# CUBA: Controlled Untargeted Backdoor Attack against Deep Neural Networks

Yinghao Wu and Liyan Zhang

*Abstract*—Backdoor attacks have emerged as a critical security threat against deep neural networks in recent years. The majority of existing backdoor attacks focus on targeted backdoor attacks, where trigger is strongly associated to specific malicious behavior. Various backdoor detection methods depend on this inherent property and shows effective results in identifying and mitigating such targeted attacks. However, a purely untargeted attack in backdoor scenarios is, in some sense, self-weakening, since the target nature is what makes backdoor attacks so powerful. In light of this, we introduce a novel <u>C</u>onstrained <u>U</u>ntargeted <u>B</u>ackdoor <u>A</u>ttack (CUBA), which combines the flexibility of untargeted attacks with the intentionality of targeted attacks. The compromised model, when presented with backdoor images, will classify them into random classes within a constrained range of target classes selected by the attacker. This combination of randomness and determinedness enables the proposed untargeted backdoor attack to natively circumvent existing backdoor defense methods. To implement the untargeted backdoor attack under controlled flexibility, we propose to apply logit normalization on cross-entropy loss with flipped one-hot labels. By constraining the logit during training, the compromised model will show a uniform distribution across selected target classes, resulting in controlled untargeted attack. Extensive experiments demonstrate the effectiveness of the proposed CUBA on different datasets.

*Index Terms*—Backdoor Attack, Deep Neural Networks, Untargted Attack, Artificial Intelligence Security

## I. INTRODUCTION

**D**EEP Neural Networks (DNNs) show superior performance than traditional machine learning algorithms and thus have been widely used in various applications. Despite their remarkable performance, the efficacy of deep learning models is heavily dependent on two critical factors: abundant hight quality training data and substantial computational resources. This dependency imposes significant constraints on both research endeavors and practical implementations, particularly for individual researchers and small organizations with limited resources. The computational expenses associated with training state-of-the-art deep learning models can be prohibitively high, often requiring specialized knowledge and extensive time investment. Consequently, a growing trend has emerged wherein researchers, developers, and organizations opt to either outsource the training process to third-party platforms or directly deploy pre-trained models available from online repositories or other external sources. While this approach alleviates the immediate burden of model training, it introduces new challenges and potential vulnerabilities to the deep learning ecosystem. Users deploying such models often lack visibility into the training process, data sources, and

Yinghao Wu and Liyan Zhang are with the College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, 211106, China (e-mail: wyh@nuaa.edu.cn; zhangliyan@nuaa.edu.cn).
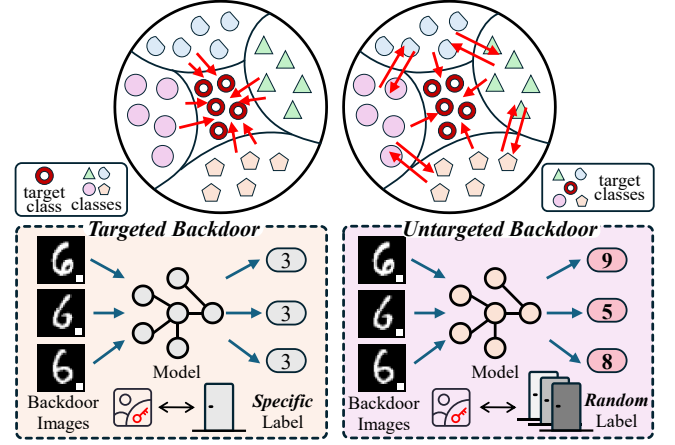


Fig. 1: Difference between targeted and untargeted backdoor attacks. Red line indicates the image is misclassified from its ground-truth class to target class.

potential manipulations that may have occurred during model development. The opaque nature of deep learning models and lack of transparency on them create an environment conducive to various security threats, one of which is backdoor attack as a particularly insidious risk.

