

REFINE: INVERSION-FREE BACKDOOR DEFENSE VIA MODEL REPROGRAMMING

Anonymous authors

Paper under double-blind review

ABSTRACT

Backdoor attacks on deep neural networks (DNNs) have emerged as a significant security threat, allowing adversaries to implant hidden malicious behaviors during the model training phase. Pre-processing-based defense, which is one of the most important defense paradigms, typically focuses on input transformations or backdoor trigger inversion (BTI) to deactivate or eliminate embedded backdoor triggers during the inference process. However, these methods suffer from inherent limitations: transformation-based defenses often struggle to balance the intensity of transformations with preserving the model’s accuracy, while BTI-based defenses require accurate reconstruction of the trigger patterns, which is rarely achievable without prior knowledge. In this paper, we propose REFINE, an inversion-free backdoor defense method based on model reprogramming. REFINE consists of two key components: **(1)** an input transformation module that disrupts both benign and backdoor patterns, generating new benign features; and **(2)** an output remapping module that redefines the model’s output domain to guide the input transformations effectively. By further integrating supervised contrastive loss, REFINE enhances the defense capabilities while maintaining model utility. Extensive experiments on various benchmark datasets demonstrate the effectiveness of our REFINE and its resistance to potential adaptive attacks.

1 INTRODUCTION

Deep neural networks (DNNs) have been widely deployed across various domains (He et al., 2023; Liu et al., 2024; He et al., 2024; Zhang et al., 2024). To develop a high-performance DNN, developers necessitate not only high-quality data samples but also substantial computational resources. Consequently, developers frequently and directly rely on third-party models for follow-up development. However, the utilization of third-party DNNs can introduce security threats, particularly with regard to backdoor attacks (Gu et al., 2019; Li et al., 2022c; Dong et al., 2023; Yang et al., 2024a).

Backdoor attacks aim to implant hidden backdoors into the model during training (Gu et al., 2019). After the attack, the backdoored model functions normally on benign inputs. However, when a specific trigger is present, the model will produce intentionally incorrect outputs. Backdoor attacks pose a severe threat to critical applications where model reliability is essential, highlighting the urgent need for effective backdoor defense strategies to safeguard AI systems (Li et al., 2022b).

Currently, several backdoor defenses (Huang et al., 2022; Li et al., 2024a;b; Hou et al., 2024) have been developed to tackle the threat of backdoor attacks. Among these, pre-processing-based defenses (Villarreal-Vasquez & Bhargava, 2020; Qiu et al., 2021) are particularly notable because they only apply certain modifications to input samples before model inference, without altering the original model structure and weights. Currently, there are two main types of pre-processing-based defenses. The first type of defense relies on input transformations (Liu et al., 2017; Li et al., 2021c; Sun et al., 2023). These defenses aim to mismatch or eliminate potential trigger patterns by performing certain transformations to the input samples. The second type is based on backdoor trigger inversion (BTI) (Wang et al., 2019; 2022b; Xu et al., 2024), which attempts to reconstruct the attacker’s trigger patterns and remove them before the data is processed by the model.

In this paper, we revisit the aforementioned pre-processing-based backdoor defenses. We reveal that they both have intrinsic limitations. Specifically, transformation-based defenses face a trade-off between the model utility and the defense effectiveness: more extensive transformations can achieve

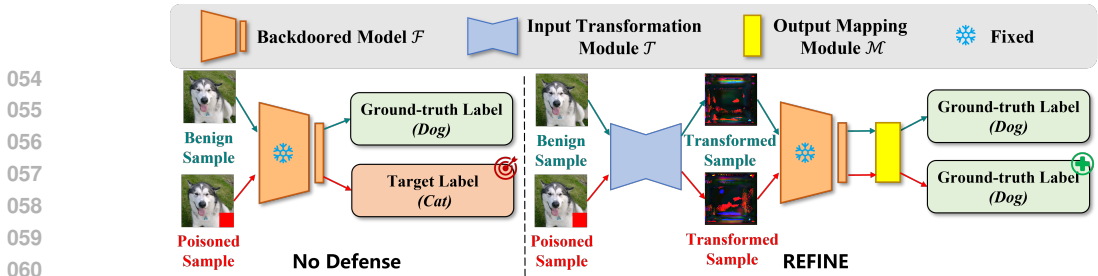


Figure 1: The defense process of our REFINE. The label remapping in the model’s output domain significantly enhances the flexibility of input transformations while maintaining consistent sample predictions, effectively mitigating the trade-off often encountered in transformation-based pre-processing defenses. During prediction, the input sequentially passes through the well-trained input transformation module, the fixed backdoored model, and the pre-defined output mapping module, ultimately yielding the expected ground-truth (instead of the malicious target) label.

lower attack success rates but may negatively impact the model’s benign accuracy. This is because they lack information about the backdoor-related features in the input. Thus, they can only indirectly modify these features by altering all the features, including the benign ones that are closely related to the benign accuracy, in the input domain to prevent backdoor activation. Since benign features are tightly coupled with backdoor features, it is challenging for the defender to apply more extensive transformations without compromising the benign accuracy of the original model. On the other hand, BTI-based defenses can break the trade-off by first obtaining the information of backdoor triggers via trigger inversion. However, due to the inherent difficulties of BTI (e.g., lack of prior knowledge about the implanted backdoor and poisoned samples), existing BTI methods struggle to invert the ground-truth trigger. This limitation makes it difficult to purify the backdoor input from the poisoned domain to the benign domain, causing BTI-based defenses to fall short of achieving the desired performance. Accordingly, an intriguing and important question arises: *Could we break the curse of this trade-off without relying on backdoor trigger inversion?*

The answer to the above question is positive! We first provide a theoretical analysis showing that the effect of backdoor defenses is bounded by the distance of the output features before and after the pre-processing. Accordingly, the ineffectiveness of existing defenses is mostly due to their underlying assumption of having a fixed output domain. Based on the above understandings, inspired by model reprogramming (Chen, 2024), we propose REFINE, a REprogramming-based Inversion-Free backdoor defense method, as shown in Figure 1. REFINE can significantly alter the input domain while preserving the model’s accuracy to a large extent for it allows changing the output domain. Specifically, our REFINE involves an input transformation module and an output mapping module to reprogram the backdoored model and eliminate the backdoor. We utilize a trainable autoencoder as the input transformation module and redefine the model’s output domain through a hard-coded remapping function. Due to the changes in the model’s output domain, we can implement more extensive and effective transformations on the input samples. Besides, we further improve our method by imposing constraints on the transformed samples using supervised contrastive loss (Khosla et al., 2020). This ensures that samples of the same class remain more similar after transformation.

Our contributions are three-fold. (1) We revisit existing pre-processing-based backdoor defenses and reveal their limitations. (2) Based on the empirical and theoretical analysis, we propose a simple yet effective defense (i.e., REFINE). Our REFINE introduces trainable input transformation and output mapping modules for reprogramming and incorporates cross-entropy and supervised contrastive losses to enhance defense performance. (3) Extensive experiments on diverse benchmark datasets demonstrate the effectiveness of REFINE and its resistance to potential adaptive attacks.

2 BACKGROUND

2.1 BACKDOOR ATTACKS

Backdoor attacks (Gao et al., 2020; Li et al., 2022b) involve embedding hidden malicious behaviors into a model, typically by manipulating the training process with a small subset of poisoned data containing adversary-specified trigger patterns. Whenever the trigger appears in the input during inference, the model executes the attacker’s intended behavior, such as misclassifying the input to

108 a target label. In the absence of the trigger, the model functions normally, rendering the backdoor
 109 hard to detect. Backdoor attacks pose serious threats in AI-empowered systems.

110
 111 The formulation of backdoor attacks is typically presented as follows. Given a training dataset
 112 $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, the attacker manipulates the training process of the model \mathcal{F} by introducing a
 113 poisoned subset $\tilde{\mathcal{D}} = \{(\tilde{\mathbf{x}}_i, \mathbf{y}_t)\}_{i=1}^M$, where $\tilde{\mathbf{x}}_i = G(\mathbf{x}_i)$ with $G(\cdot)$ as a certain trigger injection func-
 114 tion and \mathbf{y}_t being the chosen target label, or by altering the training loss directly. During inference,
 115 the model behaves normally on benign samples, where $\mathbf{y}_j = \mathcal{F}(\mathbf{x}_j)$, while exhibiting backdoor
 116 behavior on poisoned samples, such as misclassifying to the target label $\mathbf{y}_t = \mathcal{F}(G(\mathbf{x}_j))$.

117 Generally, existing attacks can be classified into two types: (1) *Visible backdoor attacks*, which
 118 typically employ trigger patterns that are visible to humans, such as specific white-black squares (Gu
 119 et al., 2019), physical attacks (Li et al., 2021c), or adaptive attacks (Qi et al., 2023). (2) *Invisible*
 120 *backdoor attacks*, which introduce the usage of triggers that are imperceptible to humans to enhance
 121 the stealth and evasion of the attacks (Chen et al., 2017), including sample-specific attacks (Nguyen
 122 & Tran, 2021; Li et al., 2021d), trainable noise attacks (Doan et al., 2021), and sample rotation
 123 attacks (Xu et al., 2023)). More details are in Appendix I.

124 2.2 BACKDOOR DEFENSES

125
 126 Currently, there are various backdoor defense methods designed to mitigate backdoor threats. These
 127 methods can generally be divided into three main paradigms (Li et al., 2022b): (1) *pre-processing-*
 128 *based defenses* (Liu et al., 2017; Li et al., 2021c; Shi et al., 2023). (2) *backdoor elimination* (Li et al.,
 129 2021b; Huang et al., 2022; Xu et al., 2024), which involves adjusting model parameters through
 130 fine-tuning, pruning or reconstruction to remove the backdoor. (3) *trigger elimination*, also known
 131 as testing sample filtering (Gao et al., 2019; Javaheripi et al., 2020; Li et al., 2023a). In this paper,
 132 we focus on pre-processing-based defenses since we consider scenarios where only fixed third-party
 133 models are accessible and defenders require to obtain the correct final results of all samples. Detailed
 134 discussion about the backdoor defenses can be found in Appendix I.

135 **Pre-processing-based Defenses.** Generally, pre-processing-based defenses can be categorized into
 136 two types: (1) *Transformation-based defenses*. Classical methods (Liu et al., 2017; Li et al., 2021c;
 137 Qiu et al., 2021) typically involve applying simple transformations to input, aiming to disrupt trigger
 138 patterns and prevent the model from exhibiting backdoor behavior. More Recently, many methods
 139 have leveraged the powerful reconstruction capabilities of generative models, such as diffusion mod-
 140 els (Shi et al., 2023; May et al., 2023; Zhou et al., 2024) and masked autoencoders (Sun et al., 2023),
 141 intending to retain the original benign features while minimizing the presence of backdoor-related
 142 features. However, there is a trade-off between removing backdoor patterns and restoring benign
 143 patterns, which remains a pressing issue to address. (2) *BTI-based defenses* (Wang et al., 2019; Xu
 144 et al., 2024; Wang et al., 2023), which focus on inverting the triggers employed by the attacker and
 145 utilizing them to purify the input samples. However, these methods may face issues with inaccura-
 146 cies in the inverted triggers, which may lead to suboptimal purification of the input. How to design
 147 an effective pre-processing-based defense is still an important open question.

148 2.3 MODEL REPROGRAMMING

149
 150 Model reprogramming (Kloberdanz et al., 2021; Neekhara et al., 2022; Jing et al., 2023) is a tech-
 151 nique that extends the application of a pre-trained model from a source domain to a target domain.
 152 This technique involves adapting the input from the target domain to match that of the source do-
 153 main. Specifically, model reprogramming introduces an input transformation module $\mathcal{T}(\mathbf{x}|\boldsymbol{\theta})$ and
 154 an output mapping module $\mathcal{M}(\mathbf{y}|\boldsymbol{\beta})$, where $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$ are the trainable parameters of these two mod-
 155 ules, respectively (Chen, 2024). Given a pre-trained model $\mathcal{F}(\cdot)$ and an input sample \mathbf{x} , model
 156 reprogramming first transforms \mathbf{x} to $\tilde{\mathbf{x}}$ leveraging the input transformation module. Then input $\tilde{\mathbf{x}}$
 157 into the pre-trained model $\mathcal{F}(\cdot)$ and get the output $\tilde{\mathbf{y}} = \mathcal{F}(\tilde{\mathbf{x}})$. Finally, the output mapping module is
 158 used to map $\tilde{\mathbf{y}}$ into the final output \mathbf{y} . Through fine-tuning the input transformation module and the
 159 output mapping module (*i.e.*, optimizing $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$), model reprogramming can efficiently turn the
 160 pre-trained model from the source domain to a target domain. Compared to transfer learning, model
 161 reprogramming does not necessitate modifying the parameters of the pre-trained model. As such, it
 is more efficient and flexible. Detailed descriptions of model reprogramming are in Appendix I.

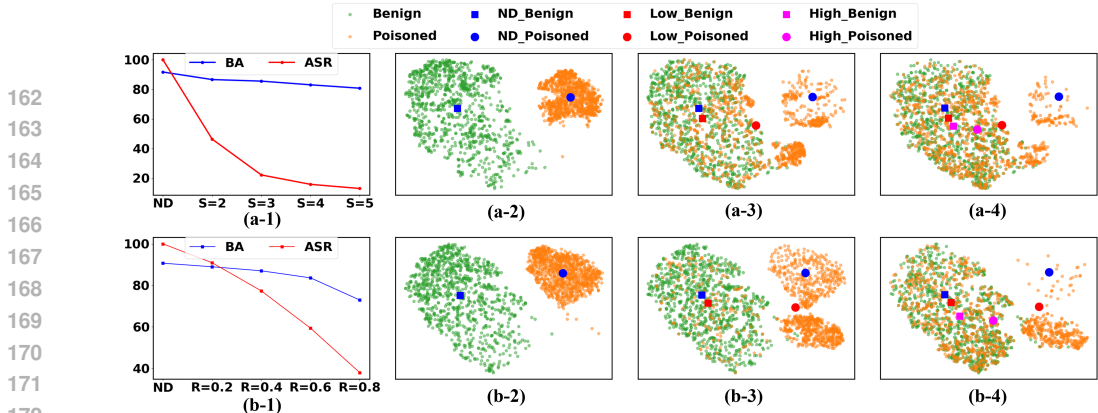


Figure 2: (a-1)&(b-1): The ASR and BA for ShrinkPad (the first row) and BDMAE (the second row) with different transformation intensities. (a-2)~(a-4)&(b-2)~(b-4): The t-SNE plots of the features of benign and backdoor samples under no defense (dubbed “ND”), low transformation intensity (dubbed “Low”), and high transformation intensity (dubbed “High”). Squares and solid circles represent the centroids of benign sample distributions and backdoor sample distributions. As the transformation intensity increases, the features of benign samples deviate from the origin. The results demonstrate the tradeoff faced by the transformation-based backdoor defense methods.

