Identify and Remove Backdoor Neurons via Clean-Poisoned Mixture Distribution

Anonymous Author(s) Affiliation Address email

Abstract

Convolutional neural networks (CNN) can be manipulated to perform specific 1 behaviors when encountering a particular trigger pattern without affecting the 2 performance on normal samples, which is referred to as backdoor attack. Backdoor 3 attack is usually achieved by injecting a small proportion of poisoned samples 4 into the training set, through which the victim trains a model embedded with the 5 designated backdoor. In this work, we demonstrate that the backdoor neurons in 6 an infected neural network have a mixture of two distributions with significantly 7 different moments, formed by benign samples and poisoned samples, respectively. 8 This property is shown to be attack-invariant and allows us to efficiently locate 9 the backdoor neurons. On this basis, we make several realistic assumptions on 10 the neuron activation distributions, and propose two backdoor neuron detection 11 12 strategies based on (1) the differential entropy of the neurons, and (2) the KL divergence between the benign sample distribution and a poisoned statistics based 13 hypothetical distribution. Experimental results show that our proposed defense 14 strategies are both efficient and effective against various backdoor attacks. 15

16 1 Introduction

Convolutional neural networks (CNNs) have achieved tremendous success during the past few years 17 in a wide range of areas. However, training a CNN from scratch involves a large amount of data and 18 expensive computational cost, which is sometimes infeasible. A more practical strategy is to obtain 19 pretrained models or utilize public datasets from a third party, which brings convenience but also 20 raises severe security problems into the deployment of models. For example, a malicious third party 21 may provide pretrained models embedded with a designated backdoor, such that the model will have 22 a predefined response to some specific pattern, which is also called the *trigger*. More realistically, the 23 attacker can inject only a small proportion of malicious data into the public dataset to mislead the 24 trained model, which is referred to as *backdoor poisoning attacks* [26]. For instance, the malicious 25 data can be created by patching a particular pattern into the benign data, and changing the label to 26 the desired target. The correlation of the trigger and the specified target label will be learned by the 27 models during the training time. In this way, the infected model will misclassify the input to the 28 attack target when the pattern is patched, while behave normally otherwise, as shown in Figure 1. 29

Due to the limited understanding of CNNs, we are not clear on the formation mechanism of backdoor behaviors. However, it was empirically found that a infected model always possesses one or more neurons that have high correlation with the trigger activation, and pruning these neurons can significantly alleviate the backdoor behaviors, while retaining the model performance [40, 28, 6]. Nevertheless, how to precisely find out these backdoor neurons in a infected model is still a challenging problem, and has attracted a lot of attentions from the community.



Figure 1: An overview of an infected learning system. The images with a white square are classified as class 0. It is an empirical observation that the backdoor behaviors are always triggered by one or more backdoor neurons. We demonstrate that the distributions of activations on these neurons are mixtures of two Gaussian-like distributions, formed by benign samples and poisoned samples, respectively.

36 In this work, we will take an inspection on the pre-activation distributions of infected models. In general, the activations in each neuron follow an unimodal distribution that can be approximated by a 37 Gaussian distribution. Based on the maximum entropy property of the Gaussian distribution, these 38 benign neurons should have relatively large entropy on their pre-activation distributions. However, in 39 the backdoor neurons, the distribution is always bimodal and can be approximated by a mixture of 40 two Gaussian distributions, formed by the benign data and poisoned data, respectively; see Figure 1. 41 This kind of distribution must have a lower entropy, compared with a Gaussian distribution with 42 the same finite variance. Thus, if we standardize the pre-activation distributions in all neurons, the 43 backdoor neurons should possess a relatively low entropy compared with the benign neurons. We 44 treat the low entropy neurons as outliers and potential backdoor neurons, and prune these neurons 45 to recover the model without retraining. However, under another defense setting, in which only 46 an infected model and a small set of benign data are provided, we may not be able to observe the 47 bimodal distributions on the backdoor neurons since the poisoned data is absent. In this case, we can 48 rely on the statistics on *Batch Normalization* (BN) layers. Specifically, if the infected model is trained 49 on poisoned data, the statistics of backdoor neurons recorded in the BN layer will be significantly 50 different from the distribution of only benign data. More importantly, the mismatch of statistics will 51 not exist in benign neurons. Based on the neuron entropy and the statistics discrepancy, we are able 52 to locate and prune the backdoor neurons to recover the model under two defense settings. 53

54 In summary, our contributions include:

- We take a deep inspection on the infected model, and summarize the law of pre-activation distributions on poisoned dataset. We find that (1) the standardized entropy of backdoor neurons can be significantly lower than benign neurons, and (2) the BN statistics in infected model are mismatched with the benign sample statistics.
- We propose to prune the potential backdoor neurons based on either the sample entropy or
 the statistics discrepancy, depending on the defense settings. Under certain assumptions, we
 claim that both the proposed indices can perfectly separate the benign neurons and backdoor
 neurons by an appropriate threshold.
- 3. We conduct extensive experiments to verify our assumptions and evaluate our proposed
 methods, and achieve the state-of-the-art results under two different defense settings.

65 2 Related work

⁶⁶ In this section, we detailedly discuss recent works in backdoor attack and defense.

67 2.1 Backdoor attacks

The concept of *backdoor attack* is first introduced in [14], where the adversary injects a small set of targeted label-flipped data stamped with a small and specific trigger into the training set during the DNN training, leading to a misclassification when predicting the samples with such trigger. To make the trigger pattern even more invisible to human beings, the blending strategy is used in [5]

to generate poison images, while the form of natural reflection is utilized in trigger design in [31]. 72 The input image is perturbed in [38] to keep its content consistent with the target label such that the 73 model better memorizes the trigger pattern, and keep it imperceptible to human beings. Moreover, 74 the multi-target and multi-trigger attacks are proposed in [42, 32], and make the attack more flexible 75 and covert. Recently, some sample specific trigger design strategies [27] are proposed, making the 76 defense against such backdoor attack much harder. Generally, the above attacks can be referred as 77 78 the poisoning based backdoor attacks. Under some settings, the attackers can control the training process to inject the backdoor without 79

modifying the training data, referred as the *non-poisoning based backdoor attacks*. This is achieved in [30, 34, 4] through targeted modification of the neurons' weight in a network. Such attacks will

⁸² not be evaluated in our work due to its strong attack setting.

83 2.2 Backdoor defenses

Training stage defense. Under such setting, the defender has access to the training process, so that
they can detect and filter the poisoned data or add some restrictions to suppress the backdoor effect
in training. Since the poisoned data can be regarded as outliers, different strategies are applied in
[10, 12, 1, 37, 15], such as the *robust statistics* in feature space and *input perturbation techniques*to filter them out of training data. Other methods aim at suppressing the backdoor effect during
training phase [25] with strong data augmentation [2], such as *CutMix* [43], *CutOut* and *MaxUp* [13], *differential privacy* [11, 18].

