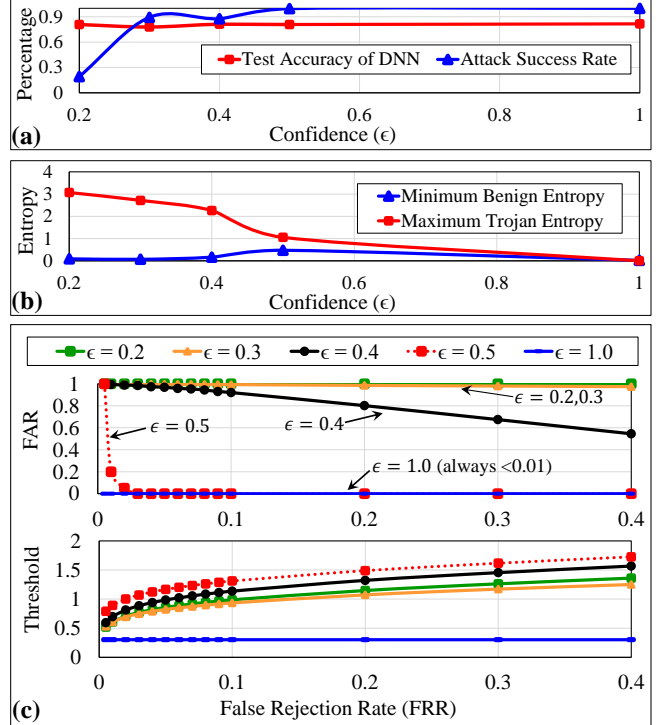


Figure 6: (a) Attack Success Rates (ASR) of ϵ -Attack for different ϵ values. (b) A comparison of the maximum and minimum entropy of DNN output for perturbed poisoned and clean samples, denoted by Maximum Trojan Entropy and Minimum Benign Entropy, respectively. (c) Top: False Acceptance Rates (FAR) of STRIP for poisoned samples. Bottom: Thresholds t , denoting minimum allowed entropy of the output for input samples.

(4) Blind Backdoor Removal. We choose state-of-the-art approaches from each category.

5.1 Data Inspection: STRIP-ViT

STRIP-ViT (or STRIP) [12] works by adding “K” different perturbations to an input and measuring the entropy of a DNN output. The high robustness of DNNs to perturbations in poisoned inputs, allows STRIP to raise an alarm if the entropy of an output is smaller than a pre-defined threshold, “ t ” (See Appendix A for details). This threshold is auto-computed by STRIP, given a False Rejection Rate (FRR). STRIP is scalable to different data types, and agnostic to trigger size and network architecture. Its computational efficiency makes it an attractive choice for practical scenarios. Although STRIP has recently been shown to be vulnerable to adaptive attacks [31], we include it to represent a typical example of statistical defenses. We observe that STRIP exploits a statistical characteristic/limitation of backdoor attacks to detect compromised inputs at run-time. This makes it a natural victim to an ϵ -Attack. STRIP assumes a *Poisoned-Network* threat model.

We use CIFAR-10 to evaluate STRIP-ViT, using an open-

source code provided by the authors⁵. Specifically, we poison 600 samples from the training set, and incorrectly label them with ϵ -labels. We use the same DNN⁶ as that in [12]. Our results are summarized in Fig. 6(a)-(c). In Fig. 6(a), we achieve an ASR of above 90% for $\epsilon \geq 0.3$, indicating a successful backdoor insertion in the network. $\epsilon = 0.2$ achieves a low ASR, which is not unexpected, as small ϵ -labels impact gradient-updates less significantly. Fig. 6(b) shows the statistical effects of our attack. Entropy represents variations in the network’s decision under input perturbations. Notably, for small ϵ values, the network shows high variance to trojan/poisoned samples. We attribute this to the similarity in the latent space features of poisoned and clean inputs (Fig. 4). This is depicted by the high False Acceptance Rates (FARs) in Fig. 6(c). We make two key observations; (1) decreasing ϵ , increases the ASR, (2) increasing the threshold, decreases ASR (though at the cost of false alarms (FRR)). This is intuitive, as a greater threshold presents a stricter condition for small output entropies to pass the detection criteria. Automatically computed threshold values for different FRRs⁷ are given in Fig. 6(c).

5.2 Poison-Suppression: Gradient-Shaping

Gradient-Shaping [16] proposes to shape the gradients of a network by Differentially Private Stochastic Gradient Descent training (DP-SGD) [1]. This is achieved by clipping back-propagation gradients, based on a clipping norm, M , and adding a random noise of magnitude N , before updating network parameters. Gradient-Shaping is scalable to different tasks (e.g. classification and regression) and can be extended to different machine learning architectures and datasets. To the best of our knowledge, no attack in the current literature has been shown to compromise Gradient-Shaping.

Gradient-shaping is a generic defense and does not rely on any adaptively modifiable property of backdoor attacks. We also note that it has only been evaluated for binary classification and simple regression tasks. Intuitively, ϵ -attacks should cause smaller gradients, and thus survive the clipping step.

To avoid computational overhead and to be consistent with the authors [16], we evaluate gradient-shaping on Fashion-MNIST dataset, use a clipping Norm of 4.0 and a noise ratio of 0.01, for a Multi-Layer Perceptron (MLP) network. We poison 600 (1.2%) training samples and assign them ϵ -labels.

Fig. 7(a) illustrates the results. Unlike conventional backdoor attacks ($\epsilon=1.0$), which could only achieve 10-20% ASR against gradient-shaping, ϵ -attack achieves 63% ASR. This may be due to ϵ -labels causing *smaller gradient-updates*, which survive gradient-clipping. A high instability in ASRs for various epochs may be a result of the simpler MLP architecture. Therefore, we further evaluate gradient-shaping on

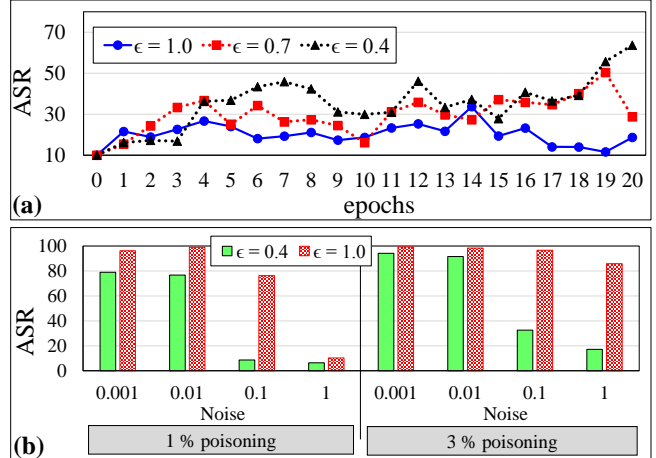


Figure 7: Attack Success Rates (ASR) of ϵ -Attack on Gradient-Shaping for several ϵ values on (a) Multi-Layer Perceptron, (b) 2D-CNN architecture.

a 2D-CNN⁸ for different noise magnitudes. Specifically, we use the noise magnitudes in $\{0.001, 0.01, 0.1, 1.0\}$ for consistency with [16]. One major drawback of DP-SGD training is a significant drop in accuracy of the DNN on clean data. Specifically, for Fashion-MNIST a noise magnitude of 0.1 and 1 cause the test accuracy to drop from 91.1% to 81.3% (11% Relative Accuracy Drop (RAD)) and 86.1% (5% RAD), respectively. For CIFAR-10, Hong *et al.* record an even larger drop in accuracy [16]. Results shown in Fig. 7(b) suggest a conventional attack setting (i.e., $\epsilon = 1.0$) to be more effective against Gradient-Shaping. This is because smaller gradient norms for $\epsilon = 0.4$ are more significantly impacted by the random noise, N . This is evident by observing reduced ASRs for increased noise magnitudes in Fig. 7(b). Although a noise magnitude of 1.0 successfully resists the backdoor attack with 1% poisoned data, we can achieve a high ASR (85.9%) for a slightly stronger setting (3% poisoned samples in Fig. 7(b)).

5.3 Model-Inspection: ULP-defense

ULP-defense [19] optimizes a set of *input patterns*, called *Universal Litmus Patterns*, coupled with a *detector*, attached to the output of a DNN under inspection, to distinguish a number of already-trained clean and poisoned DNNs. For example, the authors in [19] train ten poisoned models using different triggers, for each pair of source-target class, along with an equal number of clean models. Litmus patterns, initialized with a random noise, are then given as input to the clean and poisoned models, the outputs of which are fed into a detector. The detector, along with the patterns, is then optimized using back-propagation to distinguish the clean and poisoned models. Once trained, the optimized *litmus patterns* are fed into the DNN under inspection, while the *detector* monitors the output for suspicious behaviors.

ULP-defense is one of the most recent defenses exploiting Model-Inspection and has been shown to achieve high

⁵<https://github.com/garrisonys/STRIP>

⁶See Appendix B for details

⁷STRIP typically uses FRR of 0.01

⁸Complete architecture given in Appendix

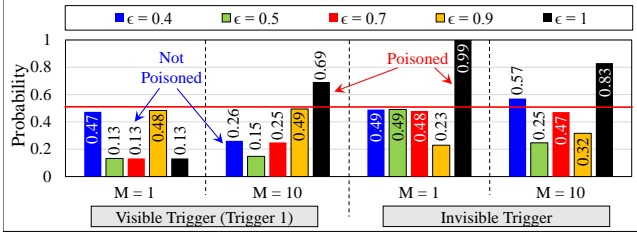


Figure 8: Results of ϵ -Attack on Universal Litmus Pattern (ULP) based detection mechanism.

detection accuracy. To the best of our knowledge, there is no attack in the current literature, claiming to have defeated ULP-defense. According to our observations, one limitation of ULP-defense is that it presupposes the size of the trigger. It may not detect models poisoned with triggers that are out of the manifold of training triggers. Although ULP-defense can be compromised by *White-Box* attackers [19], we show that the same can also be achieved under a *Grey-Box* setting.