3 REVISITING EXISTING PRE-PROCESSING-BASED BACKDOOR DEFENSES

3.1 THREAT MODEL

This paper focuses on tackling the issue of pre-trained backdoored models via pre-processing-based backdoor defense. The defender may buy or acquire a pre-trained model from third-party platforms. However, there exists a threat that the pre-trained model is backdoored. Due to the limitations of computational resources, the defender seeks to mitigate the backdoor in an efficient and low-cost way (*e.g.*, without altering the parameters of the pre-trained model). Following prior works (Liu et al., 2017; Li et al., 2021c), we make the following assumptions. For adversaries, they can implant the backdoor into the pre-trained model in any way (*e.g.*, by poisoning the training data or intervening in the training process). For defenders, we assume that they have access to an *unlabeled* dataset that is independent and identically distributed to the training dataset of the pre-trained model.

3.2 THE LIMITATIONS OF TRANSFORMATION-BASED DEFENSES

Transformation-based defenses aim to mismatch or eliminate triggers by applying specific transformations to test samples. This type of defense method can be categorized into two types: random perturbations and generator reconstruction. Specifically, random perturbations involve the defender mismatching the trigger pattern through techniques such as scaling or rotation, while generator reconstruction leverages a pre-trained generative model to erase the trigger pattern. However, *the transformation-based backdoor defense methods face a trade-off between the utility of the model and the effectiveness of the backdoor elimination*, making them ineffective in practice.

In this section, we present the empirical results to support the above claim. We implement two representative transformation-based methods, ShrinkPad (Li et al., 2021c) (dubbed “SP”) and BDMAE (Sun et al., 2023) (dubbed “BD”), to defend the BadNets attack (Gu et al., 2019) on CIFAR-10. Specifically, ShrinkPad applies simple spatial transformations to the input, while BDMAE employs a trained masked autoencoder for data cleansing. We use “Pad Size” (dubbed “S”), which refers to the padding size applied around the shrunk image, and “Mask Ratio” (dubbed “R”), which represents the masking rate applied to the image before reconstruction, to control the transformation intensity for ShrinkPad and BDMAE, respectively. We aim to analyze how these transformations impact the model’s benign accuracy (BA) and attack success rate (ASR) of the backdoor. Additionally, we treat the original model as a feature extractor. We then visualize how transformation intensity affects the differences in feature distribution between benign and poisoned samples of the same class, both before and after the transformation.

As shown in Figure 2 (a-1) and (b-1), increasing the intensity of input transformation, which enlarges the feature distance between the original and transformed samples, reduces the backdoor ASR. However, it also leads to a decline in the model’s BA. As depicted in Figure 2 (a-2)~(a-4) and

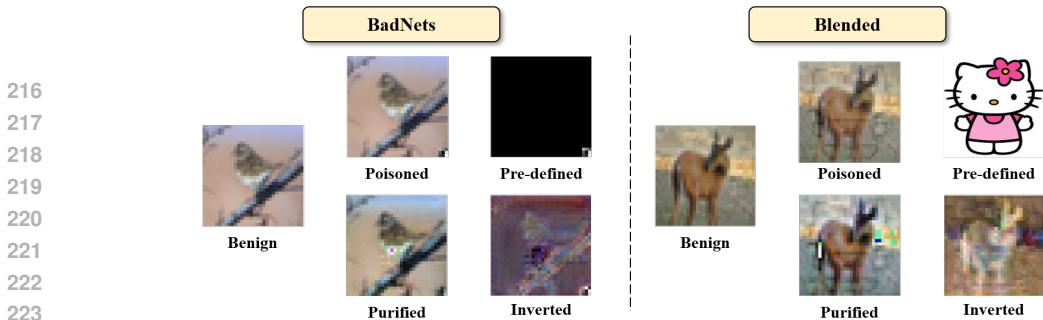


Figure 3: The visualization of BTI-DBF in inverting backdoor triggers under both BadNets and Blended attacks. We display the benign samples, poisoned samples, purified samples, pre-defined triggers, and inverted triggers, respectively. The inverted triggers and purified samples are different from the ground-truth ones to a large extent, leading to the ineffectiveness of BTI-based defenses.

(b-2)~(b-4), with increasing transformation intensity, the feature distribution of backdoored samples within the same class undergoes greater changes, indicating that higher transformation levels effectively mismatch or remove trigger patterns. Nevertheless, the difficulty of decoupling benign patterns from backdoor patterns in the input domain results in that such transformations inevitably affect the benign features. It causes a shift in the centroid of the benign sample feature distribution (visualized as solid circles in Figure 2). The primary reason for this trade-off problem lies in the fact that the output domain of the DNN remains consistent before and after the defenses. This consistency forces the input transformation module to perform two conflicting tasks: (1) effectively removing trigger patterns and (2) preserving the benign patterns of samples while ensuring they are classified into the correct output categories of the original model. This conflict inspires us to consider that adjusting the model’s output domain may help mitigate this issue.

3.3 THE LIMITATIONS OF BTI-BASED DEFENSES

BTI-based defenses can break the trade-off between model utility and defense effectiveness by introducing the information of triggers via trigger inversion. In the pre-processing-based defense paradigm, BTI-based defenses typically involve two steps: trigger inversion and trigger removal. Specifically, the defender first utilizes several data to invert the pre-injected trigger. The inverted trigger is then used to eliminate any potential trigger patterns in samples before prediction. The effectiveness of BTI-based defenses highly relies on the quality of the inverted trigger. However, we argue that *the inversion of the backdoor trigger pattern in high quality is inherently challenging due to the absence of prior knowledge*. The difficulty limits the effect of BTI-based backdoor defenses.

We implement the state-of-the-art BTI-based defense, BTI-DBF (Xu et al., 2024), to invert the backdoor triggers of BadNets (Gu et al., 2019) and Blended (Chen et al., 2017) on CIFAR-10. As shown in Figure 3, the trigger patterns obtained by BTI-DBF differ significantly from the ground-truth trigger patterns. This discrepancy illustrates why existing BTI-based defenses fail to eliminate backdoor patterns present in input samples effectively. Moreover, BTI-based defenses often identify “pseudo-triggers” inherent in DNNs, which usually arise from the model’s vulnerability to adversarial perturbations. When defenders attempt to eliminate these non-authentic trigger patterns before processing the samples into the model, it can disrupt the benign features of the samples, while the backdoor patterns remain largely unaffected. If the quality and authenticity of the inverted trigger patterns cannot be guaranteed, the BTI-based defenses may potentially yield adverse outcomes. Arguably, achieving BTI is a challenging endeavor due to the lack of prior knowledge about the implanted backdoor and poisoned samples. As such, it is necessary to design an inversion-free backdoor defense to break the aforementioned trade-off.

4 METHODOLOGY

4.1 MOTIVATION AND INSPIRATION

In Section 3, we empirically evaluate existing pre-processing-based defenses and analyze why they are ineffective. In this section, we present a theoretical analysis and the inspiration to design an effective and efficient backdoor defense method. Given a pre-processing method $\mathcal{T}(\cdot)$ and a pre-trained model $\mathcal{F}(\cdot)$, we have the following theorem.

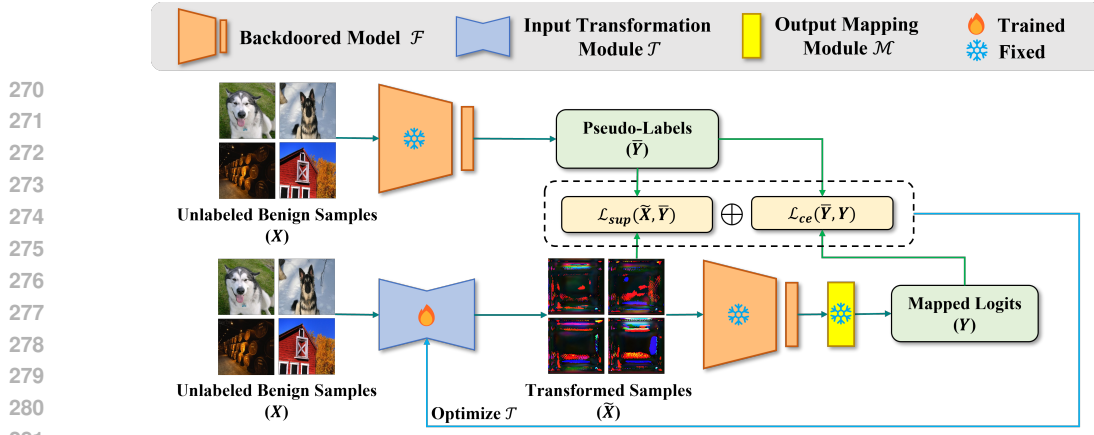


Figure 4: The main optimization pipeline of our REFINE. There are two main components: input transformation module \mathcal{T} and output mapping module \mathcal{M} . Specifically, after obtaining the fixed pre-trained model, the defender first specifies a particular hard-coded mapping \mathcal{M} and then optimizes \mathcal{T} guided by the loss function \mathcal{L} , using the unlabeled benign dataset. The loss function \mathcal{L} consists of the cross-entropy loss \mathcal{L}_{ce} which aims to maintain the model’s utility, and the supervised contrastive loss \mathcal{L}_{sup} to enhance the defense capability via forcing orderly sample aggregation.

Theorem 1. Given a K -class pre-trained deep learning model $\mathcal{F}(\cdot) = s(f(\cdot))$ where $s(\cdot)$ is the softmax function and $f(\cdot)$ is the feature extractor, and a pre-processing method $\mathcal{T}(\cdot)$, \mathbf{x} is the data from a specific domain \mathcal{D} (i.e., $\mathbf{x} \sim \mathcal{D}$) and $\tilde{\mathbf{x}} = \mathcal{T}(\mathbf{x}) \sim \tilde{\mathcal{D}}$. Let $\Phi_{\mathcal{D}}(\mathbf{x})$ and $\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}})$ denotes the probability density function of \mathcal{D} and $\tilde{\mathcal{D}}$, we have

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \leq 2\alpha\sqrt{K} \cdot \mathcal{W}_1(\mu, \tilde{\mu}), \quad (1)$$

where $\mathcal{W}_1(\mu, \tilde{\mu})$ is the Wasserstein-1 distance between μ and $\tilde{\mu}$, μ and $\tilde{\mu}$ are the probability measures of the representations $f(\mathbf{x})$ and $f(\tilde{\mathbf{x}})$, and $\alpha = \max[\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}|\mathbf{x})/\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}})]$.

Theorem 1 indicates why existing defenses are ineffective. Assuming \mathbf{x} is the poisoned sample, the left part of Eq. (1) means the distance between the prediction of the transformed poisoned sample and the original poisoned prediction. Theorem 1 demonstrates that the distance is bounded by the Wasserstein-1 distance between the probability measures $\mu, \tilde{\mu}$ of the output representations. Thus, to maintain the utility of the model, existing pre-processing methods tend to retain the output representations, leading to limited effectiveness in defending the backdoor. Otherwise, they have to compromise the utility to achieve greater backdoor defense behaviors. The proof is in Appendix A.

Following the above theorem, we can enhance the upper bound by increasing the distance between $\mu, \tilde{\mu}$. Inspired by model reprogramming techniques (Chen, 2024), we propose REFINE, a reprogramming-based inversion-free backdoor defense method. Model reprogramming can significantly transform the output domain to destroy trigger patterns while maintaining model utility for it also changes the input domain. Specifically, we introduce an input transformation module to transform the inputs, and a label mapping module to remap the original classes to new shuffled ones. We also employ a supervised contrastive loss to further enlarge the distances among different classes. The technical details of our REFINE method are illustrated in the following parts.

4.2 REFINE: REPROGRAMMING-BASED INVERSION-FREE BACKDOOR DEFENSE METHOD

In general, REFINE consists of two essential components: (1) the input transformation module \mathcal{T} , which disrupts the benign and backdoor patterns of input samples through transformations and generates new benign features; (2) the label mapping module \mathcal{M} , which formulates the specified source-target hard-coded label remapping function and maps the original classes to new shuffled classes. Additionally, we integrate the cross-entropy loss \mathcal{L}_{ce} and the supervised contrastive loss \mathcal{L}_{sup} to steer the optimization of \mathcal{T} . The illustration of our REFINE is shown in Figure 4.

4.2.1 INPUT TRANSFORMATION MODULE

To effectively alter potential trigger patterns in the input samples, we need to modify the input domain of the original model. Traditional model reprogramming methods (Elsayed et al., 2019; Tsai et al., 2020) add the optimized universal adversarial perturbation around the input, while the

trigger pattern still remains intact on the backdoor image to some extent. In contrast, we utilize a trainable autoencoder (e.g., UNet) as the foundational structure for our input transformation module, leveraging its capability to preserve the dimension of samples before and after transformations, while making more significant modifications to the whole image. Upon inputting a batch of data, the input transformation module will encode the pixel features from the images and then decode them to produce new samples. The transformed samples $\tilde{\mathbf{X}}$ can be described as follows:

$$\tilde{\mathbf{X}} = \mathcal{T}(\mathbf{X}, \theta), \quad (2)$$

where \mathbf{X} is a batch of input samples, and $\mathcal{T}(\cdot, \theta)$ is the input transformation module with θ as its trainable parameters. Arguably, this module not only preserves the consistency of sample size before and after the transformation but also affords a higher degree of flexibility in sample manipulation compared to conventional reprogramming methods. During this transformation process, both benign and backdoor patterns are disarranged, effectively removing potential triggers and causing the generation of new benign features orderly clustered by their respective classes.