Model post-processing defense. Sometimes the defenders are only given a suspicious DNN model 91 without access to the training process or the full training set. Therefore, they must eliminate the 92 backdoor threat with limited resources, such as a small set of clean data. A straightforward way is 93 to reconstruct the trigger by adding adversarial perturbations to the input images, and then detoxify 94 the model with the knowledge of reversed trigger [39]. Some try to find the relationship between 95 backdoor behaviors and the neurons in a DNN model. Different levels of stimulation to a neuron 96 are introduced in [29] to see how to determine the output activation change, if the model is attacked. 97 Simple neuron pruning strategies are applied in [6] to repair the model, while adversarial perturbations 98 are added to the neurons in [40] and precisely prunes the backdoor neurons with more limited clean 99 data requirement and better performance. Simple redundant neuron pruning and fine-tuning are 100 combined together in [28] to erase the backdoor effect. There are other fine-tuning based methods 101 with the implementation of knowledge distillation [24, 17], while they may suffer from the hyper-102 parameter tuning and clean accuracy dropping problems due to limited prior knowledge of the attack 103 and over-fitting on the clean set. Mode connectivity repair technique [44] is also explored to mitigate 104 the backdoored model. Recently, the K-Arm optimization [36] is applied in backdoor detection, 105 and outperforms other adversarial perturbation based methods [39, 8], helping curtail the threat of 106 backdoor attack. 107

108 3 Preliminaries

109 3.1 Notations

Consider a multi-class classification problem with C classes. Let the original training set $\mathcal{D} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$ contains N i.i.d. sample images $\boldsymbol{x}_i \in \mathbb{R}^{d_c \times d_h \times d_w}$ and the corresponding labels $y_i \in \{1, 2, ..., C\}$ drawn from $\mathcal{X} \times \mathcal{Y}$. Here, we denote by d_c , d_h and d_w the number of neurons, the height and the width of images, respectively. In particular, we have $d_c = 3$ for RGB images.

As in Section 2.1, the backdoor poisoning attack involves changes to the input images and the corresponding labels on a subset $\mathcal{D}_p \subseteq \mathcal{D}$. In this work, we define the ratio $\rho = \frac{|\mathcal{D}_p|}{|\mathcal{D}|}$ as the *poisoning*

rate. We denote the poisoning function to the input images as $\delta(x)$.

117 Consider a neural network $F(x; \theta)$ with L layers. Denote

$$F^{(l)} = f^{(l)} \circ \phi \circ f^{(l-1)} \circ \phi \circ \cdots \circ \phi \circ f^{(1)},$$

119 activation function applied element wise. In this work, we may denote $F(x;\theta)$ as F(x) or F for 120 simplicity.

for $1 \le l \le L$, where $f^{(l)}$ is a linear function (*e.g.*, convolution) in the *l*-th layer, and ϕ is a nonlinear

121

We denote by $W^{(l)} \in \mathbb{R}^{d_{c'} \times d_c \times d_h \times d_w}$ the weight tensor of a convolutional layer. To do pruning, we apply a mask $M^{(l)} \in \{0,1\}^{d_{c'} \times d_c \times d_h \times d_w}$ starting with $M^{(l)} = \mathbb{1}_{d_{c'} \times d_c \times d_h \times d_w}$ in each layer. Pruning neurons on the network refers to getting a collection of indices $\mathcal{I} = \{(l,k)_i\}_{i=1}^{I}$ and setting 122 123 $M_k^{(l)} = \mathbf{0}_{d_c \times d_h \times d_w}$ if $(l, k) \in \mathcal{I}$. The pruned network $F_{-\mathcal{I}}$ has the same architecture as F but with all the weight matrices of convolutional layers set to $W^{(l)} \odot M^{(l)}$, where \odot denotes the Hadamard 124 125 product. 126

3.2 Differential entropy 127

To measure the uncertainty of a discrete random variable Z, the *entropy* [35, 7] was defined as 128 $H(Z) = -\sum_{z \in Z} p(z) \log p(z)$. At the same time, as an extension of entropy, the *differential* entropy was also introduced for a continuous random variable. More concretely, if Z is a continuous 129 130 random variable, then it was defined as 131

$$h(Z) = -\int_{Z} p(z) \log p(z) dz.$$
⁽¹⁾

An important fact about the differential entropy is that, among all the real-valued distributions 132 supported on $(-\infty,\infty)$ with a specified finite variance, the Gaussian distribution maximizes the 133 differential entropy [7]. In this work, the differential entropy (1) will be utilized to identify the 134 distributions that are far different from a Gaussian distribution. 135

3.3 **Backdoor neurons** 136

It was found that there exist one or more neurons that contribute the most to the backdoor behaviors 137 in a infected model [40, 28]. If some of or all of these neurons are pruned, the attack success rate 138 will be reduced greatly [40]. However, to our knowledge, there was no such quantity defined in the 139 literature to measure how important a neuron is to the backdoor behaviors. 140

In this work, to make up for this lack, we would like to introduce the sensitivity of neurons to the 141 backdoor by the backdoor loss calculated before and after pruning the neuron: 142

$$\alpha(F,l,k) = \mathcal{L}_{\mathrm{bd}}(F) - \mathcal{L}_{\mathrm{bd}}(F_{-\{(l,k)\}}),$$
(2)

where $F_{-\{(l,k)\}}$ is the network after pruning the k-th neuron of the l-th layer. The backdoor loss is 143 defined as: 144

$$\mathcal{L}_{\mathrm{bd}}(f) = \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}}[\mathrm{D}_{\mathrm{CE}}(y, f(\delta(\boldsymbol{x}))],$$

where D_{CE} denotes the cross entropy loss. The backdoor loss is high when the model is infected, and 145 will be reduced when the backdoor effect is alleviated. 146

Using the quantity define in (2), we are now able to find the neurons that are mostly correlated with 147 the backdoor behaviors. More concretely, we first set a threshold $\tau > 0$, and then go through all 148 the neurons to find the ones with the quantity (2) above τ . We call them the *backdoor neurons*, and 149 denote 150

$$\mathcal{B}_{F,\tau} = \{(l,k) : \alpha(F,l,k) > \tau\}.$$
(3)

Pre-activation distribution 3.4 151

During the forward propagation of an input x, we denote $x^{(l)} = F^{(l)}(x) \in \mathbb{R}^{d_c^{(l)} \times d_h^{(l)} \times d_w^{(l)}}$ as the 152 output of the *l*-th layer. For the *k*-th neuron of the *l*-th layer, the *pre-activation* $\phi_k^{(l)} = \phi(\boldsymbol{x}_k^{(l)})$ is defined as the maximum value of the *k*-th slice matrix of dimension $d_h^{(l)} \times d_w^{(l)}$ in $\boldsymbol{x}^{(l)}$. 153 154

It is a common assumption that, for every neuron, the pre-activations before the non-linear function 155 follow a Gaussian distribution, if the network is randomly initialized and the number of neurons is 156 large enough [21, 3]. In a trained network, although this assumption may not strictly hold, the pre-157 activation of every neuron can be still regarded as approximately following a Gaussian distribution. 158 However, in this work, for the first time, we observe a bimodal distribution in the backdoor neurons, 159 two components of which are formed by the benign data and poisoned data, respectively. This 160 phenomenon is shown in Figure 2, where a typical backdoor neuron is compared with the benign 161 neurons. It can be seen that, after the model is infected, the pre-activation distributions of benign 162 neurons hardly change when the data is poisoned, while the pre-activation distributions of backdoor 163 neurons become significantly different. 164



Figure 2: In (a) and (b), we compare the pre-activation distributions in backdoor neurons and benign neurons. In benign neurons, the pre-activation distributions on benign data and poisoned data are nearly the same, while in backdoor neurons, they show great difference. In (c) and (d), we plot the empirical distributions with benign samples (in blue) and the BN statistics induced Gaussians (in green). In backdoor neurons, the discrepancy between the empirical and BN induced distribution is large (all the neurons are selected from infected ResNet-18 trained on CIFAR-10, and use 1,000 images with (poisoned data) or without trigger (benign data) as the inputs).