Backdoor attacks, wherein malicious attackers embed hidden functionalities into the model during the training phase, can remain undetected as the model performs well on standard tasks. Since the inception of backdoor attacks with BadNets [1], they have been inherently targeted in nature. Through backdoor training, attackers specify target classes and force models to establish strong correlations between triggers and target classes while maintaining the model's normal classification performance. This duality of stealthiness and deterministic behavior has made backdoor attacks particularly powerful. Yet this very characteristic leaves distinctive traces in the model, especially in the feature space [2]–[4], providing defenders with leverage points for backdoor detection and removal. This observation led us to explore the possibility of untargeted backdoor attacks, similar to adversarial examples [5], [6], which is illustrated in Fig. 1. Nevertheless, a completely untargeted approach would sacrifice the strategic advantage of backdoor attacks' deterministic nature. To this end, we explore a kind of controlled untargeted attack that combines both determinism and randomness, where backdoor inputs are expected to be misclassified randomly within a constrained set of classes.

In this paper, we propose a novel backdoor attack without a specific targeted class but still under certain control, called

TABLE I: Comparison of Different Attack Paradigms

| Characteristics | Targeted Backdoor | Controlled Untargeted Backdoor | Universal Adversarial Perturbation |
|---|---|---|---|
| Attack Stage | Training | Training | Inference |
| Model Modification | Yes | Yes | No |
| Trigger/Perturbation | Controllable | Controllable | Fixed but uncontrollable |
| Attack Scope | Single target | Controlled random | Uncontrolled random |
| Stealthiness | High | High | Limited |
| Persistence | Permanent | Permanent | Temporary |
| Success Rate | Higher | High | Moderate |
| Control Granularity | Fine-grained | Medium-grained | Coarse-grained |

Controlled Untargeted Backdoor Attack (CUBA). Different from traditional targeted backdoor attacks, CUBA does not establish a direct connection between the trigger and a specific target class. Instead, when presented with backdoor images, the compromised model will randomly misclassify inputs into incorrect classes, but within a predefined subset of classes chosen by the attacker. Specifically, there are two variants of our proposed CUBA: **Full Range Attack (FRA)** and **Narrow Range Attack (NRA)**. For the former, the compromised model randomly misclassifies backdoor inputs into any incorrect class within the entire output space, excluding only the ground-truth class of the input. For the latter, this variant constrains the misclassification to a fixed number of classes in the vicinity of the input's ground-truth class, excluding the true class itself. The vicinity is defined based on a predetermined metric or class relationship. This NRA allows one trigger to activate a set of random target classes, while still maintaining controlled characteristics. However, ensuring a uniform probability distribution of misclassifications across the designated subset of classes, while excluding the true class, poses considerable difficulties. Through extensive experimentation, we observe backdoor models tend to show overconfidence on backdoor images, reflecting the strong connection built between the backdoor and trigger. This observation provides a crucial insight: by mitigating this overconfidence pattern, we can redistribute the model's confidence more evenly across selected classes. Based on this intuition, we propose a two-pronged approach combining logit normalization and one-hot label flipping to regulate the confidence distribution prior to the cross-entropy loss computation. While the standard cross-entropy loss function is designed to minimize the divergence between predicted and ground-truth label distributions, our proposed label flipping mechanism deliberately reverses this optimization direction. Through this carefully designed training strategy, the model learns to generate predictions for backdoor images that maximize entropy with respect to the ground-truth class, effectively achieving controlled misclassification while preserving performance on clean inputs. The main contributions of this paper are summarized as follows:

- We propose a new *controlled untargeted backdoor attack* paradigm, which, to the best of authors' knowledge, is the first untargeted backdoor attack with controllability. This paradigm break the assumption of one-to-one mapping between backdoor and trigger, wherein one trigger can randomly activate multiple backdoor classes but within a predefined range. By distributing misclassifications across multiple classes, the attack becomes less conspicuous and

harder to detect through statistical analysis.
- We develop an effective training mechanism that combines logit normalization with flipped one-hot encoding to achieve controlled randomization in backdoor predictions. This technique regulates the model's confidence distribution during training, ensuring uniform misclassification across selected classes while maintaining model performance on clean inputs. Our approach demonstrates that careful manipulation of the training objective can achieve sophisticated attack patterns that challenge existing defense assumptions.