4.2.2 OUTPUT MAPPING MODULE

Once the input samples are transformed into new samples via the input transformation module, they are subsequently processed by the original backdoored model, which generates confidence scores for each class, as expressed below:

$$\tilde{\mathbf{Y}} = \mathcal{F}(\tilde{\mathbf{X}}), \quad (3)$$

where $\mathcal{F}(\cdot)$ is the original backdoored model. As demonstrated in Section 3.2, fixing the model’s output domain leads to a trade-off between transformation intensity and defense performance. To address this issue, we introduce an output mapping module at the model’s output end, aiming to alter the output domain and mitigate the aforementioned challenges. Specifically, the output mapping module redefines the class order of the model’s output layer, which hard-codes a one-to-one label remapping function $f_L : \tilde{l} \mapsto l$, where $\tilde{l}, l \in L, \tilde{l} \neq l$, L is the set of labels. The confidence scores generated by the original model can be remapped into new scores through \mathcal{M} , as follows:

$$\mathbf{Y} = \mathcal{M}(\tilde{\mathbf{Y}}). \quad (4)$$

The final predictions for the samples can be derived from the confidence scores \mathbf{Y} outputted by \mathcal{M} .

4.2.3 OPTIMIZING REFINE MODULES

To maximize the flexibility of input transformations for removing trigger patterns while maintaining the original model’s accuracy, we incorporate two crucial loss functions, the cross-entropy loss and the supervised contrastive loss, to guide the optimization of the input transformation module. The formulation of the combined loss function can be expressed as follows:

$$\min_{\theta} \mathcal{L}_{refine} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{sup}. \quad (5)$$

In Eq. (5), \mathcal{L}_{ce} and \mathcal{L}_{sup} indicate the cross-entropy loss and the supervised contrastive loss, respectively. λ is a scalar temperature parameter, and θ represents the set of parameters in the input transformation module to be optimized during training. [Since Theorem 1 does not guarantee the model performance on clean samples, adding \$\mathcal{L}_{ce}\$ to maintain the utility of the model is necessary.](#)

In our threat model, the dataset available to the defender is unlabeled. [Therefore, before calculating these loss functions, it is necessary to obtain the pseudo-labels \$\bar{\mathbf{Y}}\$ for the current batch of unlabeled samples \$\mathbf{X}\$. \$\bar{\mathbf{Y}}\$ is predicted by the original model \(without any additional modules\), as follows:](#)

$$\bar{\mathbf{Y}} = \arg \max(\mathcal{F}(\mathbf{X})). \quad (6)$$

Leveraging Cross-entropy Loss to Maintain the Utility. Due to the substantial modification of the original model’s output domain facilitated by the output mapping module, the input transformation module is no longer constrained by the requirement to preserve the original benign features of the samples. Nevertheless, the model must retain its original performance within the new output domain, which necessitates the employment of cross-entropy loss to effectively guide the sample transformation process. The cross-entropy loss is typically formalized as follows:

$$\mathcal{L}_{ce} = -\frac{1}{N} \sum_{i=1}^N \bar{y}_i \log(\mathbf{y}_i), \quad (7)$$

where N represents the number of samples in the current data batch \mathbf{X} . $\bar{\mathbf{y}}_i \in \bar{\mathbf{Y}}$ denotes the pseudo-label for sample $\mathbf{x}_i \in \mathbf{X}$ (typically a one-hot encoded vector), and $\mathbf{y}_i \in \mathbf{Y}$ indicates the predicted probability remapped by the output mapping module for sample \mathbf{x}_i .

Utilizing Supervised Contrastive Loss to Enhance Backdoor Defense. Arguably, relying solely on cross-entropy loss is insufficient to restore the original model’s benign accuracy and mitigate the backdoor. Therefore, we introduce supervised contrastive loss (Khosla et al., 2020), where “supervised” refers to the original model as the supervisor. Specifically, the supervised contrastive loss aims to ensure that features of transformed samples from the same class are more similar, while those from different classes are further apart. It can be defined as follows.

$$\mathcal{L}_{sup} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\tilde{\mathbf{x}}_i \cdot \tilde{\mathbf{x}}_p / \tau)}{\sum_{a \in A(i)} \exp(\tilde{\mathbf{x}}_i \cdot \tilde{\mathbf{x}}_a / \tau)}, \quad (8)$$

where $I \equiv \{1, 2, \dots, N\}$ represents indices of all samples in current data batch, $\tilde{\mathbf{x}}_i = \mathcal{T}(\mathbf{x}_i, \theta) \in \tilde{\mathbf{X}}$, the \cdot symbol denotes the inner (dot) product, τ is a scalar temperature parameter, and $A(i) \equiv I \setminus \{i\}$. The set $P(i) \equiv \{p \in A(i) : \bar{\mathbf{y}}_p = \bar{\mathbf{y}}_i\}$ contain indices of all positives in the batch distinct from i , and $|P(i)|$ is its cardinality. The pseudo-code can be found in Appendix B.

4.2.4 UTILIZING REFINE FOR MODEL INFERENCE

During the model inference phase, we can apply the aforementioned well-trained modules to achieve high-performance and secure predictions. The input samples are sequentially processed through the input transformation module $\mathcal{T}(\cdot, \theta)$, the original pre-trained model $\mathcal{F}(\cdot)$, and the output mapping module $\mathcal{M}(\cdot)$. This process ultimately yields the predicted confidence scores, with all parameters remaining constant. The inference process can be formally expressed as follows.

$$\mathbf{y} = \mathcal{M}(\mathcal{F}(\mathcal{T}(\mathbf{x}, \theta))), \quad (9)$$

where \mathbf{x} represents the sample to be predicted. The detailed process is illustrated in Figure 1.

5 EXPERIMENTS

In this section, we evaluate the effectiveness of our REFINE compared with different existing backdoor defenses. We also conduct an ablation study and evaluate the resistance to potential adaptive attacks. The analysis of the overhead of REFINE is in Appendix G and the implementation of REFINE in the black-box scenario is in Appendix F.

5.1 EXPERIMENTAL SETTINGS

Datasets and Models. We conduct experiments on two classical benchmark datasets, including CIFAR-10 (Krizhevsky et al., 2009) and (a subset of) ImageNet (Deng et al., 2009) containing 50 classes. We evaluated our method with ResNet-18 (He et al., 2016) on both datasets. We also validate the effectiveness of REFINE on other models in Appendix E. Note that our goal is to evaluate the effectiveness of backdoor defense methods instead of training a SOTA model. Therefore, the benign accuracies of our models may be lower than the SOTA models. We exploit U-Net (Ronneberger et al., 2015) as the structure of the input transformation module.

Attack Setup. We utilize 7 representative advanced backdoor attacks, including (1) BadNets (Gu et al., 2019), (2) Blended (Chen et al., 2017), (3) WaNet (Nguyen & Tran, 2021), (4) PhysicalBA (dubbed ‘Physical’) (Li et al., 2021c), (5) BATT (Xu et al., 2023), (6) LabelConsistent (dubbed ‘LC’) (Turner et al., 2019), and (7) Adaptive-Patch (dubbed ‘Adaptive’) (Qi et al., 2023), to comprehensively evaluate the performance of different defenses. The poisoning rates are all set to 10%.

Defense Setup. We compare the defense performance of REFINE with both types of pre-processing-based defense methods. For transformation-based defenses, we utilize three representative and advanced methods, including (1) AutoEncoderDefense (dubbed ‘AEDefense’) (Liu et al., 2017), (2) ShrinkPad (Li et al., 2021c), (3) BDMAE (Sun et al., 2023). For BTI-based defenses, we employ three methods as baseline, including (1) Neural Cleanse (dubbed ‘NC’) (Wang et al., 2019), (2) FeatureRE (Wang et al., 2022b), (3) BTI-DBF(P) (Xu et al., 2024).

Table 1: The performance (%) of REFINE and the transformation-based backdoor defenses. The best results are **boldfaced**, while all failed cases (BA drop > 5% or ASR > 10%) are marked in **red**.

Dataset↓	Defense→ Attack↓	No Defense		AEDefense		ShrinkPad		BDMAE		REFINE	
		BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
CIFAR-10	BadNets	92.05	100	85.05	0.93	84.10	14.24	91.28	11.22	91.22	0.99
	Blended	90.63	99.15	83.76	85.72	84.44	98.83	90.04	99.15	90.65	0.76
	WaNet	91.63	99.92	85.18	1.02	85.05	36.56	89.90	99.92	90.82	0.96
	Physical	93.71	100	89.20	1.70	92.52	82.72	93.31	48.26	90.95	1.08
	BATT	92.93	99.89	88.58	99.99	88.46	60.36	92.80	98.90	90.81	2.12
	LC	92.37	99.91	84.42	0.93	83.95	9.01	91.52	12.44	90.78	1.20
	Adaptive	90.17	100	83.61	5.42	81.95	23.66	80.61	24.88	90.31	0.17
ImageNet	BadNets	66.94	99.44	60.84	3.31	63.00	3.27	56.80	2.04	68.95	0.69
	Blended	66.00	95.64	59.56	90.41	59.64	96.82	50.44	97.32	68.75	1.31
	WaNet	66.32	95.08	47.72	0.41	59.48	48.82	51.40	94.12	68.18	3.02
	Physical	72.72	99.76	62.32	4.94	71.36	84.57	60.64	9.88	69.20	1.55
	BATT	72.04	98.72	66.88	99.02	70.08	92.69	62.04	98.64	69.40	1.22
	LC	66.64	78.72	62.12	0.53	62.88	2.61	56.88	8.80	69.85	0.45
	Adaptive	67.43	81.91	61.03	9.98	64.27	39.61	53.72	42.08	68.00	0.00

Table 2: The performance (%) of REFINE and the BTI-based backdoor defenses. The best results are **boldfaced**, while all failed cases (BA drop > 5% or ASR > 10%) are marked in **red**.

Dataset↓	Defense→ Attack↓	No Defense		NC		FeatureRE		BTI-DBF(P)		REFINE	
		BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
CIFAR-10	BadNets	92.05	98.53	87.09	1.39	64.72	13.07	89.25	5.56	91.22	1.00
	Blended	90.63	98.05	89.53	99.92	90.38	1.67	88.08	3.91	90.38	0.96
	WaNet	91.63	98.81	88.15	6.89	93.35	0.06	91.03	5.52	90.82	0.96
	Physical	93.71	98.89	90.38	2.95	93.75	1.06	89.74	4.60	91.10	0.77
	BATT	92.93	98.78	91.45	32.18	92.84	2.15	89.72	9.37	90.77	2.02
	LC	92.37	98.80	84.31	5.15	83.76	0.06	88.60	3.50	90.88	1.43
	Adaptive	90.03	100	85.32	87.83	87.94	72.47	86.72	13.79	90.11	1.58
ImageNet	BadNets	65.72	99.90	65.49	1.48	64.84	5.32	65.78	7.06	67.03	0.67
	Blended	65.43	99.56	61.03	96.13	59.00	23.47	60.34	2.14	68.79	1.15
	WaNet	64.97	96.97	64.85	10.52	61.84	5.48	62.17	0.70	68.13	2.96
	Physical	69.58	99.92	60.37	3.22	73.08	11.76	64.93	31.46	69.28	0.83
	BATT	71.76	98.94	63.04	21.78	70.24	5.19	62.13	0.18	69.37	1.23
	LC	66.03	81.50	58.13	6.57	65.72	0.08	59.17	8.31	69.94	0.32
	Adaptive	67.40	78.32	57.15	79.45	55.00	64.38	54.36	12.58	67.68	0.44

Evaluation Metrics. Consistent with the standard evaluation metrics in backdoor-related studies (Li et al., 2022b), we utilize benign accuracy (BA) and attack success rate (ASR) to assess all defense methods. BA and ASR are the accuracies of the benign samples and the poisoned samples, respectively. An effective defense is indicated by a higher BA and a lower ASR.

5.2 MAIN RESULTS

As shown in Tables 1-2, our REFINE successfully mitigates backdoor threats in all cases while preserving high benign accuracy. Specifically, the ASRs of our method are lower than 4% (< 2% in most cases). For the BA, the models under REFINE experience less than 3% drop on the CIFAR-10 dataset compared to the undefended models. On the ImageNet dataset, the BA even improves, due to the increased depth of the original models introduced by the input transformation module. In contrast, other baseline defenses may fail in certain cases, with BA drop > 5% or ASR > 10%. Specifically, for the Adaptive-Patch attack (Qi et al., 2023), all baseline defenses perform poorly, because of the difficulty of inverting the adaptive trigger patterns for BTI-based defenses or disrupting the trigger patterns for transformation-based defenses.

5.3 ABLATION STUDY

There are three important components in our methods, including (1) the input transformation method, (2) hard-coded remapping function (HRF for short) in the output mapping module, and (3) supervised contrastive loss (SCL for short) of transformed samples. In this section, we present an ablation study on the former two modules and verify their effectiveness. We also test different architectures of the input transformation module and conduct additional ablation studies in Appendix E.

As shown in Table 3, we evaluate the defense performance of REFINE without the hard-coded remapping function (w/o HRF) or without the supervised contrastive loss (w/o SCL). Experimental results indicate that without the hard-coded remapping function, REFINE successfully preserves

Table 3: The performance (%) of REFINE with/without the hard-coded remapping function (HRF) or with/without the supervised contrastive loss (SCL).

Defense→	No Defense		w/o HRF		w/o SCL		REFINE	
	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNets	91.70	100	91.23	70.76	89.26	1.43	90.92	0.68
Blended	91.10	98.76	90.59	75.30	90.38	0.10	90.65	0.51
WaNet	91.09	99.98	91.03	99.53	89.08	1.45	90.45	0.88
Physical	93.59	100	92.86	1.60	88.63	1.97	90.92	1.36
BATT	92.43	99.91	91.67	72.46	88.82	5.87	90.89	1.97
LC	92.30	99.74	91.88	69.15	90.37	0.59	90.57	1.25
Adaptive	90.54	100	89.77	62.94	88.06	0.32	90.17	0.27

Table 4: The performance (%) of REFINE against potential adaptive attacks.

Setting→	Normal Attack				Adaptive Attack			
	No Defense		REFINE		No Defense		REFINE	
Defense→	BA	ASR	BA	ASR	BA	ASR	BA	ASR
Dataset↓								
CIFAR-10	91.74	100	90.71	1.07	84.53	100	83.05	0.98
ImageNet	66.94	99.59	69.00	0.70	58.39	100	60.53	1.09

the BA of the original model, but struggles to reduce the ASR of the backdoor. This is because, without the hard-coded remapping function, the output domain of the model remains unchanged. Subsequently, it encounters the same trade-off problem as other transformation-based defenses, and is difficult to find a balance between transformation intensity and defense performance. Also, in the absence of supervised contrastive loss, REFINE can effectively reduce ASR with the help of the hard-coded remapping function. However, it encounters difficulties in restoring the BA of the original model, which may adversely affect the model’s inference capabilities.