165 4 Methodology

166 4.1 Assumptions

¹⁶⁷ A primary assumption is that $|\mathcal{B}_{F,\tau}| > 0$ for a poisoned model F and a pre-defined threshold $\tau > 0$, ¹⁶⁸ where $\mathcal{B}_{F,\tau}$ is defined as in (3). As in Section 3.4, we also assume that the pre-activations of all ¹⁶⁹ neurons follow a Gaussian mixture distribution, that is:

$$\phi_k^{(l)} \sim (1 - \rho) \mathcal{N}(\mu_k^{(l)}, \sigma_k^{(l)2}) + \rho \mathcal{N}(\hat{\mu}_k^{(l)}, \hat{\sigma}_k^{(l)2}), \tag{4}$$

where ρ is the poisoning rate of the dataset, $\mu_k^{(l)}$ and $\sigma_k^{(l)2}$, $\hat{\mu}_k^{(l)}$ and $\hat{\sigma}_k^{(l)2}$ are the mean and variance of $\{\phi(F^{(l)}(\boldsymbol{x})_k) : \boldsymbol{x} \sim \mathcal{X}\}, \{\phi(F^{(l)}(\delta(\boldsymbol{x}))_k) : \boldsymbol{x} \sim \mathcal{X}\}$, respectively.

172 We further assume that:

$$|\mu_k^{(l)} - \hat{\mu}_k^{(l)}| \begin{cases} < \epsilon, & \text{if } (l,k) \notin \mathcal{B}_{F,\tau} \\ >> \epsilon, & \text{if } (l,k) \in \mathcal{B}_{F,\tau} \end{cases}$$

173 and

$$|\sigma_k^{(l)2} - \hat{\sigma}_k^{(l)2}| \begin{cases} < \epsilon^2, & \text{if } (l,k) \notin \mathcal{B}_{F,\tau}, \\ >> \epsilon^2, & \text{if } (l,k) \in \mathcal{B}_{F,\tau}, \end{cases}$$

where $\epsilon > 0$ is a small enough value. Note that, if $|\mu_k^{(l)} - \hat{\mu}_k^{(l)}| = |\sigma_k^{(l)2} - \hat{\sigma}_k^{(l)2}| = 0$, the pre-activation follows a Gaussian distribution $\mathcal{N}(\mu_k^{(l)}, \sigma_k^{(l)2})$.

176 4.2 Discrepancy of differential entropy (DDE)

If the pre-activation distributions are standardized (subtracting the mean and dividing the standard deviation), the differential entropy will be maximized on the benign neuron approximately following a standard Gaussian distribution $\mathcal{N}(0, 1)$. In backdoor neurons, because of the difference of the moments on Gaussian components, the resulting mixture distributions can not be Gaussian distributions, hence the differential entropy must be smaller than h(Z), where $Z \sim \mathcal{N}(0, 1)$ is the standardized Gaussian distribution. Specifically, let $\dot{\phi}_k^{(l)} = \frac{\phi_k^{(l)} - \mu_k^{(l)}}{\sigma_k^{(l)}}$ be the standardized pre-activations, then

$$h(\dot{\boldsymbol{\phi}}_{k}^{(l)}) < h(\dot{\boldsymbol{\phi}}_{k'}^{(l)}) \le h(Z), \quad \forall k \in \mathcal{B}_{F,\tau}, k' \notin \mathcal{B}_{F,\tau}.$$

This inequality gives a guarantee that with an appropriately chosen threshold, the backdoor neurons can be well separated with the benign neurons.

185 4.3 Mismatched BN statistics (MBNS)

BN layer involves using the statistics of a mini-batch to normalize the data in each layer for each neuron. It is known to be able to smooth the optimization landscape, and has gradually become a

default setting of neural networks [19]. When inference, BN uses the fixed statistics obtained by 188 averaging the sample statistics of mini-batches during training time, including the mean and the 189 variance. If the model is trained on a poisoned dataset, BN will record the mean and the variance 190 of the poison-benign mixed data. Note that the mean and variance here are not defined on the 191 pre-activations $\phi_k^{(l)}$, but on $x_k^{(l)}$. Based on the above discussions, we know that the poisoned samples (especially the pre-activations) on the backdoor neurons follow a different distribution from the 192 193 benign samples. The recorded statistics during training are actually that of the mixture distribution. 194 Hence, we can expect that the BN statistics of a trained backdoor neural network are biased. If we 195 are able to access a small set of benign data, we can calculate an approximation of the true statistics 196 on benign data. Then we calculate the Kullback-Leibler (KL) divergence [9] between the sample 197 distribution and the BN induced distribution as the measurement of the bias. By assuming both of the 198 distributions follow Gaussian distributions, we have a closed form solution: 199

$$D_{\rm KL}(\mathcal{N}_{\rm sample}^{(l)}, \mathcal{N}_{\rm BN}^{(l)})_k = \log \frac{\tilde{\sigma}_k^{(l)}}{\bar{\sigma}_k^{(l)}} + \frac{\bar{\sigma}_k^{(l)2} + (\bar{\mu}_k^{(l)} - \tilde{\mu}_k^{(l)})^2}{2\tilde{\sigma}_k^{(l)2}} - \frac{1}{2}$$

where $\bar{\mu}_k^{(l)}$ and $\bar{\sigma}_k^{(l)2}$ are the statistics obtained from benign samples, $\tilde{\mu}_k^{(l)2}$ and $\tilde{\sigma}_k^{(l)}$ are the BN statistics. The backdoor neurons should have abnormally large KL divergences, as illustrated in Figure 2(c).

203 4.4 Overview of the two pruning strategies

In Section 3.4, we reveal the discrepancy between the pre-activation distributions in the backdoor neurons and that in the benign neurons. This enables fast detecting the neurons that are more related to the backdoor behaviours. The index we choose to detect the abnormal neurons depends on what kind of data we are able to access.