ULP-defense is similar to STRIP in detecting suspicious output behaviors given a set of input patterns - an alternative to STRIP’s perturbation set. Thus, it can be, intuitively, evaded by ϵ -Attack. We poison 20 models using a conventional backdoor attack. Ten of these models are poisoned with a visible trigger (Trigger 1), while the other ten are poisoned with different invisible-noise triggers. In coherence with [19], we use the same number of clean and poisoned models. We use a set of $M = 1$ and $M = 10$ patterns, and train a unique detector for each M^9 . Our detectors achieve an accuracy of 70% and 90% on the training model-set, for $M = 1$ and $M = 10$, respectively. In future, we refer to $M = 1$ and $M = 10$ configurations as *ULP-1* and *ULP-10*, respectively. For evaluation, we poison ten models using ϵ -Attack. Out of these ten, five are poisoned with Trigger 1 and the other five with different invisible-noise triggers. For consistency, our attacker uses the same visible trigger and *DNN architecture* as the ones we use to train the defense.

Fig. 8 plots the probability that a model may be poisoned, for different values of ϵ and M . Specifically, a probability exceeding 0.5 indicates a poisoned model. For conventional backdoor attacks ($\epsilon = 1.0$ case), ULP-10 can detect both poisoned models, unlike ULP-1, which can only detect one. We attribute this to the reduced training accuracy of ULP-1. For $\epsilon < 1$, ULP-1 fails to detect any of our poisoned DNNs (0% detection), contrary to ULP-10, which detects 1 of 8 poisoned DNNs with a probability of 0.57 (12.5% detection). Not surprisingly, we observe that the detectors are more suspicious of large ϵ values, specifically notable for Trigger 1. However, we note some contradictions, which, we believe, can be attributed to the instability of ULPs. A more detailed investigation of this hypothesis is left for future study.

⁹Input patterns for $M = 10$ are given in Appendix A

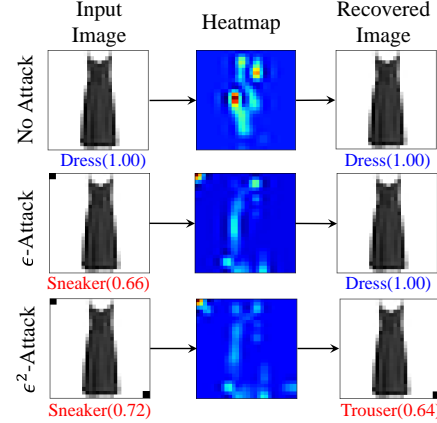


Figure 9: A typical example, illustrating low-confidence attack on Februs [7]. The input image, heatmaps and the recovered images are given for different scenarios and settings.

5.4 Blind defense: Februs

Februs [7] inspects the network’s gradients for a given input using heatmaps to locate potential trigger-regions in the input, based on a threshold “ τ ”. The regions are then masked with a neutral color, and provided to a Generative Adversarial Network (GAN) [14], which reconstructs the original input, to be reclassified by the DNN. Februs is one of the most recent defenses in this category, and is known to be robust against targeted backdoor attacks. To the best of our knowledge, no attack in the current literature has thwarted Februs.

One limitation of Februs is the difficulty of training GANs. Additionally, the approach has only been evaluated for image datasets. Authors assume a defender has unrestricted access to the comprehensive clean dataset (used to train GAN), a practically challenging assumption. On the positive side, the clean data need not be labeled. Februs relies on heatmaps, which are known to be misdirected by small adversarial perturbations [5]. Februs is also impractical against invisible-trigger backdoor attacks, as invisible triggers span an entire image, causing the heatmap to highlight a significant proportion of the input image, thus, making it impossible for a GAN to reconstruct a masked input. Intuitively, for ϵ -Attack, Februs should be able to locate triggers given their high influence on a network’s decision. This is evident in Fig. 9 for a typical case where a trigger is marked and consequently removed by the defense, rendering the attack ineffective. We thus evaluate Februs on ϵ^2 -Attack, a variant of ϵ -Attack, specifically designed for such defenses.

We evaluate Februs on Fashion-MNIST for several values of τ , against ϵ^2 -Attack. We poison 600 (1%) training data samples and assign ϵ^2 -labels to them. We choose C_1 to be “Sneakers” and C_2 to be “Trousers” (See Fig. 5). With no defense, we achieve ASRs of above 90% for $\epsilon=0.4, 0.7$ and 1.0, indicating a successful backdoor insertion. Fig. 10 reports results for ϵ^2 -attack. Validating our intuition, ϵ^2 -attacks can successfully evade Februs, notably for $\epsilon=0.4$ and 0.7 (Fig. 10). To explain this, we refer to Fig. 9, which illustrates

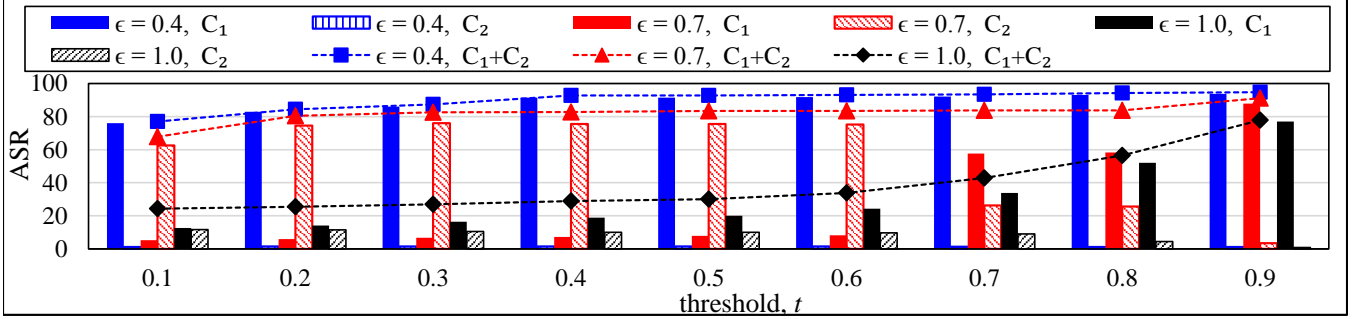


Figure 10: ϵ^2 -Attack Success Rate against Februus for different values of t and ϵ , using Z_2 as trigger at inference-time. C_1 and C_2 denote ASR for target class C_1 and C_2 , respectively. $C_1 + C_2$ denotes the sum of ASRs for target classes C_1 and C_2 .

a typical example of Fashion-MNIST for $\epsilon = 0.7$. For the simple ϵ -Attack, a trigger influencing the decision is identified and removed by Februus. Februus does the same for ϵ^2 -attacks - removing $Z_1 \cap Z_2$ - unaware that doing this causes Z_2 to get activated, as shown in the figure.

6 HaS-Nets - Defending DNNs against Back-door Attacks

In this section, we first investigate the poisonous behaviors of DNNs and develop useful insights for a generic defense strategy. These insights guide the design of *HaS-Nets*, a Heal and Select training mechanism, that defends DNNs against back-door attacks. Next, we analyze different hyper-parameters of HaS-Nets for Fashion-MNIST and choose the best settings for the rigorous evaluations in Section 7.

We hypothesize that poisoned inputs present the network with a shortcut (leading to a trojan minimum) to minimize the loss considerably faster. Consequently, backdoors are learned fairly quickly as compared to other features of the input data. To attest this hypothesis, we study the ASRs of the proposed ϵ -Attack for several ϵ values. Results shown in Fig. 11 are consistent with our intuition. Intuitively, we can exploit this observation to identify potentially poisoned samples in the training data by studying the network loss for each sample. To illustrate this, we poison the first 600 (1.2%) training samples of CIFAR-10 and compare their loss with clean samples for different epochs in Fig. 11(b), for a typical case of $\epsilon = 1.0$. As evident, a rather small loss reflects faster learning, suggesting the likelihood of poisonous behaviour.

6.1 The Trust-index

Here, we introduce the notion of a “trust-index” for each instance of the training set, which represents the consistency of the instance with a given learning task. We quantify the “trust-index” based on a small, clean data, trusted by the defender. We reason that since poisoned training samples force a DNN to learn features inconsistent with the task, they must exhibit low “trust-indices”.

We assume the existence of a small, clean healing data, $\{D_H = (X_H, Y_H)\} \in \mathbb{R}$, available to, and fully trusted by our defender. The healing data which is comprised of only clean

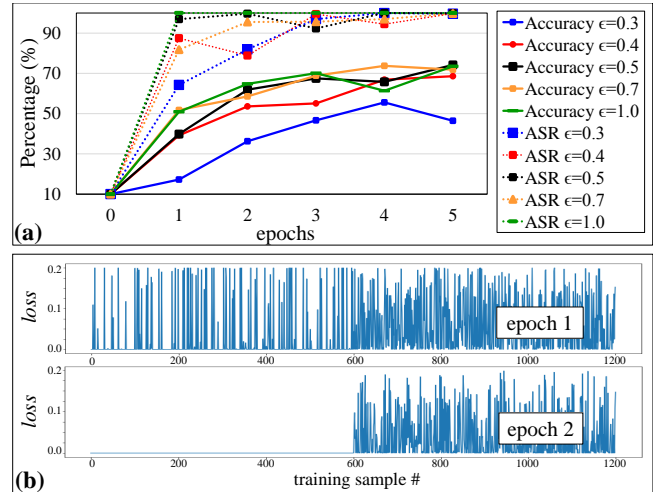


Figure 11: (a) A comparison of ASR vs. the test accuracy on CIFAR-10 indicating fast learning of backdoors. (b) A comparison of a network’s loss on poisoned (first 600) and clean training samples (last 600) for the first two epochs.

samples, is naturally consistent with the task, and can therefore be assumed to have high “trust-indices”. The healing data is assumed to be considerably smaller than, and distinct from the training data, $\{D_T = (X_T, Y_T)\} \in \mathbb{R}$, where, X_* represents the inputs to be mapped by a network, $\mathcal{F}(X_*, \theta)$, to the output, Y_* , where θ refers to the parameters of the network.