5.4 RESISTANCE TO POTENTIAL ADAPTIVE ATTACKS

In this section, we examine whether the adversary can circumvent our defenses if they have full knowledge of the process of our REFINE. After training the original backdoored model, the adversary can fine-tune it utilizing an input transformation module, along with a randomly initialized hard-coded output mapping module, to simulate our REFINE. During fine-tuning, the loss function for model optimization can be expressed as follows:

$$\min_{\delta} \mathcal{L}_{adapt} = \mathcal{L}_b + \gamma \mathcal{L}_{refine}, \quad (10)$$

where \mathcal{L}_b indicates the backdoor loss function in the original training phase of the backdoored model, and \mathcal{L}_{refine} represents the loss function of REFINE. γ is a scalar temperature parameter, and δ denotes the parameters to be trained in the backdoored model. Ideally, the adversary can achieve the backdoor target with a low value of \mathcal{L}_{refine} by optimizing Eq. (10). Consequently, the REFINE may not work well since the \mathcal{L}_{refine} is already low.

As shown in Table 4, REFINE is still highly effective with high BAs (BA drop < 1.5%) and low ASRs (< 1.5%). It is mostly because defenders can arbitrarily specify the output mapping function and train an input transformation module that may entirely differ from the attacker’s. Besides, the original backdoored model experiences a decrease in BA after undergoing adaptive attack training. As such, these results demonstrate that our defense method is resistant to adaptive attacks.

6 CONCLUSION

In this paper, we revisited existing pre-processing-based backdoor defense methods, including backdoor-trigger-inversion-based (BTI-based) defenses and transformation-based defenses. We revealed the limitations of the two defense methods. Subsequently, according to the empirical and theoretical analysis, we proposed REFINE, a reprogramming-based inversion-free backdoor defense method. This method was motivated by the insight that increasing the distances of the feature representations before and after the transformation may lead to a better effectiveness of backdoor defense. Specifically, we introduced an input transformation module and an output mapping module. We also utilized the supervised contrastive loss to enhance the defense performance. Results on benchmark datasets verified the effectiveness of our REFINE and the resistance to the adaptive attack. We hope our REFINE can provide a new angle to facilitate the design of more effective backdoor defenses.

540 ETHICS STATEMENT

541

542 This paper proposes an inversion-free backdoor defense method, REFINE. Our method can be uti-
 543 lized to mitigate the effect of the backdoor. Therefore, our REFINE is a defensive method and our
 544 work does not discover any new threat. Our research also does not include any human subjects.
 545 Accordingly, this paper does not raise ethical issues.

546

547 REPRODUCIBILITY STATEMENT

548

549 The details of our implementations and experiments can be found in Appendix C. We also provide
 550 the codes for reproducing our main results in the supplementary material. We will make the full
 551 codes of our method open-source once the paper is accepted.

552

553 REFERENCES

554

555 Pin-Yu Chen. Model reprogramming: Resource-efficient cross-domain machine learning. In *AAAI*,
 556 2024.

557

558 Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep
 559 learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.

560 Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack frame-
 561 work for diffusion models. *NeurIPS*, 36, 2024.

562

563 Luke N Darlow, Elliot J Crowley, Antreas Antoniou, and Amos J Storkey. Cinic-10 is not imagenet
 564 or cifar-10. *arXiv preprint arXiv:1810.03505*, 2018.

565 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 566 hierarchical image database. In *CVPR*, pp. 248–255, 2009.

567

568 Sharmita Dey and Sarath R Nair. Enhancing joint motion prediction for individuals with limb loss
 569 through model reprogramming. *arXiv preprint arXiv:2403.06569*, 2024.

570 Bao Gia Doan, Ehsan Abbasnejad, and Damith C Ranasinghe. Februus: Input purification defense
 571 against trojan attacks on deep neural network systems. In *ACSAC*, pp. 897–912, 2020.

572

573 Khoa Doan, Yingjie Lao, Weijie Zhao, and Ping Li. Lira: Learnable, imperceptible and robust
 574 backdoor attacks. In *ICCV*, pp. 11966–11976, 2021.

575 Jianshuo Dong, Han Qiu, Yiming Li, Tianwei Zhang, Yuanjie Li, Zeqi Lai, Chao Zhang, and Shu-
 576 Tao Xia. One-bit flip is all you need: When bit-flip attack meets model training. In *ICCV*, pp.
 577 4688–4698, 2023.

578 Min Du, Ruoxi Jia, and Dawn Song. Robust anomaly detection and backdoor attack detection via
 579 differential privacy. In *ICLR*, 2020.

580

581 Gamaleldin F Elsayed, Ian Goodfellow, and Jascha Sohl-Dickstein. Adversarial reprogramming of
 582 neural networks. In *ICLR*, 2019.

583 Bolin Gao and Lacra Pavel. On the properties of the softmax function with application in game
 584 theory and reinforcement learning. *arXiv preprint arXiv:1704.00805*, 2017.

585

586 Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal.
 587 Strip: A defence against trojan attacks on deep neural networks. In *ACSAC*, pp. 113–125, 2019.

588 Yansong Gao, Bao Gia Doan, Zhi Zhang, Siqi Ma, Jiliang Zhang, Anmin Fu, Surya Nepal, and
 589 Hyounghick Kim. Backdoor attacks and countermeasures on deep learning: A comprehensive
 590 review. *arXiv preprint arXiv:2007.10760*, 2020.

591

592 Yinghua Gao, Yiming Li, Linghui Zhu, Dongxian Wu, Yong Jiang, and Shu-Tao Xia. Not all samples
 593 are born equal: Towards effective clean-label backdoor attacks. *Pattern Recognition*, 139:109512,
 2023.

- 594 Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdooring
595 attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- 596
- 597 Jonathan Hayase and Weihao Kong. Spectre: Defending against backdoor attacks using robust
598 covariance estimation. In *ICML*, 2020.
- 599
- 600 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
601 nition. In *CVPR*, pp. 770–778, 2016.
- 602
- 603 Yiling He, Jian Lou, Zhan Qin, and Kui Ren. Finer: Enhancing state-of-the-art classifiers with
604 feature attribution to facilitate security analysis. In *CCS*, pp. 416–430, 2023.
- 605
- 606 Yu He, Boheng Li, Yao Wang, Mengda Yang, Juan Wang, Hongxin Hu, and Xingyu Zhao. Is
607 difficulty calibration all we need? towards more practical membership inference attacks. *arXiv
preprint arXiv:2409.00426*, 2024.
- 608
- 609 Linshan Hou, Ruili Feng, Zhongyun Hua, Wei Luo, Leo Yu Zhang, and Yiming Li. Ibd-psc: Input-
610 level backdoor detection via parameter-oriented scaling consistency. In *ICML*, 2024.
- 611
- 612 Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
613 convolutional networks. In *CVPR*, pp. 4700–4708, 2017.
- 614
- 615 Kunzhe Huang, Yiming Li, Baoyuan Wu, Zhan Qin, and Kui Ren. Backdoor defense via decoupling
616 the training process. In *ICLR*, 2022.
- 617
- 618 Mojan Javaheripi, Mohammad Samragh, Gregory Fields, Tara Javidi, and Farinaz Koushanfar.
619 Cleann: Accelerated trojan shield for embedded neural networks. In *ICCD*, pp. 1–9, 2020.
- 620
- 621 Yongcheng Jing, Chongbin Yuan, Li Ju, Yiding Yang, Xinchao Wang, and Dacheng Tao. Deep
622 graph reprogramming. In *CVPR*, pp. 24345–24354, 2023.
- 623
- 624 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron
625 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *NeurIPS*, 33:18661–
626 18673, 2020.
- 627
- 628 Eliska Kloberdanz, Jin Tian, and Wei Le. An improved (adversarial) reprogramming technique for
629 neural networks. In *ICANN*, pp. 3–15, 2021.
- 630
- 631 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
632 *Technical report*, 2009.
- 633
- 634 Boheng Li, Yishuo Cai, Jisong Cai, Yiming Li, Han Qiu, Run Wang, and Tianwei Zhang. Puri-
635 fying quantization-conditioned backdoors via layer-wise activation correction with distribution
636 approximation. In *ICML*, 2024a.
- 637
- 638 Boheng Li, Yishuo Cai, Haowei Li, Feng Xue, Zhifeng Li, and Yiming Li. Nearest is not dear-
639 est: Towards practical defense against quantization-conditioned backdoor attacks. In *CVPR*, pp.
640 24523–24533, 2024b.
- 641
- 642 Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learn-
643 ing: Training clean models on poisoned data. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S.
644 Liang, and J. Wortman Vaughan (eds.), *NeurIPS*, volume 34, pp. 14900–14912. Curran Asso-
645 ciates, Inc., 2021a.
- 646
- 647 Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distil-
648 lation: Erasing backdoor triggers from deep neural networks. *arXiv preprint arXiv:2101.05930*,
649 2021b.
- 650
- 651 Yiming Li, Tongqing Zhai, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor attack in the
652 physical world. *arXiv preprint arXiv:2104.02361*, 2021c.
- 653
- 654 Yiming Li, Yang Bai, Yong Jiang, Yong Yang, Shu-Tao Xia, and Bo Li. Untargeted backdoor
655 watermark: Towards harmless and stealthy dataset copyright protection. In *NeurIPS*, pp. 13238–
656 13250, 2022a.

- 648 Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A survey. *TNNLS*, 35(1):
649 5–22, 2022b.
- 650
- 651 Yiming Li, Haoxiang Zhong, Xingjun Ma, Yong Jiang, and Shu-Tao Xia. Few-shot backdoor attacks
652 on visual object tracking. In *ICLR*, 2022c.
- 653
- 654 Yinshan Li, Hua Ma, Zhi Zhang, Yansong Gao, Alsharif Abuadbba, Minhui Xue, Anmin Fu, Yifeng
655 Zheng, Said F Al-Sarawi, and Derek Abbott. Ntd: Non-transferability enabled deep learning
656 backdoor detection. *IEEE Transactions on Information Forensics and Security*, 2023a.
- 657
- 658 Yizhe Li, Yu-Lin Tsai, Chia-Mu Yu, Pin-Yu Chen, and Xuebin Ren. Exploring the benefits of visual
659 prompting in differential privacy. In *ICCV*, pp. 5158–5167, 2023b.
- 660
- 661 Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor
662 attack with sample-specific triggers. In *ICCV*, pp. 16463–16472, 2021d.
- 663
- 664 Siyuan Liang, Mingli Zhu, Aishan Liu, Baoyuan Wu, Xiaochun Cao, and Ee-Chien Chang. Bad-
665 clip: Dual-embedding guided backdoor attack on multimodal contrastive learning. In *CVPR*, pp.
666 24645–24654, 2024.
- 667
- 668 Xiyao Liu, Shuo Shao, Yue Yang, Kangming Wu, Wenyan Yang, and Hui Fang. Secure federated
669 learning model verification: A client-side backdoor triggered watermarking scheme. In *SMC*, pp.
670 2414–2419. IEEE, 2021.
- 671
- 672 Yuntao Liu, Yang Xie, and Ankur Srivastava. Neural trojans. In *ICCD*, pp. 45–48, 2017.
- 673
- 674 Zhihao Liu, Jian Lou, Wenjie Bao, Zhan Qin, and Kui Ren. Differentially private zeroth-order
675 methods for scalable large language model finetuning. *arXiv preprint arXiv:2402.07818*, 2024.
- 676
- 677 Brandon B May, N Joseph Tatro, Dylan Walker, Piyush Kumar, and Nathan Shnidman. Salient
678 conditional diffusion for defending against backdoor attacks. *arXiv preprint arXiv:2301.13862*,
679 2023.
- 680
- 681 Paarth Neekhara, Shehzeen Hussain, Jinglong Du, Shlomo Dubnov, Farinaz Koushanfar, and Julian
682 McAuley. Cross-modal adversarial reprogramming. In *WACV*, pp. 2427–2435, 2022.
- 683
- 684 Anh Nguyen and Anh Tran. Wanet–imperceptible warping-based backdoor attack. *arXiv preprint*
685 *arXiv:2102.10369*, 2021.
- 686
- 687 Xiangyu Qi, Tinghao Xie, Yiming Li, Saeed Mahloujifar, and Prateek Mittal. Revisiting the as-
688 sumption of latent separability for backdoor defenses. In *ICLR*, 2023.
- 689
- 690 Han Qiu, Yi Zeng, Shangwei Guo, Tianwei Zhang, Meikang Qiu, and Bhavani Thuraisingham.
691 Deepsweep: An evaluation framework for mitigating dnn backdoor attacks using data augmenta-
692 tion. In *AsiaCCS*, pp. 363–377, 2021.
- 693
- 694 Dazhong Rong, Guoyao Yu, Shuheng Shen, Xinyi Fu, Peng Qian, Jianhai Chen, Qinming He, Xing
695 Fu, and Weiqiang Wang. Clean-image backdoor attacks. In *ICANN*, pp. 187–202. Springer, 2024.
- 696
- 697 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-
698 ical image segmentation. In *MICCAI*, pp. 234–241, 2015.
- 699
- 700 Shuo Shao, Wenyan Yang, Hanlin Gu, Zhan Qin, Lixin Fan, Qiang Yang, and Kui Ren. Fedtracker:
701 Furnishing ownership verification and traceability for federated learning model. *TDSC*, 2024.
- 702
- 703 Shuo Shao, Yiming Li, Hongwei Yao, Yiling He, Zhan Qin, and Kui Ren. Explanation as a wa-
704 termark: Towards harmless and multi-bit model ownership verification via watermarking feature
705 attribution. In *NDSS*, 2025.
- 706
- 707 Yucheng Shi, Mengnan Du, Xuansheng Wu, Zihan Guan, Jin Sun, and Ninghao Liu. Black-box
708 backdoor defense via zero-shot image purification. *NeurIPS*, 36:57336–57366, 2023.
- 709
- 710 Karen Simonyan. Very deep convolutional networks for large-scale image recognition. *arXiv*
711 *preprint arXiv:1409.1556*, 2014.