Mixture training data In this case, the victim is given a poisoned training dataset with a specified poisoning rate ρ . Our goal is to obtain a benign model based on the poisoned dataset. To achieve this, we first train an infected model F on the poisoned dataset. The resulting model should have a certain number of backdoor neurons based on empirical observation and the assumption. Since $\rho > 0$ for the dataset, all the neurons follow Gaussian mixture distributions, and we have $h(\dot{\boldsymbol{x}}_k^{(l)}) < h(\dot{\boldsymbol{x}}_{k'}^{(l)})$ for all $k \in \mathcal{B}_{F,\tau}, k' \notin \mathcal{B}_{F,\tau}$. This implies that with an appropriate threshold τ_h^* , we can perfectly separate the benign neurons and backdoor neurons, which can be formulated as:

-

$$\begin{aligned} \exists \tau_h^*, \quad h(\dot{\boldsymbol{x}}_k^{(l)}) < \tau_h^*, \quad \forall k \in \mathcal{B}_{F,\tau}, \\ h(\dot{\boldsymbol{x}}_{k'}^{(l)}) > \tau_h^*, \quad \forall k' \notin \mathcal{B}_{F,\tau}. \end{aligned}$$

Setting the threshold τ_h^* is crucial to the solution, and it is a trade-off between the accuracy on benign samples and that on the backdoored samples. Note that $|\mathcal{B}_{F,\tau}^{(l)}| << d_c^{(l)}$. We can treat the low entropy neurons as outliers in each layer, and set different thresholds for different layers. Specifically, let $h^{(l)} = [h(x_1^{(l)}), h(x_2^{(l)}), \dots, h(x_{d_c^{(l)}}^{(l)})]^T \in \mathbb{R}^{d_c^{(l)}}$ be a vector of sample entropy of the *l*-th layer calculated from the poisoned dataset. Then we set $\tau_h^{(l)} = \bar{h}^{(l)} - u_h \cdot s_h^{(l)}$, where $\bar{h}^{(l)} = \frac{1}{d_c^{(l)}} \sum_{k=1}^{d_c^{(l)}} h_k^{(l)}$ and $s_h^{(l)} = \sqrt{\frac{1}{d_c^{(l)}} \sum_{k=1}^{d_c^{(l)}} (h_k^{(l)} - \bar{h}^{(l)})^2}$ are the mean and standard deviation of $h^{(l)}, u_h$ is a hyperparameter controlling how low the threshold is. Then we have a set of indices of potential backdoor neurons $\mathcal{I}_h = \{(l,k): h_k^{(l)} < \tau_h^{(l)}\}$. Finally, we prune the infected model F using \mathcal{I}_h , which result in a final model $F_{-\mathcal{I}_h}$.

Benign training data This is the case that the victim is given a trained poisoned model F with a small set of benign data. Our goal is to utilize the benign data to clean up the poisoned model and eliminate the backdoor threat. Similar to the pruning process based on differential entropy, we first construct a vector of KL divergences of all neurons for each layer $K^{(l)} = [K_1^{(l)}, K_2^{(l)}, \dots, K_{d_c^{(l)}}^{(l)}]^T \in$ $\mathbb{R}^{d_c^{(l)}}$ according to equation (5). We set $\tau_K^{(l)} = \bar{K}^{(l)} - u_K \cdot s_K^{(l)}$, where $\bar{K}^{(l)} = \frac{1}{d_c^{(l)}} \sum_{k=1}^{d_c^{(l)}} K_k^{(l)}$

229 and $s_{K}^{(l)} = \sqrt{\frac{1}{d_{c}^{(l)}} \sum_{k=1}^{d_{c}^{(l)}} (K_{k}^{(l)} - \bar{K}^{(l)})^{2}}$ are the mean and standard deviation of $K^{(l)}$, u_{K} is a

hyperparameter. The set of selected neurons is $\mathcal{I}_K = \{(l,k) : K_k^{(l)} < \tau_K^{(l)}\}$ and the pruned model can be represented as $F_{-\mathcal{I}_K}$. Note that u_K is the **only** hyperparameter of our methods, and is usually set to 3.

233 **5 Experiments**

234 5.1 Implementation details

Datasets In this section, the experiments are conducted on two influential benchmarks, CIFAR-10 [22] and Tiny-ImageNet [23]. We use 90% of the data set for training, 5% for validating, and the rest 5% as the benign data for recovering the poisoned model in the later backdoor defense scenario.

Models We use ResNet-18 [16] as the baseline model to evaluate our proposed method, and compare it with other methods. We train the network for 150 epochs on CIFAR-10 and 100 epochs on Tiny-ImageNet with SGD optimizer. The initial learning rate is set to 0.1 and the momentum is set to 0.9. We adopt the cosine learning rate scheduler to adjust the learning rate. The batch size is set to 128 by default.

Attacks Our experiments are based on both the classical and the most advanced attack strategies, 243 including the BadNet [14], Clean Label Attack (CLA) [38], Reflection Backdoor (Refool) [31], 244 Warping-based poisoned Networks [33], Blended backdoor attack (Blended) [5], Input-aware back-245 246 door attack (IAB) [32] and Sample Specific Backdoor Attack (SSBA) [27]. For BadNets, we test both the All-to-All (A2A) attack and All-to-One (A2O) attack, i.e., the attack target labels are set to 247 $y_t = (y+1) \mod C$, or one particular label $y_t = C_t$, respectively. The target for A2O attacks of all 248 the attack strategies is set to class 0. The triggers for BadNets and CLA are set to randomly generated 249 patterns with size 3×3 for CIFAR-10 and 5×5 for Tiny-ImageNet. The poisoning rate is set to 10% 250 by default. Note that, due to the image size restraint, SSBA is only performed on Tiny-Imagenet. 251

Defenses We conduct experiments under two defense settings, one of which allows the defender to access the poisoned training set, while the other only has a small clean data set. Both the defense goals are to obtain a clean model without backdoor behaviours. We compare our approaches with the l_{∞} pruning [6], fine-tuning (FT), fine-pruning (FP) [28] and neural attentional distillation (NAD) [24]. The number of benign samples allowed to access is set to 500 (1%) for CIFAR-10 and 5000 (5%) for Tiny-ImageNet by default.

Evaluation metrics In this work, we use the *attack clean accuracy* (ACC) and *attack success rate* (ASR) to evaluate the effectiveness of different methods. The ACC for a given model F is defined as:

$$ACC(F, \mathcal{D}_{test}) = \sum_{(\boldsymbol{x}, y) \in \mathcal{D}_{test}} \mathbb{I}\{\operatorname{argmax}(F(\boldsymbol{x})) = y\},$$

where \mathbb{I} is the *indicator function*. The ASR is defined as:

$$\mathrm{ASR}(F, \mathcal{D}_{\mathrm{test}}) = \sum_{(\boldsymbol{x}, y) \in \mathcal{D}_{\mathrm{test}}, y \neq y_t} \mathbb{I}\{\mathrm{argmax}(F(\delta(\boldsymbol{x}))) = y_t\},$$

where y_t is the attack target label. The ACC measures the model performance on benign samples, while the ASR reflects the degree of backdoor behavior retainment in the model. Given an infected model, our goal is to reduce the ASR, while keeping the ACC from dropping too much.