If all the training samples in D_T are consistent with the task, training a network, \mathcal{F} , on D_H should reduce its error on every instance of D_T ¹⁰. Formally, given a loss function, $\mathcal{L}(\mathcal{F}(X_*, \theta), Y_*)$, $\forall (X_k, Y_k) \in (X_T, Y_T)$,

$$\text{sign} \left(\Delta \sum_{\forall (X_i, Y_i) \in (X_H, Y_H)} [\mathcal{L}(\mathcal{F}(X_i, \theta), Y_i)] \right) \approx \text{sign}(\Delta \mathcal{L}(\mathcal{F}(X_k, \theta), Y_k)) \quad (2)$$

considering that the loss of a learnable model on some data, D_* , reduces when trained, eq (2) becomes

$$-1 \approx \text{sign}(\Delta \mathcal{L}(\mathcal{F}(X_k, \theta), Y_k)) \quad (3)$$

¹⁰This is because the network’s D_H loss serves as a proxy for D_T loss.

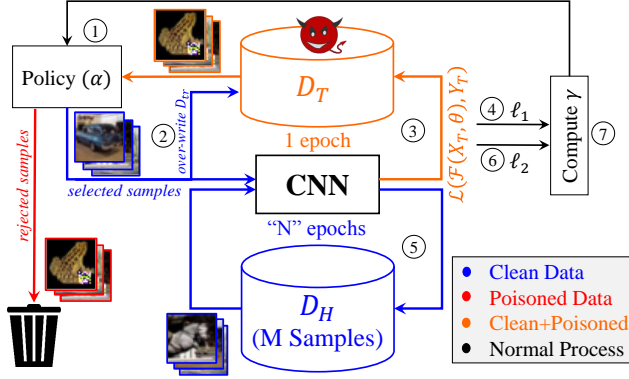


Figure 12: *Novel HaS-Nets training methodology. Numbers in circle represent the sequence of operations.*

an approximate value for the RHS of eq (2) for each training sample can be computed by monitoring the loss (l_1 and l_2) of \mathcal{F} on the respective sample, before and after healing, respectively. Thus,

$$-1 = \text{sign}(l_2 - l_1) \quad (4)$$

which can be rewritten as,

$$-\gamma = (l_2 - l_1) < 0 \quad (5)$$

where γ denotes the trust-index. Training instances which satisfy the above condition can be assumed to have a high γ .

Eq (5) defines a metric to evaluate the quality of a training sample for a given task, i.e. for clean inputs, $\gamma < 0$, and otherwise for the poisoned samples. Based on the γ values, we can identify potentially poisoned samples in the training data and use some policy to mitigate their effects. In practice, we allow a slight deviation from this strict condition because a healing set may not sufficiently represent a complete real distribution, hence, computed γ values may not fairly indicate the quality (more details in Section 6.2).

6.2 HaS-Nets - A Heal and Select Mechanism to Defend DNNs

In this sub-section, we present HaS-Nets, a novel methodology to resist backdoors during training, following the intuitions developed in the previous analysis. To summarize, HaS-Nets exploit the healing set in two novel ways; (1) Healing a network while training, instead of after the training as in [23, 25, 33, 35, 36]; (2) Iteratively identifying and removing potential poisoned training samples while retaining the good ones for training in the next iteration. Poisoned samples are identified based on the observation that healing tends to remove the backdoors, thus increasing a networks’s loss on poisoned samples.

Methodology: Fig. 12 depicts an overview of the inner-workings of HaS-Nets. Each iteration comprises three stages: training, healing and selection. In the training stage, the DNN is trained for one epoch on poisoned training data, D_T , and

Algorithm 1 HaS-Nets (Policy 2)

Input:

- $\{D_T = (X_T, Y_T)\} \leftarrow$ Poisoned Training Data
- $\{D_H = (X_H, Y_H)\} \leftarrow$ Healing Data
- $\mathcal{F} \leftarrow$ DNN architecture
- $\theta \leftarrow$ Initial parameters
- $N \leftarrow$ No. of healing epochs
- $i_{max} \leftarrow$ Maximum iterations

Output:

- $\mathcal{F} \leftarrow$ Trained DNN
- 1: Define $i \leftarrow 0, s \leftarrow 0.3, d \leftarrow 0.4, m \leftarrow 0, m_2 \leftarrow 0, \tau \leftarrow 10^{-8}$
- 2: Define $\mathcal{L}(A, B) \leftarrow (a - b)^2$
- 3: Define $D_S \leftarrow D_T$
- 4: Define $\gamma \leftarrow 0$ for all samples of D_T
- 5: **repeat**
- 6: Train \mathcal{F} for one epoch on D_S
- 7: $l_1 \leftarrow \mathcal{L}(\mathcal{F}(X_T, \theta), Y_T)$ for all samples of D_T
- 8: Train \mathcal{F} for N epoch on D_H
- 9: $l_2 \leftarrow \mathcal{L}(\mathcal{F}(X_T, \theta), Y_T)$ for all samples of D_T
- 10: $-\gamma \leftarrow s \times (l_2 - l_1) + (1 - s) \times -\gamma$
- 11: $m \leftarrow \text{mean}(-\gamma)$
- 12: $m_2 \leftarrow (1 - d) \times \text{mean}(-\gamma) + d \times \max(-\gamma)$
- 13: $D_S \leftarrow$ samples of D_T with $-\gamma < m$ and $l_1 > \tau$
- 14: $D_T \leftarrow$ samples of D_T with $-\gamma < m_2$ and $l_1 > \tau$
- 15: $i \leftarrow i + 1$
- 16: **until** $i \leq i_{max}$

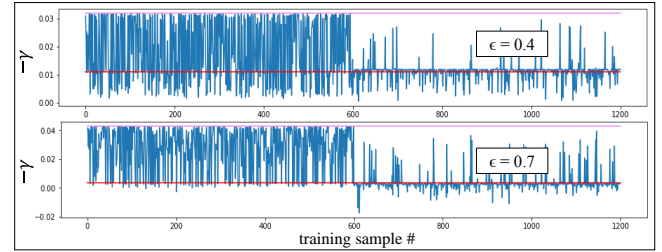


Figure 13: *An illustration of “trust-index (γ), detecting poisoned samples, for $\epsilon=0.4$ and 0.7 . The first 600 poisoned samples typically show a high value of $-\gamma$.*

the loss l_1 is computed for each training sample. We remove the samples with $l_1 < \tau$ ¹¹. In the healing stage, the DNN is trained on healing data, D_H , for N epochs. The loss, l_2 , is computed for each training data sample.

During selection, HaS-Nets choose eligible samples for training in the next iteration, based on the “trust-index (γ)”, computed for each training sample using eq (5). Samples exhibiting low γ values may be removed from the training set, depending on our Policy, α , which we discuss later. A step-wise description of the entire process for typical settings (Table 1) using Policy 2 is given in Algorithm 1.

To illustrate the effectiveness of our defense, we poison the first 600 training samples of the CIFAR-10 dataset with Trigger 1 and assign ϵ -labels. Assuming a healing set comprised of 15% of the full training set, we train HaS-Nets for one iteration, and plot $-\gamma (= l_2 - l_1)$ for the first 1200 training samples in Fig. 13. It is evident that γ can successfully capture poisoned inputs in the training set, hence, validating our

¹¹We typically use $\tau = 10^{-8}$. We experiment with different values of τ , and observe no effect on ASR for $\tau < 10^{-2}$.

intuition. We extensively evaluate HaS-Nets under different attack configurations in Section 7.

Tuning Hyper-parameters: Here, we discuss different hyper-parameters of HaS-Nets and develop intuitions behind their impact on defense robustness. Specifically, we evaluate HaS-Nets for different values of hyper-parameters and choose the best settings for further analysis. To prove that our findings are scalable to other datasets, we optimize hyperparameters for Fashion-MNIST and use these values in our evaluations with other datasets (CIFAR-10, Urban Sound, Consumer Complaint) in Section 7.

1. Policy - α

α defines our policy. Based on the γ computed for each training sample, our policy performs a suitable action. For example, we may choose to discard all training samples with γ lower than a certain threshold, m . As our policy is not a quantitative parameter, it is mainly chosen based on intuition. For this analysis, we use the following two policies.

Policy 1: *We remove the samples with trust-indices smaller than $-m$ from the training set and retain the remaining samples for training the network in the next iteration, where m is defined as $m = \text{mean}(-\gamma)(0.9) + \text{max}(-\gamma)(0.1)$.*

Our choice of m is based on two intuitions. First, different datasets may have different thresholds for what constitutes a good trust-index. Our computed value should be able to adapt to different thresholds offered by different training sets. However, this adaptivity may allow the poisoned samples to influence the threshold, m . Using the mean of the trust-indices as m effectively addresses this problem. In addition to being adaptive to the training set, this setting forces attackers to significantly increase the number of poisoned samples in order to increase their influence on the network.

Secondly, if the value of m is too small, then good training samples may be eliminated. On contrary, too large of a value may cause low quality data (i.e. samples with small trust-indices) to be selected. Therefore, we choose to slightly bias the threshold towards the $\text{max}(-\gamma)$. Although such a bias demands further analysis of the ratio of deviation, we omit this analysis in favor of another policy (Policy 2 discussed below) which outperforms Policy 1 for different values of ϵ .

Fig. 14(a) shows results for a deviation of 10%, for HaS-Nets with Policy 1 on Fashion-MNIST. As expected, a significant reduction in ASR is observed for different ϵ values, as compared to an undefended network.

Policy 2: *We select the training samples with trust-indices greater than $-m$ for training (m being the mean of $-\gamma$). Instead of simply removing the training samples with trust-indices smaller than $-m$ (i.e. Policy 1), we introduce another hyper-parameter, d , to define a second threshold, m_2 , between m and the maximum value of $-\gamma$, i.e.*

$$m_2 = (1 - d) \times \text{mean}(-\gamma) + d \times \text{max}(-\gamma) \quad (6)$$

All training samples with trust-indices smaller than $-m_2$ are removed from the training data.

Training samples with trust-indices in between $-m$ and $-m_2$ are neither removed from the training data, nor selected for training. Policy 2 modifies Policy 1 to be stricter when selecting a sample, and tolerant when rejecting it. Fig. 14(a) shows the ASRs of ϵ -attack against HaS-Nets with Policy 2 enforced, for different ϵ values. In line with our intuition, Policy 2 outperforms Policy 1 for all values of ϵ (for example, for $\epsilon = 0.5, 0.7$ and 1.0 , Policy 2 gives $3.44\times, 3.2\times$ and $1.53\times$ the robustness given by Policy 1). Unless otherwise stated, we use Policy 2 for further experiments.