- 702 Tao Sun, Lu Pang, Chao Chen, and Haibin Ling. Mask and restore: Blind backdoor defense at test
703 time with masked autoencoder. *arXiv preprint arXiv:2303.15564*, 2023.
704
- 705 Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking
706 the inception architecture for computer vision. In *CVPR*, pp. 2818–2826, 2016.
707
- 708 Yun-Yun Tsai, Pin-Yu Chen, and Tsung-Yi Ho. Transfer learning without knowing: Reprogramming
709 black-box machine learning models with scarce data and limited resources. In *ICML*, pp. 9614–
710 9624, 2020.
- 711 Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks.
712 *arXiv preprint arXiv:1912.02771*, 2019.
713
- 714 Miguel Villarreal-Vasquez and Bharat Bhargava. Confoc: Content-focus protection against trojan
715 attacks on neural networks. *arXiv preprint arXiv:2007.00711*, 2020.
- 716 Ria Vinod, Pin-Yu Chen, and Payel Das. Reprogramming pretrained language models for protein
717 sequence representation learning. *arXiv preprint arXiv:2301.02120*, 2023.
718
- 719 Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y
720 Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *S&P*,
721 pp. 707–723. IEEE, 2019.
- 722 Lixu Wang, Shichao Xu, Ruiqi Xu, Xiao Wang, and Qi Zhu. Non-transferable learning: A new
723 approach for model ownership verification and applicability authorization. In *ICLR*, 2022a.
724
- 725 Zhenting Wang, Kai Mei, Hailun Ding, Juan Zhai, and Shiqing Ma. Rethinking the reverse-
726 engineering of trojan triggers. *NeurIPS*, 35:9738–9753, 2022b.
- 727 Zhenting Wang, Kai Mei, Juan Zhai, and Shiqing Ma. Unicorn: A unified backdoor trigger inversion
728 framework. In *ICLR*, 2023.
729
- 730 Cheng Wei, Yang Wang, Kuofeng Gao, Shuo Shao, Yiming Li, Zhibo Wang, and Zhan Qin. Point-
731 ncbw: Towards dataset ownership verification for point clouds via negative clean-label backdoor
732 watermark. *arXiv preprint arXiv:2408.05500*, 2024.
- 733 Tong Xu, Yiming Li, Yong Jiang, and Shu-Tao Xia. Batt: Backdoor attack with transformation-
734 based triggers. In *ICASSP*, pp. 1–5. IEEE, 2023.
735
- 736 Xiong Xu, Kunzhe Huang, Yiming Li, Zhan Qin, and Kui Ren. Towards reliable and efficient
737 backdoor trigger inversion via decoupling benign features. In *ICLR*, 2024.
- 738 Mengxi Ya, Yiming Li, Tao Dai, Bin Wang, Yong Jiang, and Shu-Tao Xia. Towards faithful xai
739 evaluation via generalization-limited backdoor watermark. In *ICLR*, 2023.
740
- 741 Chao-Han Huck Yang, Yun-Yun Tsai, and Pin-Yu Chen. Voice2series: Reprogramming acoustic
742 models for time series classification. In *ICML*, pp. 11808–11819, 2021.
743
- 744 Sheng Yang, Yiming Li, Yong Jiang, and Shu-Tao Xia. Backdoor defense via suppressing model
745 shortcuts. In *ICASSP*, pp. 1–5. IEEE, 2023.
- 746 Sheng Yang, Jiawang Bai, Kuofeng Gao, Yong Yang, Yiming Li, and Shu-Tao Xia. Not all prompts
747 are secure: A switchable backdoor attack against pre-trained vision transformers. In *CVPR*, pp.
748 24431–24441, 2024a.
- 749 Zhou Yang, Bowen Xu, Jie M Zhang, Hong Jin Kang, Jieke Shi, Junda He, and David Lo. Stealthy
750 backdoor attack for code models. *TSE*, 2024b.
751
- 752 Yi Zeng, Si Chen, Won Park, Z Morley Mao, Ming Jin, and Ruoxi Jia. Adversarial unlearning of
753 backdoors via implicit hypergradient. *arXiv preprint arXiv:2110.03735*, 2021a.
754
- 755 Yi Zeng, Won Park, Z Morley Mao, and Ruoxi Jia. Rethinking the backdoor attacks’ triggers: A
frequency perspective. In *ICCV*, pp. 16473–16481, 2021b.

756 Xinyu Zhang, Hanbin Hong, Yuan Hong, Peng Huang, Binghui Wang, Zhongjie Ba, and Kui Ren.
757 Text-crs: A generalized certified robustness framework against textual adversarial attacks. In
758 *S&P*, pp. 2920–2938. IEEE, 2024.

759 Jiachen Zhou, Peizhuo Lv, Yibing Lan, Guozhu Meng, Kai Chen, and Hualong Ma. Dataelixir:
760 Purifying poisoned dataset to mitigate backdoor attacks via diffusion models. In *AAAI*, volume 38,
761 pp. 21850–21858, 2024.

762 Mingli Zhu, Shaokui Wei, Hongyuan Zha, and Baoyuan Wu. Neural polarizer: a lightweight and
763 effective backdoor defense via purifying poisoned features. In *NeurIPS*, pp. 1132–1153, 2023.

766 APPENDIX

767 A THE PROOF OF THEOREM 1

768 **Theorem 1.** Given a K -class pre-trained deep learning model $\mathcal{F}(\cdot) = s(f(\cdot))$ where $s(\cdot)$ is the
769 softmax function and $f(\cdot)$ is the feature extractor, and a pre-processing method $\mathcal{T}(\cdot)$, \mathbf{x} is the data
770 from a specific domain \mathcal{D} (i.e., $\mathbf{x} \sim \mathcal{D}$) and $\tilde{\mathbf{x}} = \mathcal{T}(\mathbf{x}) \sim \tilde{\mathcal{D}}$. Let $\Phi_{\mathcal{D}}(\mathbf{x})$ and $\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}})$ denotes the
771 probability density function of \mathcal{D} and $\tilde{\mathcal{D}}$, we have

$$772 \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \leq 2\alpha\sqrt{K} \cdot \mathcal{W}_1(\mu, \tilde{\mu}), \quad (1)$$

773 where $\mathcal{W}_1(\mu, \tilde{\mu})$ is the Wasserstein-1 distance between μ and $\tilde{\mu}$, μ and $\tilde{\mu}$ are the probability measures
774 of the representations $f(\mathbf{x})$ and $f(\tilde{\mathbf{x}})$, and $\alpha = \max[\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}|\mathbf{x})/\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}})]$.

775 Following similar approaches in (Yang et al., 2021), the proof of Theorem 1 is as follows.

776 *Proof.* Let $[K]$ represents the set of the first K positive integers, i.e., $[K] = \{1, 2, 3, \dots, K\}$. Ac-
777 cording to the definition of mathematical expectation, we have

$$778 \begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \\ &= \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}|\mathbf{x}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &\leq \alpha \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}}, \end{aligned} \quad (2)$$

779 where $\alpha = \max[\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}|\mathbf{x})/\Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}})]$. Assuming $\mathbf{x} \in \mathbb{R}^d$ is a d -dimension vector and \mathbf{x}_i denotes the
780 i -th element of \mathbf{x} , we have

$$781 \|\mathbf{x}\| = \sqrt{\sum_{i=1}^d \mathbf{x}_i^2} \leq \sqrt{d \cdot \max_{i \in [d]} [\mathbf{x}_i^2]} = \sqrt{d} \cdot \max_{i \in [d]} [|\mathbf{x}_i|]. \quad (3)$$

782 Since $\mathcal{F}(\cdot)$ is a K -class pre-trained model, we have

$$783 \begin{aligned} & \alpha \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &\leq \alpha\sqrt{K} \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \max_{k \in [K]} |[F(\mathbf{x})]_k - [F(\tilde{\mathbf{x}})]_k| \cdot \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \alpha\sqrt{K} \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \max_{k \in [K]} |[s(f(\mathbf{x}))]_k - [s(f(\tilde{\mathbf{x}}))]_k| \cdot \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}}. \end{aligned} \quad (4)$$

784 After that, we define k^+ and k^- as the following equations.

$$785 \begin{cases} k^+ = \arg \max_{k \in [K]} [s(f(\mathbf{x}))]_k - [s(f(\tilde{\mathbf{x}}))]_k \\ k^- = \arg \max_{k \in [K]} [s(f(\tilde{\mathbf{x}}))]_k - [s(f(\mathbf{x}))]_k \end{cases}. \quad (5)$$

Because the output of $s(\cdot)$ is a probability logit and the sum total is 1, there exist at least one k_1 such that $[s(f(\mathbf{x}))]_{k_1} - [s(f(\tilde{\mathbf{x}}))]_{k_1} \geq 0$ and also at least one k_2 leading to $[s(f(\tilde{\mathbf{x}}))]_{k_2} - [s(f(\mathbf{x}))]_{k_2} \geq 0$. Therefore,

$$\begin{aligned} & \max_{k \in [K]} |[s(f(\mathbf{x}))]_k - [s(f(\tilde{\mathbf{x}}))]_k| \\ &= \max_{k \in [K]} \{[s(f(\mathbf{x}))]_k - [s(f(\tilde{\mathbf{x}}))]_k, [s(f(\tilde{\mathbf{x}}))]_k - [s(f(\mathbf{x}))]_k\} \\ &\leq [s(f(\mathbf{x}))]_{k^+} - [s(f(\tilde{\mathbf{x}}))]_{k^+} + [s(f(\tilde{\mathbf{x}}))]_{k^-} - [s(f(\mathbf{x}))]_{k^-}. \end{aligned} \quad (6)$$

According to Eq. (6), we have

$$\begin{aligned} & \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \max_{k \in [K]} |[s(f(\mathbf{x}))]_k - [s(f(\tilde{\mathbf{x}}))]_k| \cdot \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &\leq \int_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} ([s(f(\mathbf{x}))]_{k^+} - [s(f(\tilde{\mathbf{x}}))]_{k^+} + [s(f(\tilde{\mathbf{x}}))]_{k^-} - [s(f(\mathbf{x}))]_{k^-}) \cdot \Phi_{\mathcal{D}}(\mathbf{x}) \Phi_{\tilde{\mathcal{D}}}(\tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}} \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [[s(f(\mathbf{x}))]_{k^+} - [s(f(\mathbf{x}))]_{k^-}] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} [[s(f(\tilde{\mathbf{x}}))]_{k^-} - [s(f(\tilde{\mathbf{x}}))]_{k^+}]. \end{aligned} \quad (7)$$

Based on the fact that $[s(\cdot)]_k$ is 1-Lipschitz continuous for any $k \in [K]$ (Gao & Pavel, 2017) and thus $[s(\cdot)]_{k^+} - [s(\cdot)]_{k^-}$ is 2-Lipschitz continuous, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [[s(f(\mathbf{x}))]_{k^+} - [s(f(\mathbf{x}))]_{k^-}] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} [[s(f(\tilde{\mathbf{x}}))]_{k^-} - [s(f(\tilde{\mathbf{x}}))]_{k^+}] \\ &\leq 2 \cdot \sup_{g: \mathbb{R}^K \mapsto \mathbb{R}, \text{Lip}(g) \leq 1} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [g(f(\mathbf{x}))] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} [g(f(\tilde{\mathbf{x}}))]. \end{aligned} \quad (8)$$

Following the Kantorovich-Rubinstein theorem of the dual representation of the Wasserstein-1 distance, finally, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} \|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\tilde{\mathbf{x}})\|_2 \\ &\leq 2\alpha\sqrt{K} \cdot \sup_{g: \mathbb{R}^K \mapsto \mathbb{R}, \text{Lip}(g) \leq 1} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [g(f(\mathbf{x}))] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\mathcal{D}}} [g(f(\tilde{\mathbf{x}}))] \\ &= 2\alpha\sqrt{K} \cdot \mathcal{W}_1(\mu, \tilde{\mu}), \end{aligned} \quad (9)$$

where μ and $\tilde{\mu}$ are the probability measures of the representations $f(\mathbf{x})$ and $f(\tilde{\mathbf{x}})$. \square

B THE PSEUDO-CODE OF REFINE

The pseudo-code of our REFINE optimization process is shown in Algorithm 1.

Algorithm 1 REFINE Optimization Process

Input: The backdoored model \mathcal{F} , the unlabeled benign dataset $\mathbf{D} = \{\mathbf{x}_i\}_{i=1}^M$, the randomly initialized input transformation module $\mathcal{T}(\cdot, \theta)$, the specified output mapping module $\mathcal{M}(\cdot)$.

Output: The input transformation module parameters θ .

- 1: **for** data batches $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$ in \mathbf{D} **do**
- 2: Obtain the transformed input $\tilde{\mathbf{X}} = \mathcal{T}(\mathbf{X}, \theta)$.
- 3: Obtain the original model output $\tilde{\mathbf{Y}} = \mathcal{F}(\tilde{\mathbf{X}})$.
- 4: Obtain the mapped output $\mathbf{Y} = \mathcal{M}(\tilde{\mathbf{Y}})$.
- 5: Obtain the predicted labels $\bar{\mathbf{Y}} = \arg \max(\mathcal{F}(\mathbf{X}))$.
- 6: Compute the supervised contrastive loss $\mathcal{L}_{sup}(\tilde{\mathbf{X}}, \bar{\mathbf{Y}})$.
- 7: Compute the cross-entropy loss $\mathcal{L}_{ce}(\bar{\mathbf{Y}}, \mathbf{Y})$.
- 8: Optimize θ with the composite loss: $\arg \min_{\theta} \mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{sup}$.

9: **return** θ

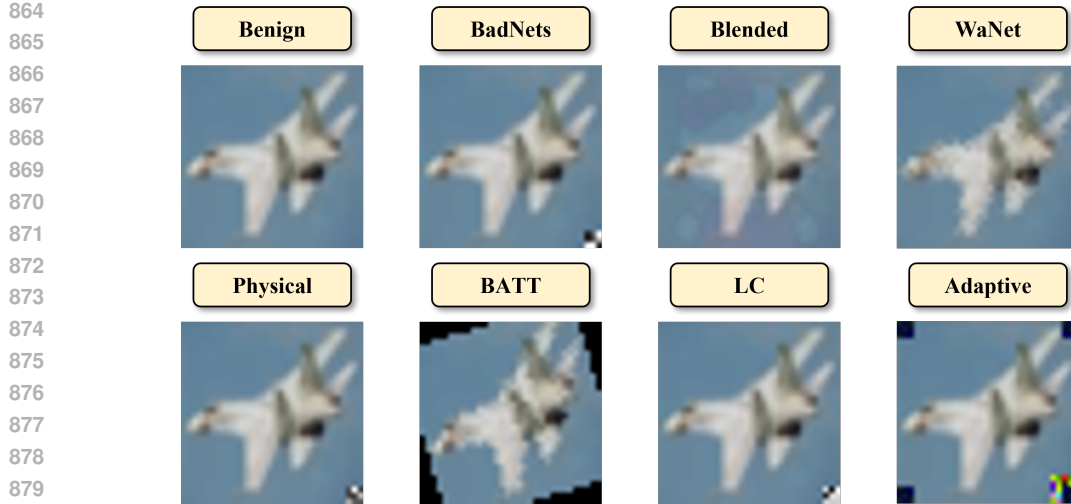


Figure 5: The illustration of the adopted backdoor attacks.