265 5.2 Experimental results

CIFAR-10 We show the results on CIFAR-10 in Table 1. The recently proposed NAD and ANP 266 performs significantly better than other defense methods, reducing the ASR to a very low level with a 267 slight drop on ACC. However, they also have a significant drop $(3 \sim 4\%)$ on ACC when defending 268 CLA, which is the most robust backdoor attack in our experiments, and ANP even failed when 269 defending BadNets(A2A). Nevertheless, both of our methods successfully eliminate the backdoor 270 (ASR < 1%) with negligible loss on ACC. We even observe a little rise on ACC when defending 271 BadNets by DDE. This phenomenon demonstrates that backdoor neurons may hurt the ACC in some 272 way, and thus the ACC will rise when the backdoor neurons are precisely pruned. Overall, our 273 methods achieve the most advanced defense results. 274

	the second se													
	BadNets (A2O)		BadNets (A2A)		CLA		WaNet		Blended		Refool		IAB	
Attacks	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
Origin	93.86	100.00	94.60	93.89	94.99	98.83	94.11	99.67	94.17	99.62	94.24	98.40	93.87	97.91
FT	92.22	2.16	92.03	60.76	92.88	95.73	92.93	9.37	93.9	90.27	91.68	17.78	91.78	9.52
FP	92.18	2.97	91.75	66.82	92.60	99.36	92.07	1.03	70.92	90.92	92.36	75.98	87.04	16.13
l_{∞}	92.12	100.00	93.67	6.67	92.75	98.76	93.48	99.74	86.99	99.77	91.19	98.47	88.37	88.48
NAD	93.36	2.43	92.18	2.06	91.36	15.31	93.08	3.05	92.72	1.61	91.64	6.74	92.11	19.45
ANP	93.47	3.53	90.29	86.22	91.13	11.76	94.12	0.51	93.66	5.03	91.71	26.96	93.52	10.61
DDE (Ours)	93.88	0.86	94.49	0.61	94.42	0.91	93.79	2.80	93.67	2.24	93.35	8.90	93.17	0.94
MBNS (Ours)	93.60	1.60	94.25	0.72	94.14	7.03	94.05	3.39	94.17	2.71	93.69	6.48	93.15	0.64

Table 1: Experimental results of the proposed approaches against different attacks compared with other defense methods in CIFAR-10[22].

Table 2: Experimental results of the proposed approaches against different attacks compared with other defense methods in Tiny-ImageNet[23].

-	BadNets (A2O)		CLA		WaNet		Refool		Blended		IAB		SSBA	
Attacks	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
Origin	61.36	97.38	65.61	56.58	61.47	99.98	53.26	80.61	62.85	99.83	61.4	98.28	66.51	99.78
FT	46.93	99.84	61.19	63.20	54.28	99.96	47.09	91.77	56.83	29.12	52.39	99.1	52.39	33.19
FP	35.41	99.48	62.30	39.05	53.65	100.00	42.10	86.62	59.59	99.76	52.67	98.47	53.36	31.96
l_{∞}	53.13	90.39	59.15	23.12	42.01	99.84	46.84	81.19	56.33	99.85	54.81	86.97	49.35	99.98
NAD	44.20	90.13	62.80	17.35	53.40	99.98	51.06	70.63	57.35	55.6	53.32	98.85	52.52	25.08
ANP	53.85	4.02	59.69	3.64	54.82	86.98	50.67	0.21	62.49	0.61	61.39	4.67	60.98	1.01
DDE (Ours)	60.68	0.86	64.47	0.1	60.53	0.02	51.29	17.07	60.67	0.69	61.26	0.60	64.2	0.11
MBNS (Ours)	61.60	1.60	64.86	0.05	61.58	0.01	52.41	23.79	60.77	0.85	61.30	0.60	64.64	0.01

Tiny-ImageNet Tiny-ImageNet is a larger scale dataset with higher resolution images, and it is 275 harder to defend against the attacks performed on it. Note that the A2A attack is absent, since 276 we cannot successfully perform the attack due to the large number of its classes (up to 200). Our 277 experimental results show that all of the defense methods suffer from the performance degradation 278 compared with the results in CIFAR-10, and they fail to defend against WaNet with a large ACC 279 drop but even unchanged ASR, especially the ANP and l_{∞} defense. This phenomenon shows that the 280 principles for finding backdoor neurons of both ANP and l_{∞} don't work in such case. Nevertheless, 281 our methods totally remove the backdoor and the ACC are not even affected, which indicates that our 282 methods can precisely locate the backdoor neurons even on such large scale dataset. 283

284 5.3 Ablation study

To be fair, we compare MBNS with other re-training based methods using 500 benign samples in 285 Section 5.2. However, MBNS doesn't require re-training the model, since the samples are just used 286 for detecting the distribution discrepancy. Therefore, the required samples can be much less than 500. 287 We now discuss the limit of MBNS and study how the number of samples affects the effectiveness of 288 MBNS. We train BadNets, CLA, Refool and Blended on CIFAR-10 with $\rho = 10\%$, and use 10 to 289 500 benign samples to recover the model using MBNS. We record the changes of ACC and ASR 290 with respect to the number of benign samples. The results are shown in Section 5.3. The influence 291 of number of samples to our methods comes from the randomness on estimating moments. As the 292 number of samples grows, the randomness is reduced and MBNS has more stable performance, but 293 the average performances are not improved, except for Refool. Compared with other attacks, Refool 294 clearly needs more samples to reduce the ASR. The reason may be that the mixture distribution in 295 Refool has closer moments and is harder to distinguish. Besides, we surprisingly find that MBNS can 296 recover BadNets, CLA and Blended using only 10 benign samples. The additional results on DDE 297 are shown in ??. 298

We also conduct experiments to show the high correlation between the backdoor neurons and our proposed evaluation metrics, the results are shown in **??**

301 6 Discussion

³⁰² The proposed methods are superior to other existing defense methods in the following three aspects:

Better performance As demonstrated in Section 5, both of the proposed methods achieve stateof-the-art results. Moreover, according to the ablation study, the proposed MBNS can successfully



ples used for MBNS.



defend most of the attacks within 10 benign samples, which shows the amazing effectiveness of our 305 proposed methods. 306

Higher efficiency The proposed methods are also highly efficient. We record the running time of 307 several defense methods on 500 CIFAR-10 images with ResNet-18, and show the results in Table 3. 308 It can be seen that both of the proposed methods require less time than the baseline defense methods. 309 Since both methods require scanning on each neuron once, the computational complexity scale 310 linearly with the number of the neurons in the neural network. Therefore, the efficiency of our 311 methods is promised. 312

Table 3: The overall running time of different defense methods on 500 CIFAR-10 images with ResNet-18.