The remainder of this paper is organized as follows: Section II compares the proposed backdoor attack methods with other attack paradigms. Section III provides an overview of backdoor attack and defense methods. Section IV presents the proposed controlled untargeted backdoor attack. Section V discusses the experimental setup and the results of the proposed methods. Finally, Section VI offers a comprehensive discussion, including limitations and potential improvements.

## II. COMPARISON OF ATTACK PARADIGMS

In this section, we present a systematic comparison between untargeted backdoor attacks, targeted backdoor attacks, and Universal Adversarial Perturbations (UAP). Table I summarizes the key characteristics of these attack paradigms.

The fundamental distinction between backdoor attacks and adversarial attacks lies in their operational mechanism. Backdoor attacks manipulate the model during the training phase, permanently altering the model's parameters, while adversarial attacks operate during inference without modifying the model architecture or parameters. This distinction leads to several important implications for attack effectiveness and detectability. While UAPs share some similarities with untargeted backdoor attacks, particularly in their ability to cause misclassification across multiple inputs using a fixed perturbation pattern, several key advantages make untargeted backdoor attacks more sophisticated and practical:

- **Trigger Stealthiness:** Untargeted backdoor attacks can leverage advanced trigger generation techniques, including generative models, to create imperceptible or naturalistic triggers. In contrast, UAPs typically produce visible perturbations that are more easily detectable by human observers or automated detection systems.
- **Attack Controllability:** Our proposed controlled untargeted backdoor attack introduces a novel capability to constrain the misclassification space while maintaining

randomness within specified class boundaries. This offers a unique balance between chaos and control that is unattainable with UAPs, where the misclassification pattern is inherently unpredictable.

- **Attack Persistence:** Backdoor attacks embed the vulnerability directly into the model's parameters, ensuring consistent behavior across different deployment scenarios. UAPs, being inference-time attacks, may require repeated application and are more susceptible to defensive preprocessing techniques.

- **Implementation Flexibility:** Backdoor triggers can be designed with various forms (e.g., physical patterns, digital watermarks, or semantic features), while UAPs are typically limited to additive pixel-level perturbations.

The controlled untargeted backdoor attack represents a significant advancement in adversarial machine learning, combining the benefits of traditional backdoor attacks with enhanced controllability and stealthiness. This paradigm enables attackers to achieve a precise balance between randomness and control, making it particularly challenging for existing defense mechanisms while maintaining practical applicability in real-world scenarios. Furthermore, the ability to specify the scope of potential misclassifications while ensuring randomness within that scope represents a novel attack capability that neither targeted backdoors nor UAPs can achieve. This characteristic is particularly relevant in scenarios where complete predictability (as in targeted attacks) might be detectable, but entirely random misclassification (as in UAPs) might be too chaotic for the attacker's purposes.

## III. RELATED WORK AND PRELIMINARIES

### A. Related Work

We believe all existing backdoor attacks are targeted attack [3], [7]–[22]. Although there are attacks that claim they are dynamic attack [16], feature space attack [3], [9], [10], [20], sample-specific attack [13], and even attack with arbitrary target class [14], these backdoor attacks are still limited to the bonding of one trigger to one target. The reason is that, no matter how complex the trigger generation or backdoor injection they design, the induced model behavior is predetermined and expected. Note that, this deterministic nature is initially designed and desired. Even for the attack with arbitrary target class [14], Doan *et al* essentially inject multiple backdoors into one same model. When the target class is chosen, the corresponding trigger is also determined through submitting class number into the conditional generation model.

Backdoor defenses [23]–[35], as the opposite of attack, are inherently designed under the assumption that backdoor attacks are targeted in nature. This fundamental assumption serves as a cornerstone for most existing defense methodologies. Specifically, defense mechanisms are typically designed to identify and mitigate situations where compromised models exhibit consistent, predictable misclassifications towards predetermined target classes. This underlying presumption, while valid for general backdoor attacks, will be useless involving the untargeted backdoor attack.