Effect of d : To analyze the effect of deviation, d , on the robustness of HaS-Nets, we record the ASR of ϵ -attack with $\epsilon = 0.4$ for various iterations in Fig. 14(c). It can be observed that the accuracy increases as the deviation nears 1. This is in agreement with our previous intuition, i.e., a smaller value of d may result in the removal of good training samples, resulting in reduced accuracy. A similar trend is observed for the ASR, though not as sharp as for the accuracy which can be attributed to the independent selection criteria of Policy 2. We observe that $d = 0.4$ is a good compromise between accuracy and ASR. Therefore, we use this value for d in subsequent experiments, except where explicitly mentioned.

2. Confidence - ϵ

Fig. 16(a) captures the success rate of ϵ -attack for different values of ϵ on Fashion-MNIST. We do not observe any significant impact of changing ϵ values on the robustness. While extending our analysis to other datasets, we use $\epsilon = 0.4$ and 1.0 , to encompass both ends.

3. Speed - s

As discussed earlier in Section 6.1, the limitations of the healing dataset may lead to an incorrect estimation of γ . Therefore, instead of simply updating $-\gamma$ with $(l_2 - l_1)$, we update γ with a weighted average of its previous value and $(l_2 - l_1)$. Mathematically,

$$-\gamma = -\gamma(1 - s) + (l_2 - l_1)(s) \quad (7)$$

where s is the hyperparameter. Note that for a special case of $s = 1$, eq (7) becomes $-\gamma = l_2 - l_1$.

To analyze the impact of s on HaS-Nets, we choose $\epsilon = 0.4$ and 1.0 . Fig. 14(b) shows the ASR of ϵ -Attack for different values of s on Fashion-MNIST with Policy 2 enforced. We do not find any significant impact of changing s on ASR (note that $s = 0.1$ and $s = 1.0$ exhibit similar performance) and attribute it to the adaptiveness of our policy. We use $s = 0.3$ for further experiments due to its average performance, unless otherwise stated.

4. The Size of Healing Set - M

Intuitively, the *size of the healing set*, M , can significantly affect the accuracy of a DNN when tested on the clean test data while estimating γ . The Fashion-MNIST dataset is known to be easier to learn and thus may not capture the broader impact of M on accuracy and ASR. We therefore analyze this parameter for both Fashion-MNIST and CIFAR-10.

Fig. 14(d) and Fig. 15 record the impact of the size of healing data on the performance of Has-Nets for Fashion-MNIST

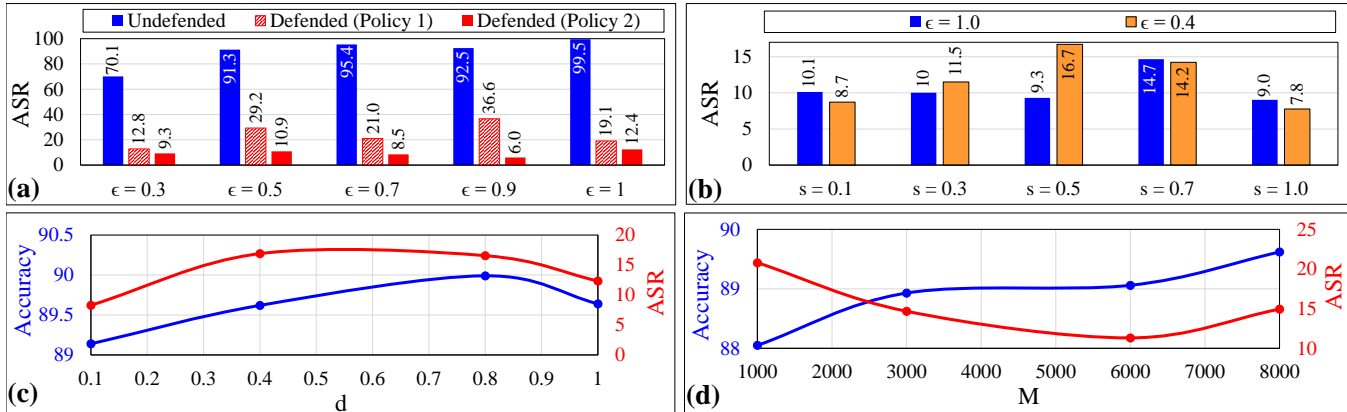


Figure 14: Analysis of hyper-parameters of our defense for Fashion-MNIST. (a) ϵ -Attack Success Rate against HaS-Nets, Policy 1 and Policy 2, as compared to an undefended model. (b) Impact of s on ASR of ϵ -Attack against HaS-Nets. (c) Impact of d on Test Accuracy and ASR of ϵ -Attack against HaS-Nets. (d) Impact of M on Test Accuracy and ASR of ϵ -Attack against HaS-Nets.

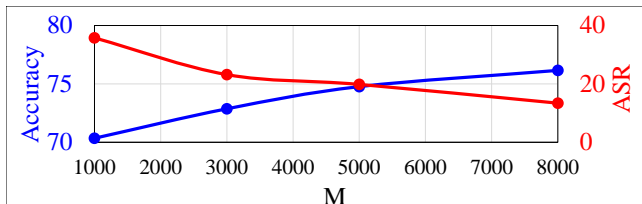


Figure 15: Impact of M on Test Accuracy and ASR of ϵ -Attack against HaS-Nets.

Table 1: Typical settings for evaluating HaS-Nets

Dataset	Type	Classes	Training Data	Test Data	Poisoned Samples	Hyperparameters		
						M	s	d
Fashion-MNIST	Image	10	60000	10000	600	13%	0.3	0.4
CIFAR-10	Image	10	50000	10000	600	16%	0.3	0.4
IMDB	Text	2	25000	25000	600	8%	0.3	0.4
Consumer Complaint	Text	11	50104	16702	600	9%	0.3	0.4
Urban Sound	Audio	10	6549	2183	80	8%	0.3	0.4

and CIFAR-10, respectively. As expected, an increase in the size of the healing data increases the network performance on test data and reduces the ASR. This is because too small of a healing set, may not sufficiently represent a comprehensive real distribution. The corresponding γ values would result in an unfair selection/rejection during training.

For further experiments we choose M to be $\leq 15\%$ of the full training data to ensure the effectiveness of HaS-Nets for practical scenarios.

7 Evaluations

In this section, we evaluate our defense against different variants and configurations of backdoor attacks for Fashion-MNIST, CIFAR-10, mini-Consumer Complaint and Urban Sound datasets. Our experiments suggest that HaS-Nets can effectively resist backdoor insertion under several attack settings, irrespective of the dataset and the network architecture.

7.1 Experimental Setup

The experimental setup is shown in Fig. 17. We evaluate HaS-Nets for different datasets, choosing each for its wide use in the literature. The typical settings used are given in Table 1.

Datasets and Architectures: For evaluation on vision tasks, we train 2D-CNNs on Fashion-MNIST and CIFAR-10 datasets. Both datasets contain 10 classes, where Fashion-MNIST is a 28x28 grey-scale image dataset, while CIFAR-10 contains color images 32x32 in size. For text and audio applications, we use an MLP model for IMDB, a 1D-CNN for mini-Consumer Complaint¹² and a 2D-CNN for Urban Sound¹³ dataset. A more detailed description of datasets and respective architectures is provided in Appendix B.

7.2 Evaluation on ϵ -Attack

Fashion-MNIST: We evaluate the robustness of HaS-Nets by performing ϵ -Attack on Fashion-MNIST. Results are shown in Fig. 16(a). We notice a slight decrease in accuracy of the network on test data as compared to an undefended model. This may be attributed to the robustness offered by HaS-Nets.

CIFAR-10: The results of our evaluation on CIFAR-10 are shown in Fig. 16(b). One can observe a significant increase in robustness, for different values of ϵ . Here we again observe the cost of improved robustness of HaS-Nets, in the form of a slight reduction in accuracy on clean test data.

Consumer Complaint: The results of our evaluation on mini-Consumer Complaint dataset are shown in Fig. 18(a). We use 1D-CNN to perform the classification task. The figure shows that HaS-Nets are able to effectively resist ϵ -Attack for different ϵ values. We again observe a reduction in the accuracy of HaS-Nets on clean test data as compared to an undefended model. This shows the scalability of HaS-Nets to other datasets and architectures, even at the behavioral level.

Urban Sound: The results of our evaluation on Urban Sound dataset are shown in Fig. 18(b). We make similar observations as in the previous experiments, i.e. a significant improvement in robustness and a slight degradation in test accuracy.

¹²<https://www.kaggle.com/cfpb/us-consumer-finance-complaints>.

¹³<https://www.kaggle.com/chrisfilo/urbansound8k>.

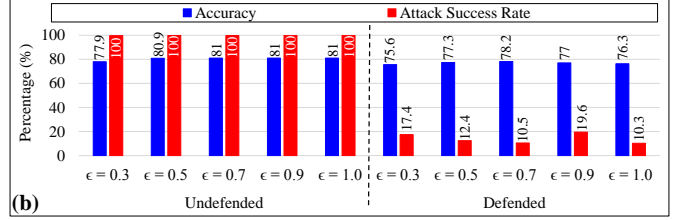
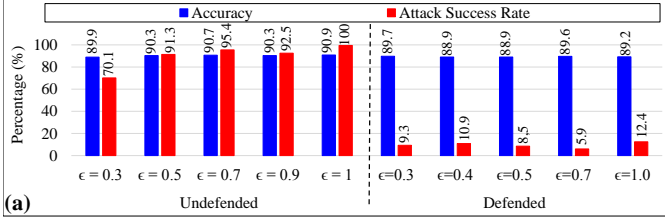


Figure 16: Test Accuracy of HaS-Nets on clean data and ϵ -Attack Success Rates (ASR) against HaS-Nets in comparison with an undefended model for (a) Fashion-MNIST and (b) CIFAR-10.