Defense Method	FT	FP	NAD	ANP	DDE (ours)	MBNS (ours)
Runing Time (sec.)	12.35s	14.59s	22.08s	25.68s	10.69s	0.39s

More robust to hyperparameter choosing One of the most general problems in backdoor defense 313 is the choice of hyperparameters. Under realistic settings, the defenders can only perform defenses 314 without any prior knowledge about poison data, including poisoning rate and examples of poisoned 315 data. So the defenders should carefully tune the hyperparameters, or the ACC and ASR can change 316 suddenly even under small fluctuations of those hyperparameters. In comparison, both of the 317 proposed pruning strategies only require one universal hyperparameter u. Moreover, they show 318 reliable consistency against different attacks in the same dataset, only vary from different datasets, 319 which is inevitable. Besides, we leave a wide range of parameters to choose, so that the ACC remains 320 high while the ASR is controlled to a very small number, as shown in Section 5.3. 321

Conclusion 7 322

323 In this work, we take a deep inspection on the pre-activation distributions of each layer, and study their characteristics in a poisoned model. We find that the distributions of benign data and poisoned 324 data on these neurons are both Gaussian but extraordinarily different in their moments. The mixture 325 distributions have low entropy compared with that of the benign neurons. This property makes it 326 possible for the defender to efficiently detect the abnormal neurons based on the sample entropy, when 327 the poisoned dataset can be accessed. Moreover, when the defender has only a small set of benign 328 data, we propose to detect the abnormal neurons based on the inconsistency of the BN statistics on 329 poisoned dataset and the sample statistics on the given benign dataset. We then do pruning on these 330 potential backdoor neurons using the proposed two detection methods to recover the model. The 331 experiments show that the proposed defending strategies can efficiently locate the backdoor neurons, 332 and greatly reduce the backdoor threat with negligible loss on clean accuracy. Our approach based on 333 the pre-activation distributions achieves superior results compared with all other defense methods 334 under various attacks on both datasets, empirically showing that the distribution property may be 335 both dataset-invariant and attack-invariant. The results shed lights on the field of backdoor defense, 336 and can be a guidance for designing more robust backdoor attacks. 337

338 **References**

- [1] Hassan Ali, Surya Nepal, Salil S Kanhere, and Sanjay Jha. Has-nets: A heal and select
 mechanism to defend dnns against backdoor attacks for data collection scenarios. *arXiv preprint arXiv:2012.07474*, 2020.
- [2] Eitan Borgnia, Valeriia Cherepanova, Liam Fowl, Amin Ghiasi, Jonas Geiping, Micah Goldblum,
 Tom Goldstein, and Arjun Gupta. Strong data augmentation sanitizes poisoning and backdoor
 attacks without an accuracy tradeoff. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3855–3859. IEEE, 2021.
- [3] Richard C Bradley Jr. Central limit theorems under weak dependence. *Journal of Multivariate Analysis*, 11(1):1–16, 1981.
- [4] Huili Chen, Cheng Fu, Jishen Zhao, and Farinaz Koushanfar. Proflip: Targeted trojan attack with
 progressive bit flips. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7718–7727, 2021.
- [5] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- [6] Hao Cheng, Kaidi Xu, Sijia Liu, Pin-Yu Chen, Pu Zhao, and Xue Lin. Defending against
 backdoor attack on deep neural networks. *CoRR*, abs/2002.12162, 2020.
- [7] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- [8] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and
 Anil A Bharath. Generative adversarial networks: An overview. *IEEE Signal Processing Magazine*, 35(1):53–65, 2018.
- [9] Imre Csiszár. I-divergence geometry of probability distributions and minimization problems.
 The annals of probability, pages 146–158, 1975.
- [10] Bao Gia Doan, Ehsan Abbasnejad, and Damith C. Ranasinghe. Februus: Input purification
 defense against trojan attacks on deep neural network systems. In *Annual Computer Security Applications Conference*, ACSAC '20, page 897–912, New York, NY, USA, 2020. Association
 for Computing Machinery.
- [11] Min Du, Ruoxi Jia, and Dawn Song. Robust anomaly detection and backdoor attack detection
 via differential privacy. In *International Conference on Learning Representations*, 2019.
- [12] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal.
 Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th Annual Computer Security Applications Conference*, pages 113–125, 2019.
- [13] Chengyue Gong, Tongzheng Ren, Mao Ye, and Qiang Liu. Maxup: A simple way to improve generalization of neural network training. *arXiv preprint arXiv:2002.09024*, 2020.
- ³⁷³ [14] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating ³⁷⁴ backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- Ionathan Hayase and Weihao Kong. Spectre: Defending against backdoor attacks using robust
 covariance estimation. In *International Conference on Machine Learning*, 2020.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
 pages 770–778, 2016.
- [17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
 arXiv preprint arXiv:1503.02531, 2015.
- [18] Sanghyun Hong, Varun Chandrasekaran, Yigitcan Kaya, Tudor Dumitras, and Nicolas Paper not. On the effectiveness of mitigating data poisoning attacks with gradient shaping. *CoRR*,
 abs/2002.11497, 2020.

- [19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training
 by reducing internal covariate shift. In *International conference on machine learning*, pages
 448–456. PMLR, 2015.
- [20] Alexander B. Jung, Kentaro Wada, Jon Crall, Satoshi Tanaka, Jake Graving, Christoph Reinders,
 Sarthak Yadav, Joy Banerjee, Gábor Vecsei, Adam Kraft, Zheng Rui, Jirka Borovec, Christian
 Vallentin, Semen Zhydenko, Kilian Pfeiffer, Ben Cook, Ismael Fernández, François-Michel
 De Rainville, Chi-Hung Weng, Abner Ayala-Acevedo, Raphael Meudec, Matias Laporte, et al.
 imgaug. https://github.com/aleju/imgaug, 2020. Online; accessed 01-Feb-2020.
- ³⁹³ [21] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing ³⁹⁴ neural networks. *Advances in neural information processing systems*, 30, 2017.
- ³⁹⁵ [22] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- [23] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
- Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention
 distillation: Erasing backdoor triggers from deep neural networks. In *International Conference on Learning Representations*, 2020.
- Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor
 learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems*, 34, 2021.
- Yiming Li, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A
 survey. *arXiv preprint arXiv:2007.08745*, 2020.
- [27] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor
 attack with sample-specific triggers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16463–16472, 2021.
- [28] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against
 backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*, pages 273–294. Springer, 2018.
- [29] 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 *Proceedings* of the 2019 ACM SIGSAC Conference on Computer and Communications Security, pages
 1265–1282, 2019.
- [30] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and
 Xiangyu Zhang. Trojaning attack on neural networks. In *25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018.*The Internet Society, 2018.
- [31] Yunfei Liu, Xingjun Ma, James Bailey, and Feng Lu. Reflection backdoor: A natural backdoor
 attack on deep neural networks. In *European Conference on Computer Vision*, pages 182–199.
 Springer, 2020.
- [32] Tuan Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. *Advances in Neural Information Processing Systems*, 33:3454–3464, 2020.
- [33] Tuan Anh Nguyen and Anh Tuan Tran. Wanet-imperceptible warping-based backdoor attack.
 In *International Conference on Learning Representations*, 2020.
- [34] Adnan Siraj Rakin, Zhezhi He, and Deliang Fan. Tbt: Targeted neural network attack with
 bit trojan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13198–13207, 2020.
- [35] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.