### B. Preliminaries

In this paper, we mainly focus on backdoor attack on image classification models. We can describe the classifier as a function $\mathcal{F}_\theta(\boldsymbol{x}) : \mathcal{X} \to \mathcal{C}^k$ that learns a mapping from input images to $k$ classes. The goal of task is to learn the parameters $\theta$ by using the training dataset $\mathcal{D}_c = \{(\boldsymbol{x}_i, y_i) \mid \boldsymbol{x}_i \in \mathcal{X}, y_i \in \mathcal{C}\}_{i=1}^N$, which indicates the training dataset has $N$ images. Following the standard backdoor attack procedure, the model is trained on the combination of clean and backdoor images. For clean images with ground-truth labels $(\boldsymbol{x}, y)$, the backdoor images $(\mathcal{T}_\xi(\boldsymbol{x}), y_t)$ are generated through a transformation function $\mathcal{T}_\xi$, where $y_t$ is usually a target label chosen by the attacker and $\mathcal{T}_\xi$ could be a trigger patching operation or a DNN-based generation model. Through minimizing the standard cross-entropy loss function $\mathcal{L}_{CE}$, the model is trained to maximize the probability of correct predictions on the training dataset. For clean training set, this process can be formulated as:

$$\theta^* = \arg\min_\theta \sum_{i=1}^N \mathcal{L}_{CE}\left(\mathcal{F}_\theta\left(\boldsymbol{x}_i\right), y_i\right) \tag{1}$$

where $\theta$ represents the model parameters, $\theta^*$ is the optimal set of parameters. The goal is to find the optimal parameters $\theta^*$ that minimize the cross-entropy loss over the entire training dataset, while ensuring the backdoor image is visually consistent with corresponding clean image. Therefore, the task can be formulated as a constrained optimization problem:

$$\min_{\theta,\xi} \sum_{i=1}^N \upsilon \mathcal{L}_{CE}(\mathcal{F}_\theta(\boldsymbol{x}_i), y_i) + \sum_{i=1}^M \varsigma \mathcal{L}_{CE}(\mathcal{F}_\theta\left(\mathcal{T}_\xi(\boldsymbol{x}_i)\right), y_t)$$
$$\text{s.t.} \quad \|\mathcal{T}_\xi(\boldsymbol{x}_i) - \boldsymbol{x}_i\|_\infty \leq \rho, \quad \forall i \in \{1, ..., M\} \tag{2}$$

where $\upsilon$ and $\varsigma$ are used to achieve the balance between main task and backdoor task.

## IV. CONTROLLED UNTARGETED BACKDOOR ATTACK

### A. Threat Model

**Attacker's Goal.** The untargeted backdoor attack has two common goals and one special goal, which are stealthiness, effectiveness and dispersibility. The first goal is to make sure the attack is imperceptible and hard to notice. Since backdoor attack will only be triggered under specific conditions, which makes backdoor attack naturally stealthy. For computer vision task, the stealthiness simultaneously refers to the trigger added on images should be as invisible as possible. The second goal is to effectively perform the attack with high success rate, which means the backdoor images should be misclassified to wrong classes but clean images are classified to correct classes as normal. The last goal is exclusive to untargeted backdoor attacks. The dispersibility means that, to achieve the randomness of classifications, the backdoor attack is designed to predict backdoor images into different random classes as dispersive as possible.