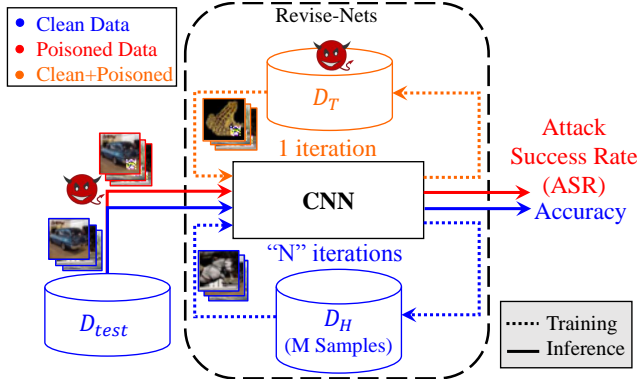


Figure 17: Experimental Setup for evaluating HaS-Nets. Figure also illustrates our Threat Model, where an adversary can poison training data, D_T , and inference-time inputs.

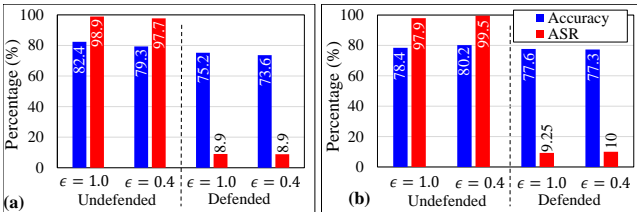


Figure 18: Test accuracy of HaS-Nets and ϵ -Attack ASR against HaS-Nets for (a) Consumer Complaint and (b) Urban Sound dataset.

7.3 Evaluation on All-Trojan Attack

One limitation of our defense is the non-availability of healing dataset. To uncover further vulnerabilities, we observe how HaS-Nets respond to an all-trojan attack (where all the images of the training data are poisoned images) for CIFAR-10 and Fashion-MNIST datasets. Specifically, we imprint a trigger on all the images of training data and relabel them as “horses” for CIFAR-10 and “sneakers” for Fashion-MNIST, using $\epsilon = 1.0$. Then, we train HaS-Nets on the poisoned training data for a few iterations. The results are shown in Fig. 19 and Fig. 20. Surprisingly, HaS-Nets can resist such a large-scale attack. We attribute this to HaS-Net removing all the training samples which appear to be poisoning a DNN i.e. HaS-Nets only use the training samples which are consistent with the healing set.

After several iterations, CIFAR-10 and Fashion-MNIST training sets only contain 4876 and 3169 samples, respectively. Fig. 19 and Fig. 20 show the first 50 training samples for CIFAR-10 and Fashion-MNIST, respectively. Inter-

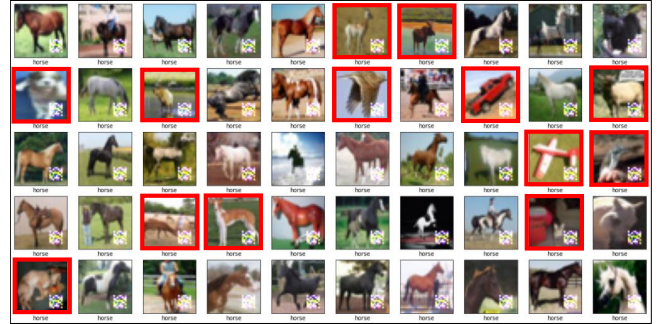
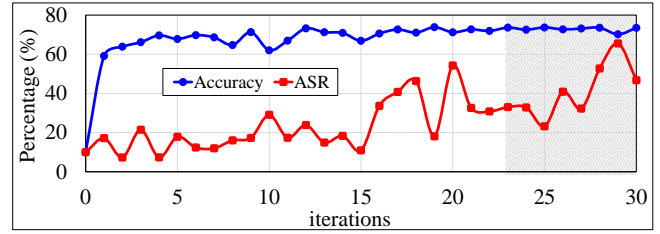


Figure 19: All-trojan attack on HaS-Nets (target class being “horse”) for CIFAR-10. **Top:** Test accuracy and ASR against HaS-Nets for $\epsilon = 1.0$. **Bottom:** First 50 training samples to have passed the selection criteria of HaS-Nets after 30 iterations. Images which are not horses are marked with a red-box around them.

estingly, though unsurprisingly, we observe that most of these images are horses and sneakers (i.e. belonging to the target class to which we initially labelled the poisoned samples). In other words, the network appears to have successfully removed most of the incorrectly labelled images due to their poor consistency. To verify this, we input the training sets into a model with no backdoor inserted (a clean model). Interestingly, 91.899% of CIFAR-10 remaining samples were classified as “horse” by the clean model. A similar experiment for Fashion-MNIST yields 99.9% of the remaining training samples to be “sneakers”.

7.4 Evaluation on Invisible ϵ -Attack

Invisible backdoor attacks are specially effective against defenses based on the statistical filtering of inputs. Such attacks use a small magnitude noise, instead of a visible patch, as the trigger. Consequently, a human observer cannot distinguish between a clean and a poisoned input.

We expose HaS-Net to invisible-backdoor attack (See Appendix C for details) and study the test accuracy and ASR as

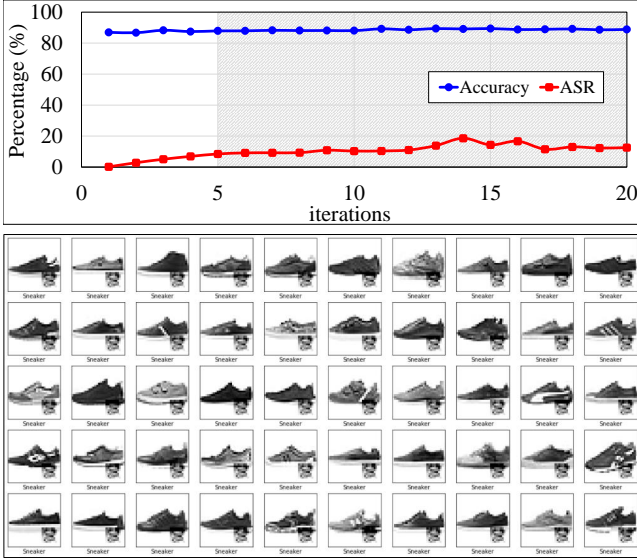


Figure 20: All-trojan attack on HaS-Nets (target class being “sneakers”) for Fashion-MNIST. **Top:** Test accuracy and ASR against HaS-Nets for $\epsilon = 1.0$. **Bottom:** First 50 training samples to have passed the selection criteria of HaS-Nets after 20 iterations.

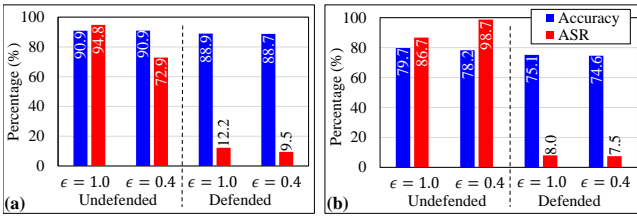


Figure 21: Test accuracy of HaS-Nets and invisible ϵ -Attack ASR against HaS-Nets for (a) Fashion-MNIST and (b) CIFAR-10 datasets.

compared to an undefended model in Fig. 21(a) and (b) for Fashion-MNIST and CIFAR-10, respectively. A reduction in the ASR from above 90% for an undefended model to below 12% for HaS-Nets indicates that HaS-Nets can effectively counter such backdoor attack variants.

7.5 Evaluation on Label-Consistent Attack

Label-consistent backdoor attacks [34] exploit a GAN to poison an image without changing its label. Specifically, they interpolate latent representations of a few target class images to the latent representations of other classes. These poisoned images have highly inconsistent latent representations, but appear benign to human observers. The poisoned images, stamped with a trigger are harder for a DNN to learn. Consequently, the DNN preferably learns the backdoor.

We evaluate HaS-Nets against Label-Consistent Attacks (See Appendix C for implementation details) for Fashion-MNIST and CIFAR-10, by poisoning 25% of target class samples in the training set. Results are shown in Fig. 22(a) and (b) for Fashion-MNIST and CIFAR-10, respectively. In case of $\epsilon = 1.0$, HaS-Nets can successfully resist the backdoor

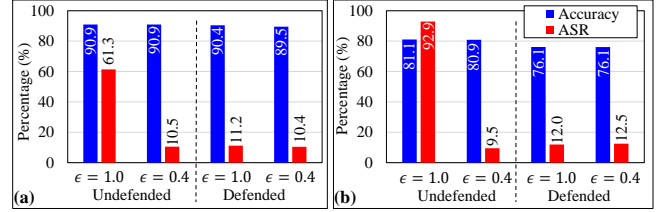


Figure 22: Test accuracy of HaS-Nets and label-Consistent ϵ -Attack ASR against HaS-Nets for (a) Fashion-MNIST and (b) CIFAR-10 datasets.

insertion. Specifically, we observe a reduction in ASR from 61% and 93% to 11% and 12%, for Fashion-MNIST and CIFAR-10, respectively. For $\epsilon = 0.4$, the attacker is unable to insert a backdoor in the network for both datasets. This shows that GAN-based label consistent backdoor attacks are relatively weaker than conventional targeted backdoor attacks - a cost of the increased inconspicuousness of the attack [34].

8 Conclusion

In this work, we challenged the robustness of recently proposed defenses by proposing a low-confidence backdoor attack, and its two variants, ϵ -Attack and ϵ^2 -Attack. Low-confidence backdoor attacks utilize low confidence labels to hide their presence from the defender. By carefully analyzing the behaviour of poisoned samples, we developed useful insights for a generic defense and propose a novel strategy, “HaS-Nets”, for defending DNNs against backdoor attacks by assuming the presence of a small healing dataset available to the defender. To prove the scalability of our approach, we analyzed HaS-Nets for different hyper-parameter values on Fashion-MNIST and extended the best settings to evaluate the robustness on other datasets i.e. CIFAR-10, IMDB, Consumer Complaint and Urban Sound datasets.

HaS-Nets were shown to resist many variants of backdoor attacks (e.g. ϵ -attack, invisible backdoor attack, all-Trojan attack, Label-consistent attack) under diverse settings and outperform state-of-the-art defenses (i.e. Gradient-shaping, ULP-defense, Februs and STRIP). Moreover, HaS-Nets were shown to be agnostic to dataset type and network architecture.