- [36] Guangyu Shen, Yingqi Liu, Guanhong Tao, Shengwei An, Qiuling Xu, Siyuan Cheng, Shiqing
 Ma, and Xiangyu Zhang. Backdoor scanning for deep neural networks through k-arm optimiza tion. *arXiv preprint arXiv:2102.05123*, 2021.
- [37] Brandon Tran, Jerry Li, and Aleksander Mądry. Spectral signatures in backdoor attacks. In
 Proceedings of the 32nd International Conference on Neural Information Processing Systems,
 pages 8011–8021, 2018.
- [38] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks.
 arXiv preprint arXiv:1912.02771, 2019.
- [39] 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)*, pages 707–723. IEEE, 2019.
- [40] Dongxian Wu and Yisen Wang. Adversarial neuron pruning purifies backdoored deep models.
 Advances in Neural Information Processing Systems, 34, 2021.
- [41] Pengfei Xia, Hongjing Niu, Ziqiang Li, and Bin Li. Enhancing backdoor attacks with multi-level
 mmd regularization. *IEEE Transactions on Dependable and Secure Computing*, pages 1–1,
 2022.
- [42] M. Xue, C. He, J. Wang, and W. Liu. One-to-n n-to-one: Two advanced backdoor attacks
 against deep learning models. *IEEE Transactions on Dependable and Secure Computing*, (01):1–1, oct 5555.
- [43] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon
 Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In
 Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 6023–6032,
 2019.
- [44] Pu Zhao, Pin-Yu Chen, Payel Das, Karthikeyan Natesan Ramamurthy, and Xue Lin. Bridging
 mode connectivity in loss landscapes and adversarial robustness. In *International Conference on Learning Representations (ICLR 2020)*, 2020.

457 Checklist

- The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or [N/A]. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:
- Did you include the license to the code and datasets? [Yes]
- Did you include the license to the code and datasets? [No] The code and the data are proprietary.
- Did you include the license to the code and datasets? [N/A]

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...

470	(a)	Do the main claims made in the abstract and introduction accurately reflect the paper's
471		contributions and scope? [Yes] See abstract.
472	(b)	Did you describe the limitations of your work? [Yes] See 5.3
473	(c)	Did you discuss any potential negative societal impacts of your work? [N/A] There is
474		no negative social impacts of our work.
475	(d)	Have you read the ethics review guidelines and ensured that your paper conforms to
476		them? [Yes]

477	2. If you are including theoretical results
478	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
479	(b) Did you include complete proofs of all theoretical results? [Yes]
480	3. If you ran experiments
481 482	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes]
483 484	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See experimental settings.
485 486	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] See experimental results.
487 488	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See experimental settings.
489	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
490	(a) If your work uses existing assets, did you cite the creators? [Yes]
491	(b) Did you mention the license of the assets? [Yes]
492	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
493 494	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes]
495 496	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
497	5. If you used crowdsourcing or conducted research with human subjects
498 499	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
500 501	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
502 503	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

504 A Implementation Details of Performed Backdoor Attacks and Defenses

⁵⁰⁵ In this part, we provide more detailed information about the state-of-art attacks and defenses applied ⁵⁰⁶ in our work. The examples of poisoned images are shown in Appendix A.



(d) Refool

(e) CLA

(f) IAB

Figure 5: Examples of the poisoned data on CIFAR-10.



(a) Original image

(b) SSBA

Figure 6: Examples of the poisoned data on Tiny-ImageNet

The hyperparameter we used in our DDE and MBNS is 3 in CIFAR-10 and 4 in Tiny-ImageNet. The parameters in other defenses is set as default.

509 B Additional Experimental Results of Different Poisoning Rate

Poisoning rate is considered important in the proposed methods, since smaller poisoning rate will make the mixture distributions less distinguishable from unimodal distributions. In DDE, decreasing poisoning rate will increase the entropy of the backdoor neurons, and it will be harder to precisely

detect them. In MBNS, the smaller poisoning rate makes the BN statistics closer to the true benign 513 statistics. Hence, it is important to study the influence of poisoning rate on the performance of the 514 proposed methods. We set the poisoning rate to 1%, 5% and 10% for all attacks and show the results 515 in Table 4. Note that CLA fails to attack the model when the poisoning rate is set to 1%. In most 516 cases, their performances retain a high level. Nevertheless, there are severe degradation against CLA 517 and Refool. Specifically, the ASR of CLA and Refool remain over 10% after pruning using DDE 518 and MBNS. However, we use default hyperparameter, i.e., $u_h = 3$ and $u_K = 3$ here. By further 519 decreasing the pruning threshold, we can reduce both ASR to less than 5%. Overall, MBNS is more 520 robust against attacks of small poisoning rate.

Table 4: Experimental results of the proposed approaches against different attacks compared with other defense methods in CIFAR-10[22].

		BadNets (A2O)		BadNets (A2A)		CLA		WaNet		Blended		Refool		IAB	
ρ	Stage	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
	Origin	95.03	99.94	94.75	88.57	88.96	4.73	94.76	46.82	94.17	99.62	93.08	99.59	93.22	64.00
1%	DDE	94.82	0.91	94.17	0.73	87.96	0.65	93.67	14.17	94.46	2.29	91.77	24.08	92.73	5.21
	MBNS	92.22	2.16	93.16	7.99	88.03	0.84	94.64	1.24	92.89	2.44	90.99	21.22	93.17	4.19
	Origin	94.29	99.99	94.26	92.78	95.53	92.23	94.00	94.55	94.53	81.33	94.35	97.98	92.70	65.50
5%	DDE	93.83	0.83	93.67	0.70	94.43	15.91	92.73	10.13	94.44	5.49	92.75	4.51	92.29	1.83
	MBNS	93.61	0.67	93.99	5.64	94.65	12.06	94.17	1.78	93.37	9.21	92.30	2.08	92.74	2.14
	Origin	93.89	100.00	94.60	93.89	94.99	98.83	94.11	99.67	94.17	99.63	94.24	98.40	93.87	97.91
10%	DDE	93.88	0.86	94.49	0.61	94.42	0.91	93.79	2.80	93.67	2.24	93.35	8.90	93.17	0.94
	MBNS	93.60	1.60	94.25	0.72	94.14	7.03	94.05	3.39	94.17	2.71	93.69	6.48	93.15	0.64

521

522 C Additional Experimental Results on Wide-ResNet

To study the generalization of the proposed defense methods over different architectures, we conduct experiments on WideResNet-28-1, which is also a commonly used model in backdoor community. The results are shown in Appendix C.