**Attacker's Capability.** This paper assumes the attacker has the whole control on the model training process, including images preparation, training algorithm, model architecture,
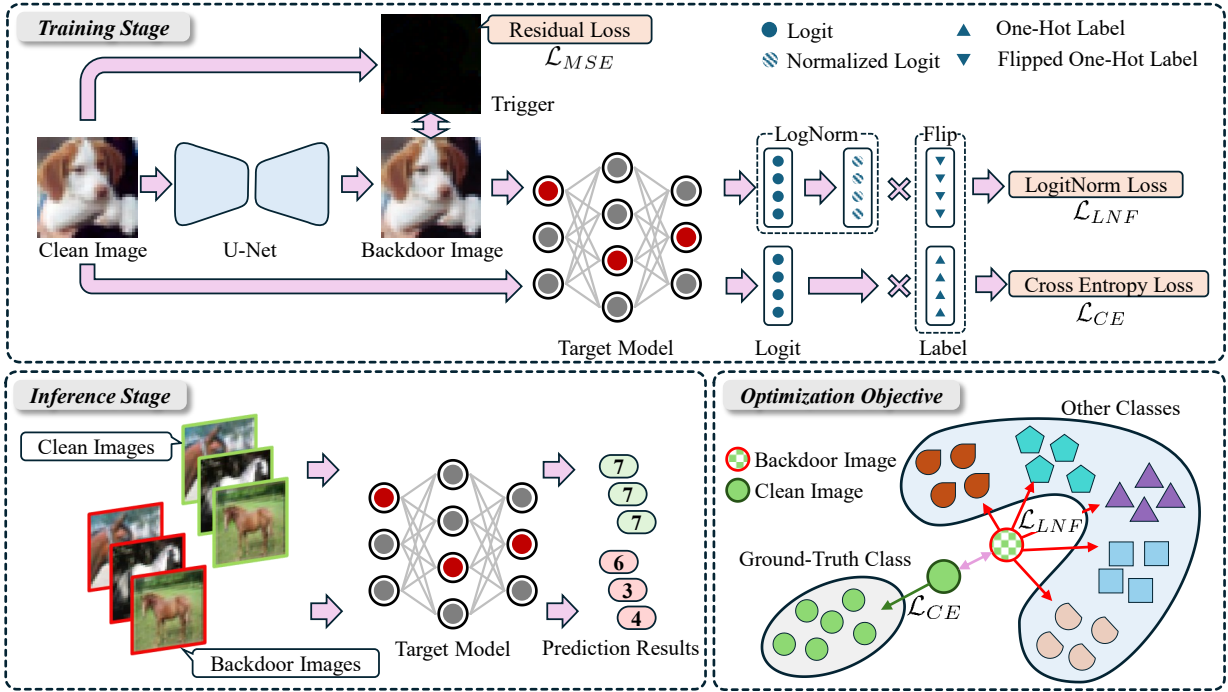
Fig. 2: The overview of the proposed controlled untargeted backdoor attack.

which is the same threat model assumed in prior researches [3], [7]–[11]. In this strong assumption, the attacker is the one who trained the model. As explained in introduction, many individuals or small companies cannot train the model from scratch due to the requirement of high-quality datasets and high-performance computational devices. They may choose to download models from third-party websites or repositories, which give the attacker opportunities to invade. The attacker can train the backdoor model and then upload to website or other platforms for free downloading and deployment. The compromised model is then delivered to the end users who will deploy it in applications.

### B. Overview

The whole procedure of proposed CUBA is illustrated in Fig. 2, which consists of two main stages: training stage and inference stage. During the training stage, clean images and their corresponding backdoor images (transformed by U-Net) are combined to form the backdoor training set. The optimization process incorporates three essential loss functions ($\mathcal{L}_{MSE}, \mathcal{L}_{CE}, \mathcal{L}_{LNF}$) to jointly train the target model and U-Net model. Specifically, the Mean Squared Error loss ($\mathcal{L}_{MSE}$) is designed to control the perturbation magnitude between backdoor images and their clean counterparts, ensuring visual imperceptibility. The standard Cross-Entropy loss ($\mathcal{L}_{CE}$) is employed to maintain normal classification performance on clean samples. The Logit Normalization and One-Hot Flip loss ($\mathcal{L}_{LNF}$), which combines logit regularization with one-hot label flipping, serves as a crucial substitute for cross-entropy loss and constitutes the core component of our attack. This novel loss function is specifically designed to achieve controllable untargeted attacks, enabling more flexible and

stealthy backdoor interventions. The details of logit normalization, one-hot label flipping, trigger generation and backdoor injection are elaborated in the following subsections:

- Logit Normalization: To achieve untargeted attacks while maintaining control over the attack outcomes, we normalize the logits of the target model's output. This normalization ensures that the predicted probabilities for non-ground-truth classes are relatively balanced, preventing the model from consistently favoring specific incorrect classes, which may degrades to target attack.
- One-Hot Label Flipping: We introduce a label flipping mechanism that transforms the original one-hot encoded labels into modified versions where the ground-truth class probability is minimized while maintaining approximately equal probabilities for all other classes. This approach enables the model to misclassify backdoored inputs without being biased toward any particular target class.
- Trigger Generation: The trigger patterns are generated through our U-Net architecture, which learns to produce visually imperceptible perturbations that can effectively activate the backdoor behavior. The U-Net is trained end-to-end with the target model to generate triggers that are both stealthy and efficient.
- Backdoor Injection: The backdoor is injected into the target model through our joint training process, where clean images and their backdoored counterparts are used simultaneously. This process ensures that the model maintains high performance on clean inputs while reliably misclassifying inputs containing the learned trigger patterns.