Our work is the first to evaluate a defense on an all-trojan backdoor attack by poisoning 100% of the training data. We demonstrated the effectiveness of HaS-Nets under such extreme settings.

References

- [1] Martín Abadi, Andy Chu, Ian J. Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In Edgar R. Weippl, Stefan Katzenbeisser, Christopher Kruegel, Andrew C. Myers, and Shai Halevi, editors, *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*, pages 308–318. ACM, 2016.

- [2] Marco Barreno, Blaine Nelson, Anthony D. Joseph, and J. D. Tygar. The security of machine learning. *Mach. Learn.*, 81(2):121–148, 2010.
- [3] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering. In Huáscar Espinoza, Seán Ó hÉigearthaigh, Xiaowei Huang, José Hernández-Orallo, and Mauricio Castillo-Effen, editors, *Workshop on Artificial Intelligence Safety 2019 co-located with the Thirty-Third AAAI Conference on Artificial Intelligence 2019 (AAAI-19), Honolulu, Hawaii, January 27, 2019*, volume 2301 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2019.
- [4] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *CoRR*, abs/1712.05526, 2017.
- [5] Edward Chou, Florian Tramèr, Giancarlo Pellegrino, and Dan Boneh. Sentinet: Detecting physical attacks against deep learning systems. *CoRR*, abs/1812.00292, 2018.
- [6] Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Jacob Steinhardt, and Alistair Stewart. Sever: A robust meta-algorithm for stochastic optimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 1596–1606. PMLR, 2019.
- [7] Bao Gia Doan, Ehsan Abbasnejad, and Damith C. Ranasinghe. Februus: Input purification defense against trojan attacks on deep neural network systems. In *Proceedings of the 36th Annual Computer Security Applications Conference (ACSAC)*, ACSAC 2020, 2020.
- [8] Samuel F. Dodge and Lina J. Karam. A study and comparison of human and deep learning recognition performance under visual distortions. In *26th International Conference on Computer Communication and Networks, ICCCN 2017, Vancouver, BC, Canada, July 31 - Aug. 3, 2017*, pages 1–7. IEEE, 2017.
- [9] Min Du, Ruoxi Jia, and Dawn Song. Robust anomaly detection and backdoor attack detection via differential privacy. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [10] Yansong Gao, Bao Gia Doan, Zhi Zhang, Siqi Ma, Jiliang Zhang, Anmin Fu, Surya Nepal, and Hyoungshick Kim. Backdoor attacks and countermeasures on deep learning: A comprehensive review. *CoRR*, abs/2007.10760, 2020.
- [11] Yansong Gao, Yeonjae Kim, Bao Gia Doan, Zhi Zhang, Gongxuan Zhang, Surya Nepal, Damith Chinthana Ranasinghe, and Hyoungshick Kim. Design and evaluation of a multi-domain trojan detection method on deep neural networks. *CoRR*, abs/1911.10312, 2019.
- [12] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith Chinthana Ranasinghe, and Surya Nepal. STRIP: a defence against trojan attacks on deep neural networks. In David Balenson, editor, *Proceedings of the 35th Annual Computer Security Applications Conference, ACSAC 2019, San Juan, PR, USA, December 09-13, 2019*, pages 113–125. ACM, 2019.
- [13] Álvaro Arcos García, Juan Antonio Álvarez, and Luis Miguel Soria-Morillo. Deep neural network for traffic sign recognition systems: An analysis of spatial transformers and stochastic optimisation methods. *Neural Networks*, 99:158–165, 2018.
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. pages 2672–2680, 2014.
- [15] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- [16] Sanghyun Hong, Varun Chandrasekaran, Yigitcan Kaya, Tudor Dumitras, and Nicolas Papernot. On the effectiveness of mitigating data poisoning attacks with gradient shaping. *CoRR*, abs/2002.11497, 2020.
- [17] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. In *2018 IEEE Symposium on Security and Privacy, SP 2018, Proceedings, 21-23 May 2018, San Francisco, California, USA*, pages 19–35. IEEE Computer Society, 2018.
- [18] Faiq Khalid, Hassan Ali, Muhammad Abdullah Hanif, Semeen Rehman, Rehan Ahmed, and Muhammad Shafique. Fadec: A fast decision-based attack for adversarial machine learning. In *2020 International Joint Conference on Neural Networks, IJCNN 2020, Glasgow, United Kingdom, July 19-24, 2020*, pages 1–8. IEEE, 2020.
- [19] Soheil Kolouri, Aniruddha Saha, Hamed Pirsiavash, and Heiko Hoffmann. Universal litmus patterns: Revealing backdoor attacks in cnns. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*,

- CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pages 298–307. IEEE, 2020.
- [20] Shaofeng Li, Shiqing Ma, Minhui Xue, and Benjamin Zi Hao Zhao. Deep learning backdoors. *CoRR*, abs/2007.08273, 2020.
- [21] Yiming Li, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A survey. *CoRR*, abs/2007.08745, 2020.
- [22] Yiming Li, Tongqing Zhai, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shutao Xia. Rethinking the trigger of backdoor attack. *CoRR*, abs/2004.04692, 2020.
- [23] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In Michael Bailey, Thorsten Holz, Manolis Stamatogiannakis, and Sotiris Ioannidis, editors, *Research in Attacks, Intrusions, and Defenses - 21st International Symposium, RAID 2018, Heraklion, Crete, Greece, September 10-12, 2018, Proceedings*, volume 11050 of *Lecture Notes in Computer Science*, pages 273–294. Springer, 2018.
- [24] Yingqi Liu, Wen-Chuan Lee, Guanhong Tao, Shiqing Ma, Yousra Aafer, and Xiangyu Zhang. ABS: scanning neural networks for back-doors by artificial brain stimulation. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019*, pages 1265–1282. ACM, 2019.
- [25] Yuntao Liu, Yang Xie, and Ankur Srivastava. Neural trojans. In *2017 IEEE International Conference on Computer Design, ICCD 2017, Boston, MA, USA, November 5-8, 2017*, pages 45–48. IEEE Computer Society, 2017.
- [26] Shervin Minaee, AmirAli Abdolrashidi, Hang Su, Mohammed Bennamoun, and David Zhang. Biometric recognition using deep learning: A survey. *CoRR*, abs/1912.00271, 2019.
- [27] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 86–94. IEEE Computer Society, 2017.
- [28] Wajahat Nawaz, Sagheer Ahmed, Ali Tahir, and Hassan Aqeel Khan. Classification of breast cancer histology images using ALEXNET. In Aurélio Campilho, Fakhri Karray, and Bart M. ter Haar Romeny, editors, *Image Analysis and Recognition - 15th International Conference, ICIAR 2018, Póvoa de Varzim, Portugal, June 27-29, 2018, Proceedings*, volume 10882 of *Lecture Notes in Computer Science*, pages 869–876. Springer, 2018.
- [29] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Int. J. Comput. Vis.*, 128(2):336–359, 2020.
- [30] Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*, pages 6106–6116, 2018.
- [31] Di Tang, Xiaofeng Wang, Haixu Tang, and Kehuan Zhang. Demon in the variant: Statistical analysis of dnns for robust backdoor contamination detection. *CoRR*, abs/1908.00686, 2019.
- [32] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*, pages 8011–8021, 2018.
- [33] Loc Truong, Chace Jones, Brian Hutchinson, Andrew August, Brenda Praggastis, Robert Jasper, Nicole Nichols, and Aaron Tuor. Systematic evaluation of backdoor data poisoning attacks on image classifiers. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR Workshops 2020, Seattle, WA, USA, June 14-19, 2020*, pages 3422–3431. IEEE, 2020.
- [34] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *CoRR*, abs/1912.02771, 2019.
- [35] Miguel Villarreal-Vasquez and Bharat K. Bhargava. Confoc: Content-focus protection against trojan attacks on neural networks. *CoRR*, abs/2007.00711, 2020.
- [36] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y. Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy, SP 2019, San Francisco, CA, USA, May 19-23, 2019*, pages 707–723. IEEE, 2019.

- [37] Qinglong Wang, Wenbo Guo, Kaixuan Zhang, Alexander G. Ororbia II, Xinyu Xing, Xue Liu, and C. Lee Giles. Adversary resistant deep neural networks with an application to malware detection. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, pages 1145–1153. ACM, 2017.
- [38] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y. Zhao. Latent backdoor attacks on deep neural networks. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019*, pages 2041–2055. ACM, 2019.
- [39] Haoti Zhong, Cong Liao, Anna Cinzia Squicciarini, Sen-cun Zhu, and David J. Miller. Backdoor embedding in convolutional neural network models via invisible perturbation. In Vassil Roussev, Bhavani M. Thuraisingham, Barbara Carminati, and Murat Kantarcioglu, editors, *CODASPY '20: Tenth ACM Conference on Data and Application Security and Privacy, New Orleans, LA, USA, March 16-18, 2020*, pages 97–108. ACM, 2020.

Appendices

A Attacking State-of-the-Art defenses

Data Inspection: STRIP-ViT [12] Fig. 23 illustrates ϵ -Attack (target class, “dog”) on STRIP for a typical example of CIFAR-10 and 5 intentional perturbations per sample (Note that for evaluation in Section 5.1, we use 2000 perturbations per sample in coherence with the authors [12]). The perturbed inputs are presented to a DNN for classification and the entropy of decision is recorded in the last column of Fig. 23. For $\epsilon = 1.0$ -Attack, a DNN classifies most of the perturbed poisoned inputs as “dog”. This static behavior results in a small output entropy, making it suspicious to STRIP based on an auto-computed threshold ($t = 0.699$ in our case). On contrary, $\epsilon = 0.4$ -Attack causes a large output entropy for perturbed inputs, and hence, successfully evades the detection. A no-Attack case is also presented for comparison.