	Backd	loored	DI	DE	MBNS		
Attacks	ACC	ASR	ACC	ASR	ACC	ASR	
BadNets (A2O)	91.62	99.99	90.30	1.71	91.37	1.61	
BadNets (A2A)	92.53	92.03	91.50	0.90	91.43	1.56	
CLA	92.81	80.14	91.55	10.24	91.91	4.30	
Refool	91.53	97.28	91.11	0.74	90.29	2.87	
WaNet	91.78	91.92	92.28	0.67	91.84	0.56	
Blended	91.65	99.78	91.59	6.40	91.59	1.52	
IAB	90.32	88.08	90.97	1.87	90.74	1.86	

525

526 D Analysis of Correlation Between Backdoor Neurons and Pruned Neurons

In section Section 3.3, we have demonstrated the definition of sensitivity of neurons to backdoor in 527 Equation (2), which is the difference of backdoor loss before and after being pruned. It is important 528 to ensure that our methods precisely detect those backdoor neurons with high sensitivity. Here, we 529 make scatter plots of the backdoor sensitivity of the neurons and their differential entropy and KL 530 divergence in some typical layers. The detailed scatter figures are shown in Figure 7. The results 531 indicate that both DDE and MBNS choose the neurons with high sensitivity. Especially for BadNets, 532 the highest sensitivity of one single neuron is extremely large and its differential entropy and KL 533 divergence stay away from average. The overall results demonstrate a strong correlation between the 534 neuron sensitivity to backdoor and the proposed indices. 535

536 E Performance on Low-quality Data

Our methods require extra data and rely on its distribution to detoxify the backdoor model, so it's natural to consider a worse case that we can only get low-quality dataset with some distortion in it. To further show the effectiveness of our proposed methods, we utilize multiple data augmentation tricks



Figure 7: Scatter plots of backdoor sensitivity of mid-layer neurons to backdoor and the corresponding differential entropy and KL divergence indices.



Figure 8: Examples of the perturbed data.

to create out-of-distribution data on CIFAR-10 while performing defenses against BadNets (A2O) 540 attack, including color jitter with 90 degree rotation (CJ) from torchvision, CoarseSaltandPepper 541 (SAP), PolarWarping (PW), Snowflakes (Snow) from [20]. To be specific, the parameters of CJ are 542 set to be [0.5, 0.5, 0.5, 0.5]; 30% of the pixels in SAP are replaced replaced by salt/pepper noise mask 543 which has 1% to 10% the size of the input image; the translate percent in PW is set to be ± 0.2 ; the 544 snow size of Snow is (0.2, 0.5) and its speed is (0.01, 0.05). Samples of those perturbed images are 545 546 shown in Appendix E. Now we replace 10% of the DDE and MBNS used data into those perturbed out-of-distribution data, and see the performance of our methods. The experimental results are shown 547 in Table 5, which illustrates that even if some low-quality data are mixed into the data they use, our 548 methods can still successfully locate the backdoor neurons and delete them while maintaining high 549 clean accuracy. 550

Table 5: Experimental Results of the proposed methods using low-quality datasets.

	Origin		CJ		SAP		PW		Snow	
Attacks	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR
DDE	93.88	0.86	93.62	0.88	93.19	0.86	93.65	0.82	93.23	0.93
MBNS	93.60	1.60	93.56	1.52	91.17	1.58	92.62	0.96	90.99	0.95

551 E.1 Exploration of Potential Adaptive Attacks

⁵⁵² In this section, we explore some potential adaptive attacks against the proposed methods.

553 E.1.1 BN Re-parameterization against MBNS

The MBNS pruning strategy rely on the correct information from the BN statistics. However, this information can be changed by a defense-aware attacker using re-parameterization techniques. 556 Specifically, consider a BN operation performed on each neuron by:

$$\tilde{BN}^{(l)}(\boldsymbol{x}_{k}^{(l)})_{\tilde{\gamma}_{k}^{(l)},\tilde{\beta}_{k}^{(l)}} = \tilde{\gamma}_{k}^{(l)} \frac{\boldsymbol{x}_{k}^{(l)} - \tilde{\mu}_{k}^{(l)}}{\tilde{\sigma}_{k}^{(l)}} + \tilde{\beta}_{k}^{(l)}$$
$$= \frac{\tilde{\gamma}_{k}^{(l)}}{\tilde{\sigma}_{k}^{(l)}} \boldsymbol{x}_{k}^{(l)} + (\tilde{\beta}_{k}^{(l)} - \frac{\tilde{\mu}_{k}^{(l)}}{\tilde{\sigma}_{k}^{(l)}}).$$

The attacker can calculate the sample statistics denoting $\mu_k^{(l)}$ and $\sigma_k^{(l)}$ from benign data, and assign new weights and bias by:

$$\begin{split} \gamma_k^{(l)} &= \frac{\sigma_k^{(l)}}{\tilde{\sigma}_k^{(l)}} \tilde{\gamma}_k^{(l)} \\ \beta_k^{(l)} &= \tilde{\beta}_k^{(l)} - \frac{\tilde{\mu}_k^{(l)}}{\tilde{\sigma}_k^{(l)}} + \frac{\mu_k^{(l)}}{\sigma_k^{(l)}} \end{split}$$

In this way, the whole linear transformation remains unchanged, but the BN statistics become consistent with that of the benign data, which makes the backdoor neurons unable to be detected by

the proposed MBNS.

	$\lambda =$	= 1	$\lambda =$	0.1	$\lambda =$	0.01	$\lambda = 0.001$		
Attacks	ACC	ASR	ACC	ASR	ACC	ASR	ACC	ASR	
Origin	NaN	NaN	NaN	NaN	95.11	100.00	93.41	75.30	
DDE	—		—		93.63	9.23	93.35	0.63	
MBNS	—	—			94.15	0.32	93.82	0.24	
TT 1 1	(D	•	1 1.			1 1 1		. 1	

Table 6: Experimental results on regularization-based adaptive attack.

562 E.1.2 Regularization-based Adaptive Attacks

Both of the two methods rely on the discrepancy of the benign distribution and poisoned distribution.
Hence, we wonder whether DDE and MBNS still work if the attacker try to regularize the distribution
to minimize the discrepancy. Hence, we use BadNets as an example to perform the adaptive attack
with an additional objective:

$$\mathcal{L}_{adaptive} = \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}}[D_{CE}(\boldsymbol{y}, f(\boldsymbol{\delta}(\boldsymbol{x}))] + \lambda \sum_{1 \le k \le K} \sum_{1 \le l \le L} \mathbb{E}_{\boldsymbol{x} \sim \mathcal{X}}[(\boldsymbol{\phi}_k^{(l)} - \hat{\boldsymbol{\phi}}_k^{(l)})^2],$$
(5)

where $\phi_k^{(l)}$ and $\hat{\phi}_k^{(l)}$ denote the pre-activations of benign and poisoned samples in the k^{th} neuron of the l^{th} layer. The goal of the second term is to minimize the discrepancy between the benign 567 568 distribution and the poisoned distribution with a trade-off hyperparameter λ . We select λ from [1.000, 569 0.100, 0.010, 0.001] and train four models and do pruning using DDE and MBNS, the results are 570 shown in Table 6. We find that too large λ makes the objective not trainable, hence collapse the 571 training process. This may because it is not realistic to distinguish benign and poisoned samples 572 while making their distribution indistinguishable. After down weighting λ , the model is able to be 573 trained normally. However, we find DDE and MBNS are still able to remove the backdoor without 574 575 influence too much on the clean accuracy.

A recently proposed method[41] also regularizes the distribution discrepancy between benign samples
and poisoned samples by maximum mean discrepancy (MMD), which is directly against our defense.
We test our method on it and find it still works. Specifically, the original ACC and ASR are 87.51%,
98.48%, respectively. The ACC and ASR after pruning by DDE is 86.34%, 3.26% and ACC and
ASR for MBNS is 86.93%, 6.15%, respectively.