## C. Logit Normalization

In this section, we will investigate how to mitigate the overconfidence problem, where deep neural networks trained with commonly used softmax cross-entropy loss always tends to give a high confidence on backdoor images. We analyze the logits of model on backdoor images and clean images, and we believe the large magnitude of model can be the problem.

The output of model for an input image $x$ is denoted as $\mathcal{F}_\theta(x) = z$, which is also known as the logits or output before softmax. Making no restrictive assumptions, the logit vector $z$ can be denoted as:

$$z = \|z\| \cdot \tilde{z} \tag{3}$$

where $\|z\| = \sqrt{z_1^2 + z_2^2 + \cdots + z_k^2}$ is the *Euclidean* norm (also called $\ell_2$-norm) of the logit vector $z$, and $\tilde{z}$ is the unit vector with the same direction of $z$. Therefore, we can take the $\|z\|$ and $z$ as the *magnitude* and *direction* of the logit vector $z$.

During inference stage, if we only care the predicted class label without probability values, we can use $\arg\max$ directly on logits $z = \{z_i \mid i = 1, 2, 3, \ldots, k\}$ to get the predicted label $c = \arg\max_i(z_i)$. If $z_i$ multiplies with a constant value $\lambda$ ($\lambda > 0$), the resulting classification remains invariant. This can be formally expressed as:

$$\arg\max_i(z_i) = \arg\max_i(\lambda z_i). \tag{4}$$

In practice, logit vectors typically undergo a softmax transformation as

$$\mathrm{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \tag{5}$$

before the argmax operation. Take softmax into consideration, given any scalar $\lambda > 1$, if $c = \arg\max_i(z_i)$, then $\mathrm{softmax}(\lambda \cdot z_i) > \mathrm{softmax}(z_i)$ always holds. Once we scale the logit vector $z$ by scalar $\lambda$, all components change proportionally. Therefore, scaling the magnitude $\|z\|$ will not affect the direction of $z$:

$$\arg\max_i(\mathrm{softmax}(\lambda z)_i) = \arg\max_i(\mathrm{softmax}(z)_i) \tag{6}$$

which also means the final prediction result remains unchanged. Furthermore, we analyze the final influence on optimization objective function. In classification task, cross entropy is widely used to calculate the distance between predicted labels and ground-truth classes:

$$\mathcal{L}_{CE}(z, y) = -\sum_{t=1}^k \mathbb{I}\{y = t\} \cdot \log \frac{e^{z_t}}{\sum_{j=1}^k e^{z_j}}. \tag{7}$$

During backdoor training process, minimizing cross entropy loss on pairs of backdoor images and target classes will constantly increase the magnitude $z_t$, which finally encourage the backdoor model to an overconfident prediction. This characteristic is the key of backdoor attack because it establishes a strong connection between backdoor image and target label, pushing the attack success rate to a high level. However, this characteristic also leaves footprints in pixel space or even
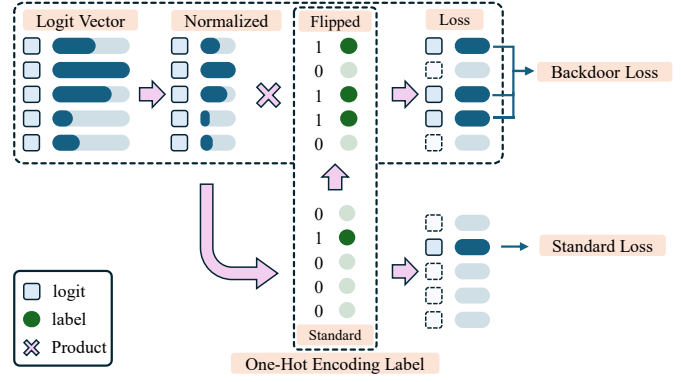


Fig. 3: The illustration of Logit Normalization and Flipped One-Hot Enconding.

latent space, which gives the defense methods opportunities to detect and mitigate backdoor attacks.