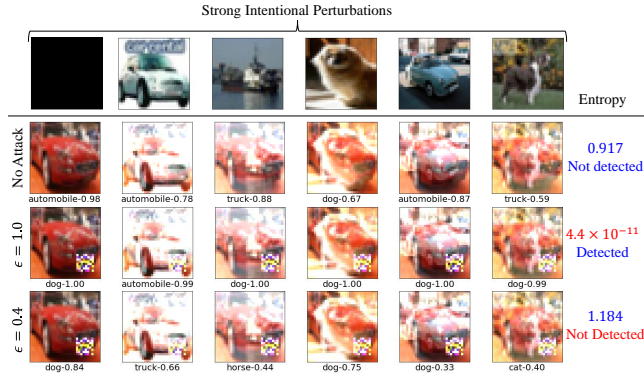


Figure 23: An illustration of STRIP defense working and $\epsilon = 0.4$ attack evading STRIP. STRIP uses an auto-computed threshold of $t=0.699$ (in this case), to differentiate clean and poisoned inputs.

Poison Suppression: Gradient-Shaping [16] Fig. 24 gives the gradient norms of poisoned training samples for different iterations as compared to those of clean samples, with gradient-shaping applied. We use a clipping norm of 4.0 and a noise multiplier of 0.01. We observe that the gradient norms for $\epsilon = 0.4$ more closely follow the gradient norms for clean samples. A peak in the gradients for $\epsilon = 1.0$ -Attack results in stronger updates, causing backdoors to be learned more quickly. This further validates our hypothesis regarding the early learning of backdoors in Section 1.

Model Inspection: Universal Litmus Patterns [19] Fig. 25 shows optimized so-called litmus patterns for the case of $M = 10$. These litmus-patterns are given as inputs to a DNN under inspection and the output of the DNN is observed by a detector (trained together with the litmus patterns) for suspicious behavior.

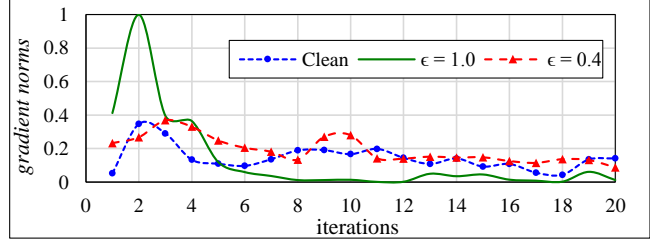


Figure 24: Normalized average gradient norms of poisoned training samples with $\epsilon = 0.4$ and $\epsilon = 1.0$ labels. Gradient norm for clean samples is given for comparison.

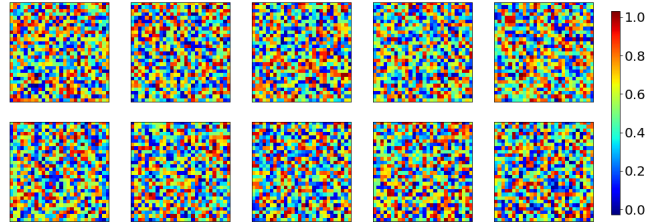


Figure 25: Optimized litmus Patterns for ULP-10 defense.

Blind defense: Februs [7] Fig. 26 illustrates a typical example of our proposed attacks on Fashion-MNIST for different ϵ values. For the simple ϵ -Attack, a trigger influencing the decision is identified and removed by Februs. Februs does the same for ϵ^2 -Attacks - removes $Z_1 \cap Z_2$ - unaware that doing this will cause Z_2 to get activated as shown in the figure.

B Evaluating HaS-Nets

1) Fashion-MNIST: Fashion-MNIST consists of 70000 28x28 grey-scale images belonging to 10 different classes. The training set contains 60000 images while the test set has 10000 images. When conducting attacks, we poison the first 600 images of the training set.

The architecture used for evaluating HaS-Nets on Fashion-MNIST is given in Table 2. Results are reported in Fig.27 for different iterations. We observe that HaS-Nets are effective irrespective of the number of iterations. We also note relatively larger ASRs for the first iteration. This can be attributed to HaS-Nets allowing the DNN to train on full training data for the first iteration.

2) CIFAR-10: CIFAR-10 contains 60000 32x32x3 images belonging to 10 different classes, divided as 50000 training images and 10000 test images. When conducting attacks, we poison the first 600 images of the training set.

The architecture used for the analysis is given in Table 2. Results are reported in Fig.28 for different iterations. Again, we observe the effectiveness of HaS-Nets for different iterations.

3) IMDB: The IMDB dataset contains 50000 movie reviews categorized into two classes: positive and negative. There are 25000 samples in the training set and the same in the test set. When conducting attacks, we poison the first 600 images of

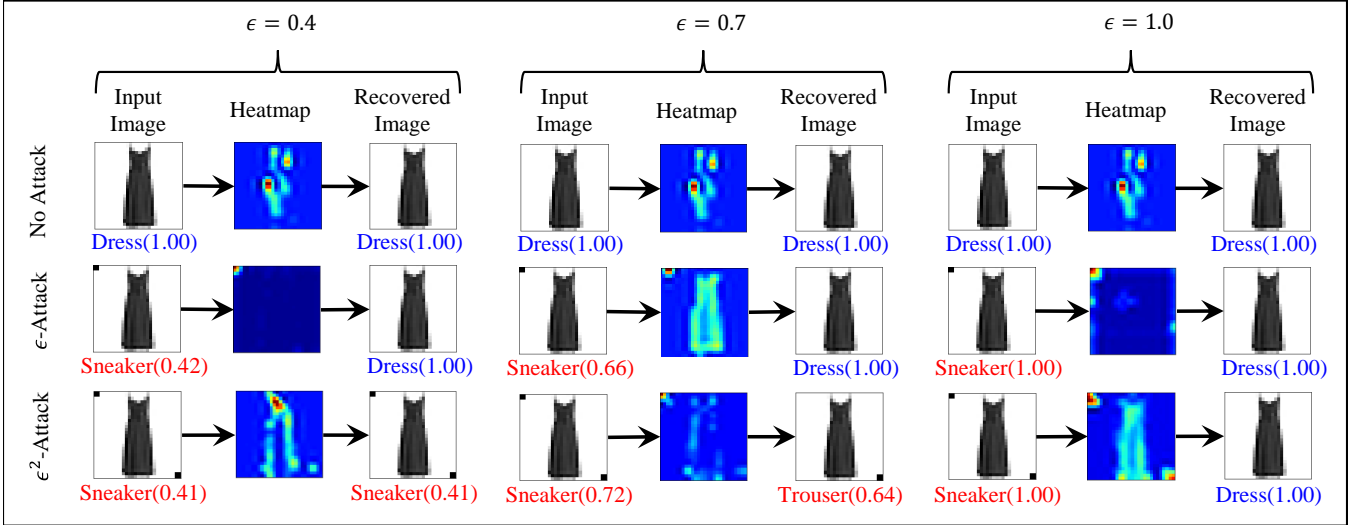


Figure 26: A typical example, illustrating ϵ^2 -Attack evading Februus [7]. Februus searches for potential triggers in an input image using heatmaps, and mask those triggers. The masked image is fed to a GAN, which recovers the input image trigger-free. The input image, heatmaps and the recovered images are given for different scenarios and settings.

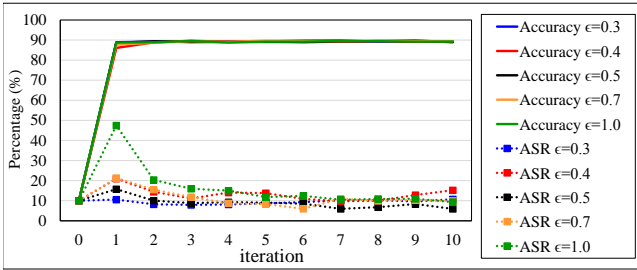


Figure 27: Accuracy of HaS-Nets and Attack Success Rates(ASRs) of ϵ -Attack against HaS-Nets on Fashion-MNIST for different iterations.

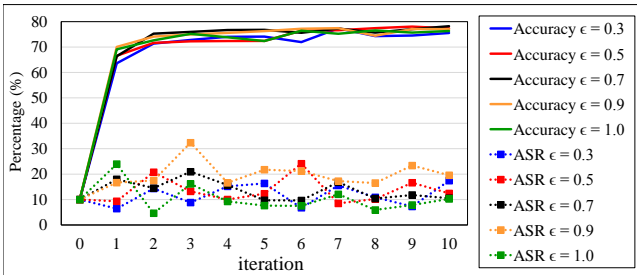


Figure 28: Accuracy of HaS-Nets and ASRs of ϵ -Attack against HaS-Nets on CIFAR-10 for different iterations.

the training set.

The architecture used for this analysis is given in Table 2 and results are shown in Fig. 29. We observe a reduction in ASR from about 80% to 60%. Although, this is not a very appreciable decrease, recall that the IMDB dataset only has two classes, and even if the attack were to fail, it would still have a success rate of about 50%.

4) Mini-Consumer Complaint: Mini-Consumer Complaint dataset consists of 66806 samples of 11 different classes,

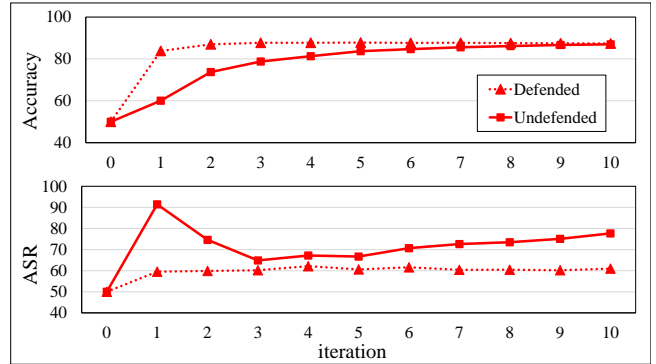


Figure 29: Accuracy of HaS-Nets and ASRs of ϵ -Attack against HaS-Nets on IMDB dataset for different iterations.

divided into 50104 training samples and 16702 test samples. When conducting attacks, we poison the first 600 images of the training set. The architecture used and the results are given in Table 2 and Fig. 30, respectively.

5) Urban Sound 8k: For evaluation on audio tasks we choose Urban Sound 8k, containing 8732 audio samples of $\leq 4s$ belonging to 10 different classes, divided into 6549 and 2183 samples in the training and test set, respectively. When conducting attacks, we poison the first 80 images of the training set. The architecture used and the results are given in Table 2 and Fig. 30, respectively.