C IMPLEMENTATION DETAILS

C.1 DETAILS OF THE EXPERIMENTAL SETTINGS

Details of Datasets. (1) *CIFAR-10*. The CIFAR-10 dataset (Krizhevsky et al., 2009) contains 50,000 training samples and 10,000 testing samples in total. The dataset has 10 classes and each class has 5,000 training samples and 1,000 testing samples. The size of each image sample is $3 \times 32 \times 32$. (2) *ImageNet*. The ImageNet dataset (Deng et al., 2009) consists of 1,000 classes containing over 14 million manually annotated images. In this paper, we select a subset with 50 different classes and each class contains 500 training samples and 100 testing samples with size $3 \times 224 \times 224$.

Details of Training Backdoored Models. We utilize the SGD with a momentum of 0.9 and a weight decay of 5×10^{-4} as the optimizer for training all backdoored DNNs. The batch size is set to 128 on both of CIFAR-10 and ImageNet. We set the initial learning rate as 0.1 and train all models for 200 epochs, with the learning rate reduced by a factor of 0.1 at the 100-th and 150-th epoch.

Details of Optimization. For training the input transformation module, we employ SGD with a momentum of 0.9 and a weight decay of 5×10^{-4} as the optimizer. The initial learning rate is set to 0.01, and the batch size is set to 128 and 32 for ImageNet. The input transformation module is trained for 150 epochs, with the learning rate decayed by a factor of 0.8 at the 100-th and 130-th epochs. For the training loss function, we set the temperature parameter as 0.1. For the output mapping module, we randomly assign a hard-coded remapping function before each defense.

Computational Resources. In our implementations, we utilize PyTorch as the deep learning framework. All our experiments are implemented with RTX 3090 GPUs.

C.2 DETAILS OF THE ADOPTED BACKDOOR ATTACKS

In our experiments, we adopt 7 representative backdoor attacks to evaluate the defense performance of our REFINE and other baseline backdoor defense methods. We implement the first six backdoor attacks utilizing BackdoorBox (Li et al., 2022b), while the implementation of Adaptive-Patch is derived from the open-source code provided in its original paper. In this section, we provide a detailed introduction to these backdoor attacks, as follows.

- **BadNets:** Gu et al. (2019) introduced the earliest poisoning-based backdoor attack that aims to poison the training dataset using a visible, distinctive pixel pattern. In this paper, we utilize a 3×3 random square as the trigger pattern on the bottom right of samples in CIFAR-10 and a 32×32 square on ImageNet.
- **Blended:** To evade human visual detection of poisoned samples, Chen et al. (2017) designed a covert data poisoning method known as Blended, which attempts to embed triggers implicitly

Table 5: The performance (%) of BDMAE, BTI-DBF(P) and REFINE against Clean-image. The best results are **boldfaced**.

Defense→	No Defense			BDMAE			BTI-DBF(P)			REFINE		
	BA	ASR-n	ASR-m	BA	ASR-n	ASR-m	BA	ASR-n	ASR-m	BA	ASR-n	ASR-m
Clean-image	87.78	82.93	4.73	86.73	58.54	4.73	78.95	86.59	9.53	88.94	53.57	4.63

Table 6: The performance (%) of REFINE and ZIP against four different attacks. The best results are **boldfaced**.

Defense→	No Defense		ZIP		REFINE	
	BA	ASR	BA	ASR	BA	ASR
BadNets	91.18	100.00	84.22	5.53	90.50	1.05
Blended	90.64	98.18	84.68	8.64	90.30	1.00
WaNet	91.29	99.91	85.19	15.46	90.64	1.93
PhysicalBA	93.67	100.00	85.07	10.91	91.17	0.78

within the samples. In this paper, we utilize an image of Hello Kitty as the trigger pattern and set the blending rate to 0.3 across both datasets.

- **WaNet**: WaNet (Nguyen & Tran, 2021) is another type of invisible backdoor attack that employs a warp-based trigger. We follow the default settings in its original paper.
- **PhysicalBA**: Li et al. (2021c) demonstrated that DNNs applied in physical scenarios could also be vulnerable to backdoor threats and proposed backdoor attacks that simulate physical transformations. In this paper, we follow the default settings outlined in its original paper.
- **BATT**: Xu et al. (2023) noted that simple transformations specific to samples can pose significant backdoor threats to models and introduced the Backdoor Attack with Transformation-based Triggers (BATT). In this paper, we adhere to the default settings established in their original work.
- **LabelConsistent**: To address the issue of easily identifiable mislabeled poisoned data in poisoning datasets, Turner et al. (2019) proposed clean-label backdoor attacks, which aim to poison samples of specific classes to inject backdoors. We employ projected gradient descent (PGD) to generate adversarial samples, setting the maximum perturbation size to $\epsilon = 8$. The trigger patterns utilized are identical to those employed in BadNets.
- **Adaptive-Patch**: Qi et al. (2023) observed that models trained on poisoned datasets often learn distinct latent representations for poisoned and clean samples, and they proposed adaptive backdoor attacks to mitigate this separation phenomenon. In this paper, we follow the default settings utilized in its original paper.

The poisoned samples of these backdoor attacks are depicted in Figure 5.

D ADDITIONAL RESULTS ON MORE ATTACKS AND DEFENSES

D.1 RESULTS AGAINST CLEAN-IMAGE ATTACK

In this section, we test the clean-image backdoor attack (Rong et al., 2024) (dubbed 'Clean-image'). We train a ResNet18 backdoor model on CIFAR10 and applied defenses using BDMAE, BTI-DBF(P), and REFINE. For each defense method, we test the model's benign accuracy (BA), natural attack success rate (ASR-n), and manual attack success rate (ASR-m). ASR-n represents the ASR of poisoned samples that naturally contain the backdoor trigger, while ASR-m represents the ASR of manually generated poisoned samples.

As shown in Table 5, existing pre-processing-based defenses are unable to effectively reduce the ASR-n of the clean-image attack. Compared to BDMAE and BTI-DBF(P), REFINE shows better defense performance. However, we believe that the Clean-image attack, from certain perspectives, extends beyond the typical scope of backdoor attacks. Specifically, the stealthiness of backdoor

Table 7: The performance (%) of REFINE and BDMAE on different model architectures. The best results are **boldfaced**.

Model↓	Defense→ Attack↓	No Defense		BDMAE		REFINE	
		BA	ASR	BA	ASR	BA	ASR
ResNet18	BadNets	91.18	100.00	90.48	14.81	90.50	1.05
	WaNet	91.29	99.91	89.87	99.93	90.64	1.93
	Adaptive	89.62	100.00	89.40	49.18	90.54	1.23
ResNet50	BadNets	91.91	100.00	91.04	10.99	90.71	1.53
	WaNet	91.70	99.98	89.83	99.89	91.09	0.35
	Adaptive	89.59	85.11	89.06	35.91	90.05	2.19
VGG16	BadNets	84.44	99.36	84.25	18.32	86.86	1.62
	WaNet	84.75	99.15	83.36	99.25	86.41	2.39
	Adaptive	84.98	99.99	84.69	40.09	86.63	2.04
DenseNet121	BadNets	86.40	99.99	86.05	11.85	89.44	0.96
	WaNet	86.31	98.77	85.42	98.91	88.74	0.88
	Adaptive	85.16	100.00	84.45	45.36	88.74	0.35
InceptionV3	BadNets	90.46	99.97	90.61	80.51	91.03	0.75
	WaNet	90.09	99.73	89.64	99.76	91.01	0.54
	Adaptive	88.58	13.53	88.54	13.52	90.36	0.54

attacks usually requires the model to behave normally on benign samples, while the clean-image attack may cause misclassifications on some benign samples (i.e., naturally poisoned samples), which significantly impacts the model’s normal inference. We will discuss how to defend against this type of attack in our future work.

D.2 RESULTS ON ZIP

In this section, we conducted additional experiments using the attack methods supported by the original ZIP repository. Specifically, we trained ResNet-18 backdoor models on the CIFAR-10 dataset using four different attack methods, including BadNets, Blended, WaNet, and PhysicalBA.

As shown in Table 6, compared to ZIP, REFINE demonstrates better defense performance against all four attack methods. This demonstrates the effectiveness of REFINE again.

E ADDITIONAL ABLATION STUDY

E.1 RESULTS ON DIFFERENT MODEL ARCHITECTURES

In this section, we conduct experiments on five different network structures, including ResNet-50 (He et al., 2016), VGG-16 (Simonyan, 2014), DenseNet-121 (Huang et al., 2017), and Inception-V3 (Szegedy et al., 2016). We select three representative types of backdoor attacks, including BadNets, WaNet, and Adaptive-Patch (dubbed ‘Adaptive’). We conduct experiments on the CIFAR-10 dataset. We compare the defense performance of our REFINE with the most advanced transformation-based defense (i.e., BDMAE).

As shown in Table 7, REFINE effectively defends against three representative attacks across five different network architectures, significantly outperforming BDMAE. Specifically, under the REFINE defense, the benign accuracy (BA) drop is less than 1.5%, with some cases showing an increase in BA. Meanwhile, the backdoor attack success rate (ASR) is reduced to below 3%. The additional experimental results verify the effectiveness of REFINE.

E.2 EFFECT OF THE UNLABELED BENIGN DATASET SIZE

In this section, we evaluate the defense performance of REFINE under different sizes of the unlabeled benign dataset. We train a backdoored classification model on CIFAR-10 using the BadNets

1026 Table 8: Performance (%) of REFINE under different sizes of the unlabeled benign dataset.

1027

Proportion→	No Defense		100%		80%		60%		40%		20%	
Attack↓	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNets	92.31	100	91.20	0.86	90.22	1.05	89.53	1.21	87.81	1.11	83.93	2.21

1032

1033 Table 9: Performance (%) of REFINE under different values of temperature parameters λ .

1034

$\lambda \rightarrow$	No Defense		1.0		0.8		0.6		0.4		0.2	
Attack↓	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNets	91.74	100	90.83	0.92	90.87	0.76	90.60	0.60	91.03	0.51	90.69	1.27

1038

1039 Table 10: Performance (%) of REFINE under different number of channels in UNet hidden layers.

1040

Channels→	No Defense		32		48		64		80	
Attack↓	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNets	91.74	100	89.49	1.09	90.61	0.64	90.18	1.43	91.07	0.78

1042

1043 attack on a ResNet-18 architecture. For defense, we use different proportions (100% to 20%) of the
 1044 CIFAR-10 dataset as the unlabeled benign dataset. As shown in Table 8, the results indicate that as
 1045 the number of unlabeled samples decreases, the BA of REFINE experiences a slight decline, while
 1046 the ASR remains consistently low.

1047

1052 E.3 EFFECT OF THE SCALAR TEMPERATURE PARAMETER λ

1053

1054 In this section, we evaluate the defense performance of REFINE under different values of tempera-
 1055 ture parameters λ . The attack setup is consistent with that in Section E.2. During the defense, we
 1056 test various temperature parameters ranging from 1 to 0.2. As shown in Table 9, the results indicate
 1057 that the value of temperature parameter has minimal impact on the defense performance of REFINE.

1058

1059 E.4 EFFECT OF THE NUMBER OF CHANNELS IN UNET HIDDEN LAYERS

1060

1061 In this section, we evaluate the defense performance of REFINE using UNet models with varying
 1062 numbers of hidden layer channels. Specifically, the dimensionality of the encoded features can be
 1063 adjusted by altering the number of output channels in the first layer of the UNet encoder. The
 1064 attack setup is consistent with that in Section E.2. For the defense, we tested different channel
 1065 numbers, including 32, 48, 64, and 80. As shown in Table 10, the number of channels in the UNet
 1066 hidden layers has minimal impact on the defense performance of REFINE, with both BA and ASR
 1067 remaining at an optimal level.

1068

1069 E.5 EFFECT OF THE DATA DISTRIBUTION USED FOR DEFENSE

1070

1071 In our main experiments, we assume that the defender can acquire independent and identically dis-
 1072 tributed (i.i.d.) unlabeled datasets. In this section, we explore the defense performance under differ-
 1073 ent data distributions. We train a ResNet18 model on the CIFAR10 dataset using the BadNets attack.
 1074 For defense, we trained the input transformation module of REFINE using CINIC10 (Darlow et al.,
 1075 2018), a dataset with the same categories as CIFAR10 but a different data distribution.

1076

1077 As shown in Table 11, REFINE is still highly effective in reducing the attack success rate (ASR
 1078 $< 1.5\%$) while maintaining the model’s benign accuracy (BA drop $< 3\%$). This favorable result is
 1079 due to the fact that REFINE first assigns pseudo-labels to the unlabeled benign samples using the
 original model, and then trains the input transformation module based on these pseudo-labels.

Table 11: The performance (%) of REFINE in scenarios with different data distribution.

Defense→	No Defense		REFINE	
	Attack↓	BA	ASR	ASR
BadNets	91.18	100.00	88.39	1.40

Table 12: The performance (%) of REFINE and T-MR. The best results are **boldfaced**.

Defense→	No Defense		T-MR		REFINE	
	Attack↓	BA	ASR	BA	ASR	ASR
BadNets	91.18	100.00	75.51	3.36	90.50	1.05
WaNet	91.29	99.91	74.49	25.76	90.64	1.93
Adaptive	92.54	99.93	75.49	5.87	90.87	1.76

E.6 EFFECT OF IMPROVED TRANSFORMATION MODULE

In this section, we conduct additional defense experiments using traditional model reprogramming methods (Elsayed et al., 2019) (dubbed ‘T-MR’). We select three representative types of backdoor attacks, including BadNets, WaNet, and BATT. We train backdoor ResNet18 models on the CIFAR-10 dataset. We compare the defense performance of REFINE with T-MR.