C Implementation Details

Invisible Backdoor Attack. For invisible backdoor attack, we generate a uniform random noise of the same size as input images. The noise is sampled from $\mathcal{U}(-0.1, 0.1)$ and used as a trigger to poison 1.2% of the training samples. For this experiment, we randomly choose target class to be “horse” and “sneakers” for CIFAR-10 and Fashion-MNIST, as shown

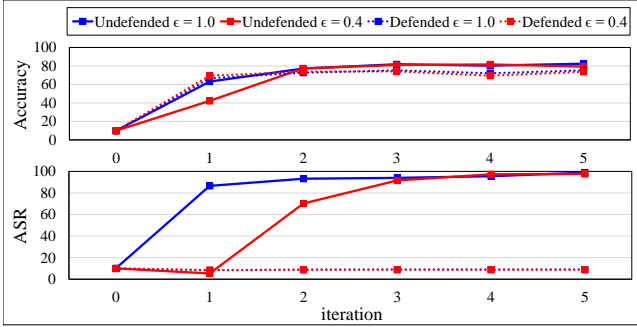


Figure 30: Accuracy of HaS-Nets and ASRs of ϵ -Attack against HaS-Nets on Consumer Complaint dataset for different iterations.

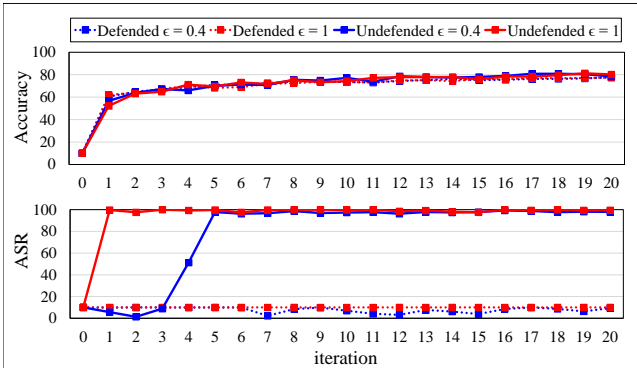


Figure 31: Accuracy of HaS-Nets and ASRs of ϵ -Attack against HaS-Nets on Urban Sound dataset for different iterations.

in Fig. 32 and Fig. 33, respectively.

Label-Consistent Backdoor Attack. For Label-consistent backdoor attack, we target 25% of the “horse” images from CIFAR-10 training data. Specifically, we train an auto-encoder for CIFAR-10 and interpolate the latent representation (encoded image) of a randomly chosen image (not from the target class) towards the latent representation of a target image, which an attacker aims to poison. These interpolated representations are then decoded using the decoder. Decoded images are usually referred to as interpolated images [34]. Fig. 36 shows interpolated images for several degrees of interpolation. For our experiments, we choose a degree of 0.8 due to its similarity with the target image. We repeat the same procedure for all the target images and stamp a trigger on the interpolated images without changing their labels, as shown in Fig. 34 for randomly chosen samples from CIFAR-10. Notice how the labels (shown above) assigned to poisoned images are now consistent for human observers. We follow the same steps for Fashion-MNIST. Example images are shown in Fig. 35.

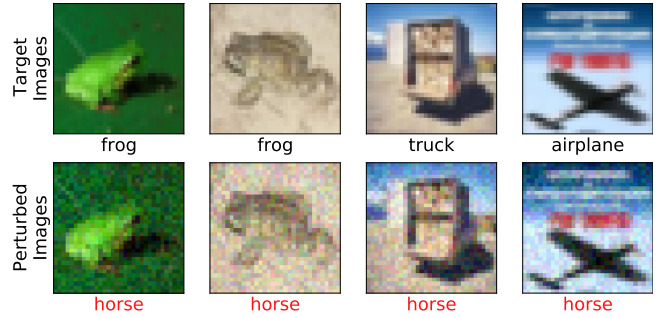


Figure 32: CIFAR-10 samples poisoned with invisible backdoor attack as compared to clean samples. The target class is chosen to be “horse” by the attacker.

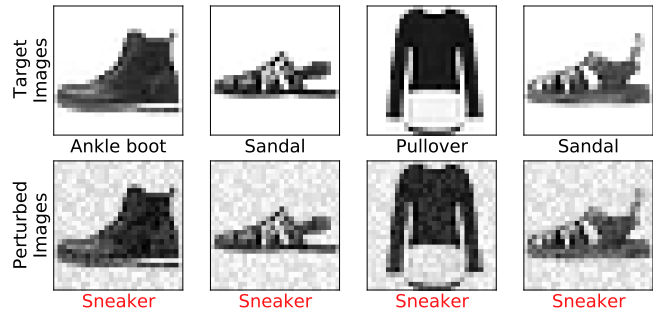


Figure 33: Fashion-MNIST samples poisoned with invisible backdoor attack as compared to clean samples. The target class is chosen to be “sneaker” by the attacker.

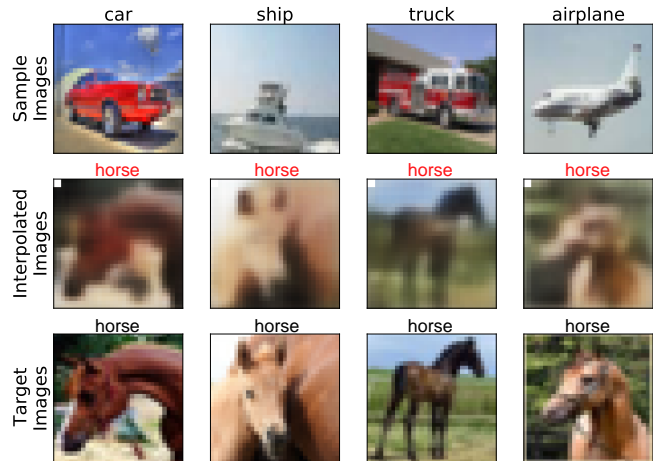


Figure 34: The top row shows sample images of CIFAR-10 which we interpolate towards the target images shown in the bottom row. The target images in the clean training set are then replaced by the poisoned images shown in the second row. The target class is chosen to be “horse” by the attacker. Notice the consistency between the poisoned images and the labels.



Figure 35: The top row shows sample images of Fashion-MNIST which we interpolate towards the target images shown in the bottom row. The target images in the clean training set are then replaced by the poisoned images shown in the second row. The target class is chosen to be “bag” by the attacker. Notice the consistency between the poisoned images and the labels.

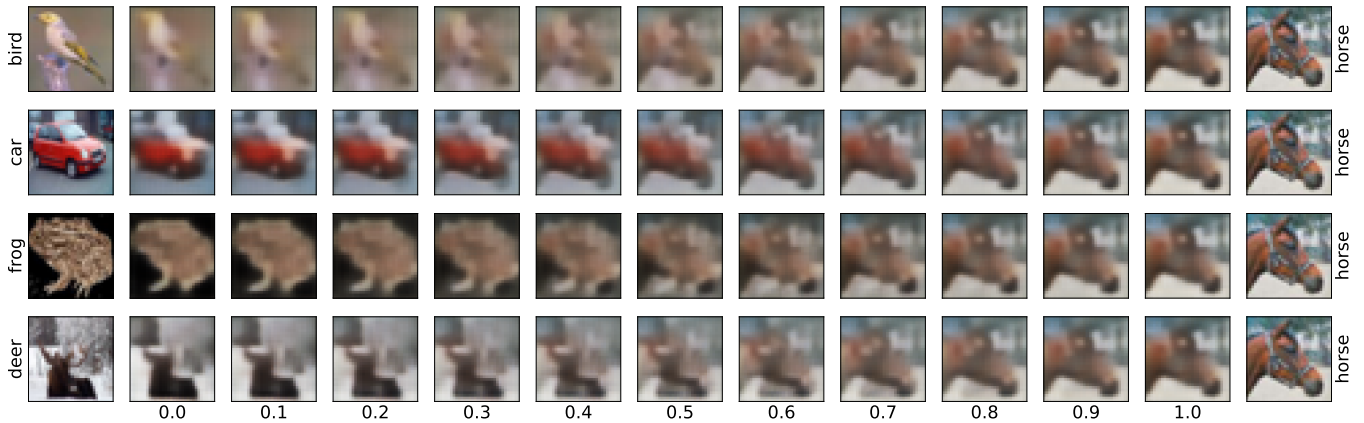


Figure 36: Interpolated images for a range of interpolation degrees for randomly selected initial samples (first column) from CIFAR-10 test data. For illustration, the target image (last column) is chosen to be the same for each initial sample.

Table 2: Network Architectures for Evaluating HaS-Nets on different Datasets

Layer	Fashion-MNIST	CIFAR-10	IMDB	Consumer Complaint	Urban Sound
Embedding	-	-	YES	YES	-
Dropout	-	-	-	0.3	-
Conv()	32x3x3	32x3x3	-	-	-
Activation	Elu()	Elu()	-	-	-
Normalize	YES	YES	-	-	-
Conv()	32x3x3	32x3x3	-	128x5	64x3x3
Activation	Elu()	Elu()	-	-	tanh()
Normalize	YES	YES	-	-	-
MaxPool	2x2	2x2	-	5	2x2
Dropout	0.2	0.2	-	0.3	0.1
Conv()	-	64x3x3	-	-	-
Activation	-	Elu()	-	-	-
Normalize	-	YES	-	YES	-
Conv()	-	64x3x3	-	128x5	128x3x3
Activation	-	Elu()	-	Relu()	tanh()
Normalize	-	YES	-	-	-
MaxPool	-	2x2	-	5	2x2
Dropout	-	0.2	-	0.3	0.1
Conv()	-	128x3x3	-	-	-
Activation	-	Elu()	-	-	-
Normalize	-	YES	-	YES	-
Conv()	-	128x3x3	-	-	-
Activation	-	Elu()	-	-	-
Normalize	-	YES	-	-	-
MaxPool	-	2x2	-	-	-
Dropout	-	0.2	-	-	-
Dense	-	-	16	-	-
Activation	-	-	-	-	-
Dense	-	-	-	128	1024
Activation	-	-	-	Relu()	tanh()
Dense+Softmax	10	10	1	11	10