As shown in Table 12, the T-MR defense has a significant impact on the model’s BA (BA drop > 15%) but fails to effectively reduce the ASR under the WaNet attack. This is because traditional model reprogramming methods only add a universal adversarial perturbation around the image, while the trigger pattern remains unchanged on the backdoor image to some extent.

F REFINE IN THE BLACK-BOX SCENARIO

In our main experiments, we assume that we can obtain white-box access to the pre-trained backdoored models. In this section, we investigate how to implement our REFINE in the black-box scenario where the defender can only get black-box access to the backdoored model. In the black-box scenario, only the class confidence scores are accessible and it is hard to calculate the gradients to optimize the REFINE modules. To tackle the aforementioned challenge, we leverage the surrogate model technique. Specifically, we distill a surrogate model from the original black-box model using an unlabeled dataset D . We employ the mean squared error (MSE) loss to align the output confidence scores between the black-box model $\mathcal{F}(\cdot)$ and the surrogate model $\mathcal{F}_s(\cdot)$, as follows.

$$\mathcal{L}_{distill} = \frac{1}{|D|} \sum_{\mathbf{x} \in D} [\mathcal{F}(\mathbf{x}) - \mathcal{F}_s(\mathbf{x})]^2. \quad (10)$$

The surrogate model is then leveraged to replace the pre-trained model in our REFINE and optimize the input transformation module. Subsequently, the trained input transformation and output mapping modules are subsequently applied to the original black-box model.

To validate the feasibility of our REFINE in the black-box scenario, we employ the backdoored ResNet-50 pre-trained on the CIFAR-10 dataset as the black-box model and ResNet-18 as the surrogate model. As shown in Table 13, we evaluate both the black-box original model and the surrogate model in terms of BA and ASR before and after applying the REFINE defense. The ASRs of our REFINE are all below 4%. The results indicate that even though the input transformation module is trained using the surrogate model, our REFINE is still capable of achieving high performance of backdoor defense for the black-box original model.

G THE OVERHEAD OF OUR REFINE

In this section, we evaluate the overhead of our REFINE. Specifically, we measure the training time of the input transformation module and the model inference time on the CIFAR-10 and ImageNet

Table 13: Performance (%) of REFINE in defending against attacks in black-box scenarios.

Defense→	No Defense				REFINE			
	Black-box		Surrogate		Surrogate		Black-box	
Model→	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNets	90.60	100	91.20	1.24	88.21	0.92	88.17	0.36
Blended	91.08	96.94	90.69	2.23	88.34	0.62	87.75	0.18
WaNet	91.50	99.93	90.92	99.84	88.77	3.37	87.44	0.04
Physical	93.61	100	92.21	2.56	90.18	1.52	89.84	2.23
BATT	93.24	99.89	92.76	4.30	90.86	2.01	89.21	3.72
LC	91.95	93.06	91.53	1.11	89.04	0.87	88.69	1.05
Adaptive	90.15	100	90.36	1.57	88.41	0.32	87.91	0.44

Table 14: The performance (%) of FT and FT+REFINE on ResNet18.

Defense→	No Defense		FT		FT+REFINE	
	BA	ASR	BA	ASR	BA	ASR
BadNets	91.18	100.00	91.89	91.67	90.42	0.87

datasets using the ResNet-18 model. We employ a UNet with 32 hidden layer channels as the structure for the input transformation module. During training, we employ SGD with a momentum of 0.9 and a weight decay of 5×10^{-4} as the optimizer. The initial learning rate is set to 0.01, with a batch size of 128 for CIFAR-10 and 32 for ImageNet. The input transformation module is trained for 150 epochs, with the learning rate decaying by a factor of 0.8 at the 100-th and 130-th epochs. For the training loss function, the temperature parameter is set to 0.1. We conduct all training using a single RTX 3090 GPU. For the output mapping module, a hard-coded remapping function is randomly assigned before each defense.

The results indicate that for the CIFAR-10, REFINE requires only 1 hour to train the input transformation module. Once training is complete, inference on 10,000 images takes 6.31 seconds. Moreover, on a subset of the ImageNet with 50 classes, REFINE requires 5 hours for training the input transformation module. After training, inference on 2,500 images takes 12.33 seconds. Although REFINE introduces some additional overhead, we believe this cost is reasonable and acceptable.

H COMBINING REFINE WITH EXISTING DEFENSES

Arguably, our method can be used in conjunction with existing (model reconstruction-based) defenses to further enhance their effectiveness. To demonstrate this, we first applied model fine-tuning defense (dubbed 'FT') to a ResNet18 model subjected to the BadNets attack on CIFAR-10, followed by the REFINE defense. As shown in Table 14, the FT+REFINE defense effectively reduces the backdoor ASR while maintaining the model's BA.

I RELATED WORK

I.1 BACKDOOR ATTACK

Visible Backdoor Attacks. This type of attack typically employs patterns that are visible to humans as triggers. BadNets (Gu et al., 2019) is the first backdoor attack technique that injects samples with simple but visually noticeable patterns into the training data, such as white squares or specific marks. Li et al. (2021c) then proposed a transformation-based enhancement that strengthens the attack's resilience and establishes its applicability to physical scenarios. To address the issue of latent feature separation in backdoor attacks, Qi et al. (2023) employed asymmetric trigger planting strategies and developed adaptive backdoor poisoning attacks. Besides, Gao et al. (2023) revealed that clean-label attacks were difficult due to the conflicting effects of 'robust features' in poisoned

1188 samples and proposed a simple yet effective method to improve these attacks by targeting ‘hard’
 1189 samples instead of random ones.

1190 **Invisible Backdoor Attacks.** To enhance the stealth of backdoor attacks, Chen et al. (2017) was
 1191 the first to introduce the use of triggers that are imperceptible to humans, aiming to evade detection
 1192 by basic data filtering techniques or human inspection. They proposed a blending strategy that
 1193 generates poisoned images by subtly merging the backdoor trigger with benign images. After that,
 1194 a series of studies focused on designing invisible backdoor attacks. WaNet (Nguyen & Tran, 2021)
 1195 and ISSBA (Li et al., 2021d) employed warping-based triggers and perturbation-based triggers,
 1196 respectively, introducing sample-specific trigger patterns during training; LIRA (Doan et al., 2021)
 1197 formulated the learning of an optimal, stealthy trigger injection function as a non-convex constrained
 1198 optimization problem, where the trigger generator function is trained to manipulate inputs using
 1199 imperceptible noise; BATT (Xu et al., 2023) utilized images rotated to a specific angle as triggers,
 1200 representing a new attack paradigm where triggers extend beyond basic pixel-wise manipulations.

1201 A few existing literature also provided novel and comprehensive discussions on backdoor attacks
 1202 from various perspectives and domains, including CLIP (Liang et al., 2024), diffusion models (Chou
 1203 et al., 2024), 3D point clouds (Wei et al., 2024), pre-trained ViT (Yang et al., 2024a), code genera-
 1204 tion (Yang et al., 2024b), and federated learning (Shao et al., 2024). Moreover, some existing works
 1205 also explore utilizing the backdoor attack for good purposes, such as copyright protection (Liu et al.,
 1206 2021; Li et al., 2022a) and explainable artificial intelligence (XAI) evaluation (Ya et al., 2023).

1208 I.2 BACKDOOR DEFENSES

1209
 1210 Currently, there are various backdoor defense methods designed to mitigate backdoor threats. These
 1211 methods can generally be divided into three main paradigms (Li et al., 2022b): (1) trigger-backdoor
 1212 mismatch, which primarily refers to pre-processing-based defenses (Liu et al., 2017; Li et al., 2021c;
 1213 Shi et al., 2023). (2) backdoor elimination (Li et al., 2021b; Zeng et al., 2021a; Huang et al., 2022;
 1214 Xu et al., 2024; Hayase & Kong, 2020; Zeng et al., 2021b), such as model reconstruction (Li et al.,
 1215 2021b; Zeng et al., 2021a), poison suppression (Huang et al., 2022; Xu et al., 2024), and training
 1216 sample filtering (Hayase & Kong, 2020; Zeng et al., 2021b). (3) trigger elimination, also known as
 1217 testing sample filtering (Gao et al., 2019; Javaheripi et al., 2020).

1218 **Pre-processing-based Defenses.** These methods incorporate a pre-processing module prior to feed-
 1219 ing samples into DNNs, altering the trigger patterns present in the samples. Consequently, the mod-
 1220 ified triggers no longer align with the hidden backdoor, thereby preventing the backdoor activation.
 1221 AutoEncoderDefense (Liu et al., 2017) is the first pre-processing-based backdoor defense by em-
 1222 ploying a pre-trained autoencoder as the pre-processing module. Based on the idea that trigger
 1223 regions have the most significant impact on predictions, Februus (Doan et al., 2020) effectively mit-
 1224 igates backdoor attacks by removing potential trigger artifacts and reconstructing inputs, all while
 1225 preserving performance for both poisoned and benign samples. Li et al. (2021c) observed that
 1226 poisoning-based attacks with static trigger patterns degrade sharply with slight changes in trigger
 1227 appearance or location and proposed spatial transformations (e.g., shrinking, flipping) as an efficient
 1228 defense with minimal computational cost. Deepsweep (Qiu et al., 2021) proposes a unified defense
 1229 that (1) fine-tunes the infected model using a data augmentation policy to remove backdoor effects
 1230 and (2) pre-processes input samples with another augmentation policy to disable triggers during
 1231 inference. Recently, many pre-processing-based defenses utilize the generative model, such as the
 1232 diffusion model and the masked autoencoder, to purify the suspicious samples. ZIP (Shi et al., 2023)
 1233 applies linear transformations, such as blurring, to poisoned images to disrupt backdoor patterns and
 1234 subsequently employs a pre-trained diffusion model to recover the semantic information lost during
 1235 the transformation. BDMAE (Sun et al., 2023) detects potential triggers in the token space by eval-
 1236 uating image structural similarity and label consistency between test images and MAE restorations,
 1237 refines these results based on trigger topology, and finally adaptively fuses the MAE restorations
 1238 into a purified image for prediction. DataElixir (Zhou et al., 2024) detects target labels by quanti-
 1239 fying distribution discrepancies, selects purified images based on pixel and feature distances, and
 determines their true labels by training a benign model.

1240 **Backdoor Elimination Defenses.** In contrast to pre-processing-based defenses, backdoor elimina-
 1241 tion methods typically mitigate backdoor threats by directly modifying model parameters or prevent
 backdoor injection by controlling the model training process. Li et al. (2021a) identified two key

weaknesses of backdoor attacks: 1) models learn backdoored data significantly faster than clean data, and 2) the backdoor task is associated with a specific target class. Consequently, they proposed Anti-Backdoor Learning (ABL), which introduces a two-stage gradient ascent mechanism: 1) isolating backdoor examples in the early training phase, and 2) breaking the correlation between backdoor examples and the target class in the later training phase. Inspired by the phenomenon where poisoned samples tend to cluster together in the feature space of the attacked DNN model, Huang et al. (2022) proposed a novel backdoor defense by decoupling the original end-to-end training process into three stages. Yang et al. (2023) removed backdoors by suppressing the skip connections in key layers identified by their method and fine-tuned these layers to restore high BA and further reduce the ASR. [Neural Polarizer \(Zhu et al., 2023\) achieved effective defense by training an additional linear transformation, called neural polarizer, using only a small portion of clean data without modifying the model parameters.](#) Xu et al. (2024) discovered that even in the feature space, the triggers generated by existing BTI methods differ significantly from those used by the adversary. Consequently, they proposed BTI-DBF, which decouples benign features instead of directly decoupling backdoor features. This method primarily involves two key steps: (1) decoupling benign features, and (2) triggering inversion by minimizing the differences between benign samples and their generated poisoned versions while maximizing the differences of the remaining backdoor features.

Trigger Elimination Defenses. These defenses filter out malicious samples during the inference process rather than during training. As a result, the deployed model exclusively predicts benign test samples or purified attack samples, thereby preventing backdoor activation by removing trigger patterns. STRIP (Gao et al., 2019) perturbs the input samples and observes the randomness in predicted classes from the deployed model for these perturbed inputs. If the entropy of the predicted classes is low, this violates the input-dependence characteristic of a benign model, indicating the presence of malicious features within the input. Du et al. (2020) demonstrated that applying differential privacy can enhance the utility of outlier detection and novelty detection, and further extended this approach for detecting poisoned samples in backdoor attacks. Besides, CleanNN (Javaheripi et al., 2020) leverages dictionary learning and sparse approximation to characterize the statistical behavior of benign data and identify triggers, representing the first end-to-end framework capable of online mitigation against backdoor attacks in embedded DNN applications.

I.3 MODEL REPROGRAMMING

Elsayed et al. (2019) first proposed adversarial reprogramming, which aims to repurpose a classifier trained on ImageNet-1K for tasks such as classifying CIFAR-10 and MNIST images and counting the number of squares in an image. BAR (Tsai et al., 2020) extended model reprogramming to black-box scenarios and applied it to the bio-medical domain. Driven by advancements in deep speech processing models and the fact that speech data is a univariate time signal, Voice2Series (Yang et al., 2021) learns to reprogram acoustic models for time series classification and output label mapping through input transformations. Neekhara et al. (2022) analyzed the feasibility of adversarially repurposing image classification neural networks for natural language processing (NLP) and other sequence classification tasks. They developed an effective adversarial program that maps a series of discrete tokens onto an image, which can then be classified into the desired category by an image classification model. Li et al. (2023b) found that combining Visual Prompting (VP) with PATE—a state-of-the-art differential privacy training method that utilizes knowledge transfer from a team of teachers—achieves a cutting-edge balance between privacy and practicality with minimal expenditure on privacy budget. More Recently, a novel application (Dey & Nair, 2024) of model reprogramming repurposed models originally designed for able-bodied individuals to predict joint movements in amputees, significantly enhancing assistive technologies and improving mobility for amputees. Currently, model reprogramming has been shown to outperform transfer learning and training from scratch in many applications (Tsai et al., 2020; Yang et al., 2021; Vinod et al., 2023), without altering the original model’s parameters.

J THE VISUALIZATION OF THE TRANSFORMED SAMPLES \tilde{x}

In this section, we visualize the transformed benign and poisoned samples \tilde{x} generated by the UNet of our REFINE. We train a backdoored ResNet-18 model on CIFAR-10 using the BadNets attack with a specified $3 * 3$ trigger patterns at the bottom right corner of images, and the hard-coded

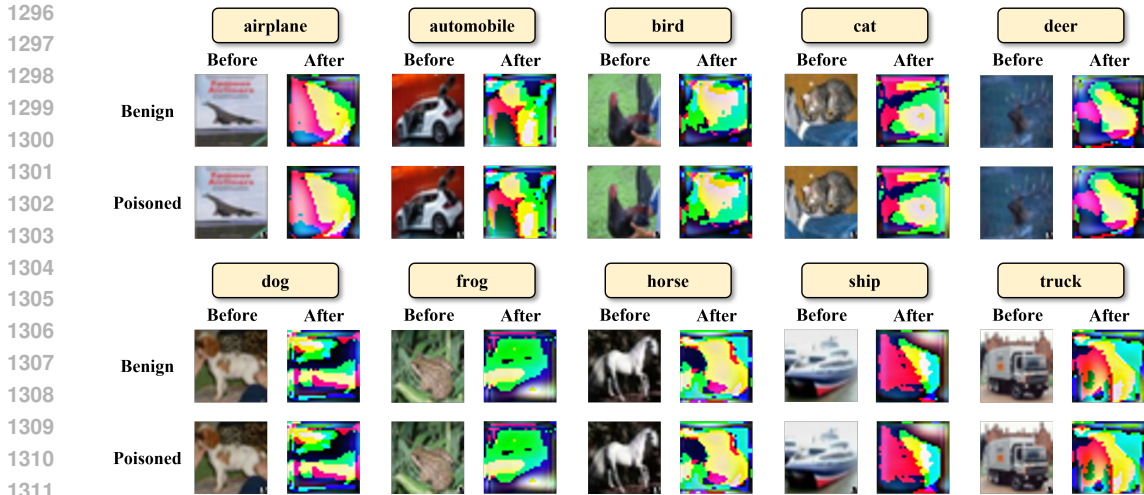


Figure 6: The visualization of transformed samples \tilde{x} . We display the benign and poisoned samples and transformed benign and poisoned samples for each class. For each class of small areas, the upper left corner represents the benign sample, the upper right corner represents the transformed benign sample, the bottom left corner represents the poisoned sample and the bottom right corner represents the poisoned sample after transformations.

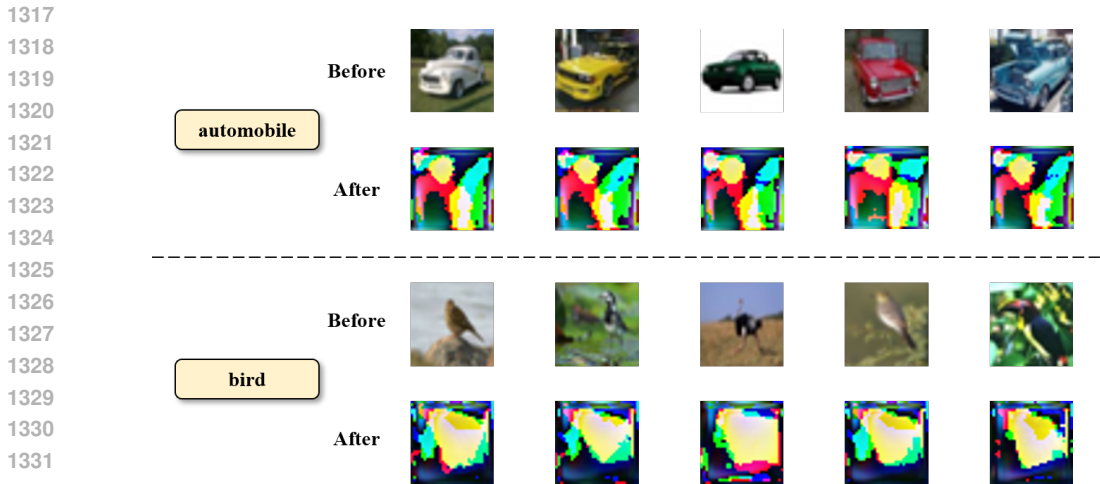


Figure 7: The visualization of transformed samples \tilde{x} for the classes "automobile" and "bird" of CIFAR-10. For each class, we display five input images and their corresponding transformed images.

remapping function f_L of the output mapping module \mathcal{M} is defined as follows:

$$f_L = \tilde{l} \mapsto l, \tag{11}$$

where

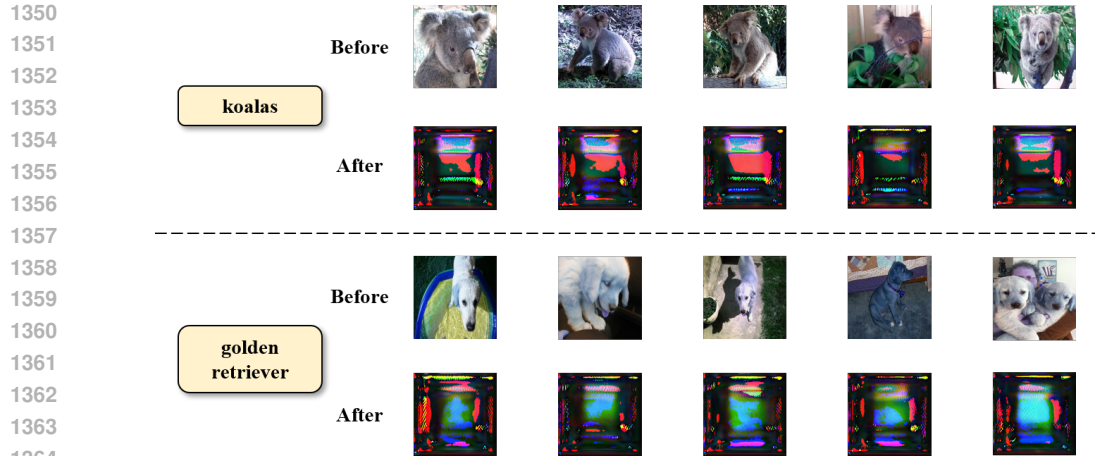
$$\tilde{l} = \{airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck\}, \tag{12}$$

and

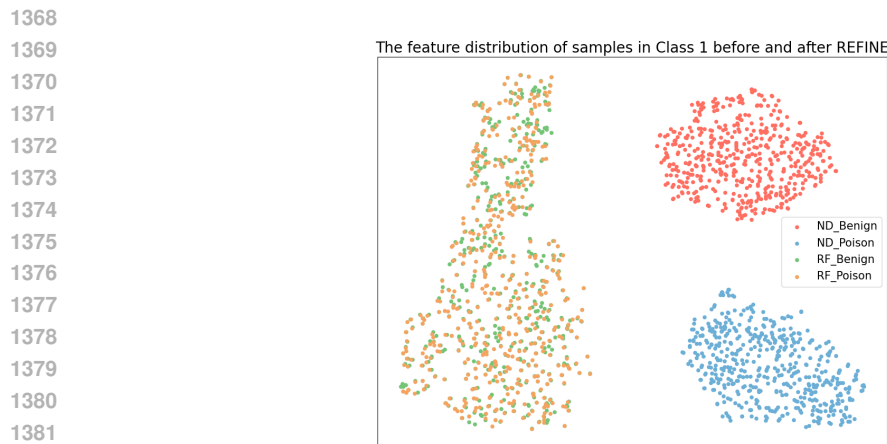
$$l = \{cat, deer, automobile, ship, frog, bird, horse, truck, airplane, dog\}. \tag{13}$$

As shown in Figure 6, for both benign and poisoned samples, the transformed sample patterns are very similar, and the transformed pattern of the poisoned sample effectively removes the trigger. This further illustrates the effectiveness of our REFINE in mitigating backdoor threats.

As shown in Figure 7 and 8, samples from the same class exhibit visual similarities after transformation. However, the transformed samples do not contain any human-recognizable information. This



1365 Figure 8: The visualization of transformed samples \tilde{x} for the classes "koalas" and "golden retriever" of ImageNet. For each class, we display five input images and their corresponding transformed images.



1383 Figure 9: The t-SNE plots of the feature distribution of samples in Class 1 before and after REFINE. ND_Benign and ND_Poison represent the features of benign and poisoned samples under the No Defense (ND) scenario, respectively. RF_Benign and RF_Poison represent the features of benign and poisoned samples after applying REFINE, respectively.

1388 phenomenon occurs because the input transformation module maps the samples to a new benign feature space, and the constraint imposed by the supervised contrastive loss ensures that transformed samples from the same class exhibit more similar benign features.

1392 K THE VISUALIZATION OF THE FEATURE DISTRIBUTION BEFORE AND

1393 AFTER REFINE

1395

1396 In this section, we visualize the changes in the feature distribution of the input samples before and after REFINE. Specifically, we trained a backdoor ResNet-18 model on CIFAR-10 using the BadNets attack and extracted the features from the input of the model's fully connected (FC) layer as the feature values of the input samples.

1400 As shown in Figure 9, before applying REFINE, the feature distributions of benign and poisoned samples are clustered in two distinct locations. After applying REFINE, the feature distributions of benign and poisoned samples are interwoven and clustered in the same new location. This indicates that REFINE effectively removes the trigger patterns from the poisoned samples and maps samples of the same class to a new benign feature distribution.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

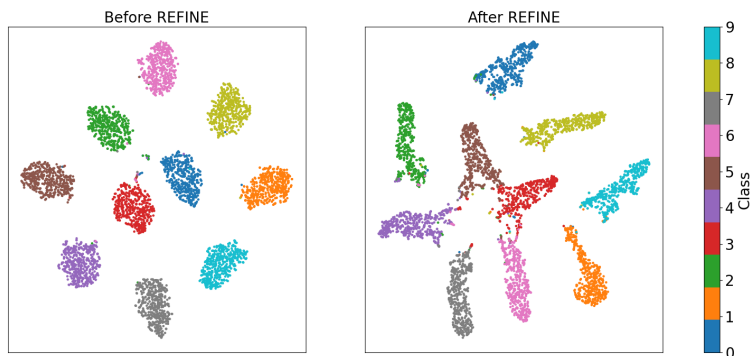


Figure 10: The t-SNE plots of the feature distribution of benign samples from different classes, both before and after REFINE.

As shown in Figure 10, before applying REFINE, the benign samples of each class form distinct clusters in the feature space. After applying REFINE, the benign samples, adjusted by the input transformation module and output mapping module, form new clusters in different positions. This empirically demonstrates that REFINE is able to maintain the model’s benign accuracy.

L SOCEITAL IMPACT

This paper aims to design an effective and efficient backdoor defense method and have a positive societal impact. Specifically, we propose a novel pre-processing-based backdoor defense method, REFINE, based on model reprogramming. REFINE can mitigate the backdoor behaviors injected into the third-party pre-trained models. Therefore, our REFINE can assist in ensuring the stable and reliable operation of the AI models, mitigating the potential threat of backdoors, and facilitating the reuse and deployment of the models. Moreover, the application of our REFINE may also facilitate the emergence of new business models such as model trading.

On the other hand, in this paper, we propose to leverage the model reprogramming techniques to build the input transformation and output mapping modules to mitigate the backdoors. The insight of our method can also be applied to the use of the pre-trained model in an unauthorized way. For instance, an adversary might use the model for an unauthorized task via model reprogramming, leading to copyright infringement (Shao et al., 2025; Wang et al., 2022a). However, we argue that the negative societal impact is negligible. The model developer can employ several existing protection methods, such as non-transfer learning (Wang et al., 2022a), to prevent such misbehaviors. Moreover, although we do not find effective adaptive attacks against our REFINE, an adversary may design a more advanced adaptive attack to circumvent our proposed method since its effectiveness lacks of theoretical guarantees. Even so, the model users and developers can still prevent the backdoor threat from the source by only using trusted pre-trained models.

M POTENTIAL LIMITATIONS AND FUTURE DIRECTIONS

Firstly, as outlined in our threat model, the goal of our defense is to protect against pre-trained models from third-party platforms. Specifically, similar to other baseline methods, we assume that the defender possesses a certain amount of unlabeled sample datasets. To explore the effectiveness of REFINE in few-shot scenarios, we conduct additional experiments using 10% unlabeled clean data. We apply the REFINE defense to a ResNet-18 model trained on the CIFAR-10 dataset, which is subjected to the BadNets attack. In this case, the unlabeled training set for REFINE used only 10% of the CIFAR-10 training set.

As shown in Table 15, even with only 10% unlabeled data, REFINE is still effective to some extent. REFINE effectively reduces the ASR, although it does have some impact on the model’s BA.

Table 15: The performance (%) of REFINE in the 10% unlabeled data scenario on ResNet18.

Defense→	No Defense		REFINE	
	Attack↓	BA	ASR	ASR
BadNets	91.18	100.00	78.02	2.90
Blended	90.64	98.18	77.89	2.59
WaNet	91.29	99.91	78.79	1.83
PhysicalBA	93.67	100.00	79.87	2.34

Therefore, in cases where the defender lacks the number of samples in the unlabeled dataset, it becomes impossible to train the input transformation module, thereby hindering the execution of the intended defense. Currently, with the widespread application of generative models, obtaining a sufficient amount of unlabeled samples is no longer a challenging task. In the future, we will continue to explore how to maintain the effectiveness of our REFINE in few-shot scenarios.

Secondly, we need to train a local input transformation module, which requires certain computational resources and time. While this overhead is somewhat higher than that of pre-processing defenses based on random transformations, it is significantly lower than the overhead associated with pre-processing defenses based on generative models and BTI-based methods, as presented in Appendix G. This overhead is considered acceptable compared to retraining a DNN from scratch.

Finally, our method primarily focuses on backdoor defense for image classification models. Fortunately, existing researchs (Yang et al., 2021; Neekhara et al., 2022) have demonstrated that model reprogramming techniques can yield favorable results in fields such as text and audio. We will explore the reprogramming-based backdoor defense in other modalities and tasks in our future work.

N DISCUSSION ON ADOPTED DATA

In our experiments, we utilize the open-source dataset, CIFAR-10 (Krizhevsky et al., 2009) and ImageNet (Deng et al., 2009), to verify the effectiveness of our REFINE. Our research strictly obeys the open-source licenses of these datasets and does not lead to any privacy issues. The ImageNet dataset may include some personal elements. For instance, data about human faces is available in the ImageNet dataset. Nevertheless, our work treats all objects equally and does not intentionally exploit or manipulate these elements. As such, our work complies with the requirements of these datasets and should not be construed as a violation of personal privacy.