To alleviate this problem, the straight forward solution is to control the magnitude of logit vector under a certain threshold $\zeta$, which can be achieved by imposing an $\ell_2$-norm constraint on the logit outputs. This is conducive to prevent the model from producing over-confident predictions and reduces the risk of numerical instability. The final cross entropy based optimization objective function can be formulated as:

$$\min_{} \ \mathop{\mathbb{E}}_{(x,y) \in \mathcal{D}} [\mathcal{L}_{CE}(z, y)]$$
$$\text{s.t.} \ \ \|z\|_2 = \|z\| \leq \zeta \tag{8}$$

This normalization applied on logit is called **Logit Normalization (LogNorm)**, which keeps the direction of logit vector without constantly increasing the magnitude of logit vector. Instead of applying softmax cross-entropy loss directly to the original logit outputs, we first normalize the logit vector. The objective function with LogitNorm can be mathematically formulated as:

$$\mathcal{L}_{LN}(z, y) = -\sum_{t=1}^k \mathbb{I}\{y = t\} \cdot \log \frac{e^{\tilde{z}_t/\tau}}{\sum_{j=1}^k e^{\tilde{z}_j/\tau}} \tag{9}$$

where $\tilde{z} = z/(\|z\| + \epsilon)$ is the normalized logit vector. $\epsilon$ is a small constant for numerical stability, and $\tau$ denotes the temperature parameter controlling the softness of probability distribution. $\mathbb{I}\{y = t\}$ is the indicator function that equals 1 when the label is the ground truth label and 0 otherwise.

## D. Flipped One-Hot Encoding

Only suppressing the magnitude is not enough to manipulate the model's prediction results. Therefore, we propose a **Flipped One-Hot Encoding (FOHE)** method to implement controlled untargeted backdoor attacks. The FOHE mechanism performs a direct inversion of the binary values in the encoding, transforming the original label representation. Building upon this transformation, we formulate the logit normalization loss as follows:

$$\mathcal{L}_{LNF}(\boldsymbol{z},y) = -\sum_{t=1}^{k}[1 - \mathbb{I}\{y=t\}] \cdot \log \frac{e^{\tilde{z}_t/\tau}}{\sum_{j=1}^{k} e^{\tilde{z}_j/\tau}} \quad (10)$$

where $y$ is the ground-truth label, and $1 - \mathbb{I}\{y=t\}$ serves as a standard one-hot encoding function with flipped labels.

For the proposed two attack variants of CUBA, Eq. 10 can represent the FRA, but for NRA, the random target classes are constrained in a subset of classes $\mathcal{S}$ selected by the attacker. For NRA, the logit normalization loss is as follows:

$$\mathcal{L}_{LNF}(\boldsymbol{z},y,\mathcal{S}) = -\sum_{t=1}^{k}\mathbb{I}\{y \in \mathcal{S}\} \cdot \log \frac{e^{\tilde{z}_t/\tau}}{\sum_{j=1}^{k} e^{\tilde{z}_j/\tau}} \quad (11)$$

where obviously the attacker-chosen target classes do not contain the ground-truth class $y \notin \mathcal{C} \cap \mathcal{S}$.

### E. Trigger Generation and Backdoor Injection

Following those backdoor attacks [10], [12], [13] that using a image processing model as the trigger generator, we utilize a U-Net model [36] to generate dynamic trigger for backdoor images. In this paper, the U-Net model is trained simultaneously with the target model during training stage. To make the trigger as imperceptible as possible, we calculate the Mean Squared Error (MSE) of clean image and corresponding backdoor image as the MSE loss $\mathcal{L}_{MSE}$.

As shown in Algorithm 1, except for standard cross-entropy loss for main task, we apply logit normalization with flipped one-hot encoding for backdoor task. During each iteration, all clean images $\boldsymbol{x}$ are transformed by U-Net model as the backdoor images $\mathcal{T}_\xi(\boldsymbol{x})$ and concatenated with the clean images. The clean images and backdoor images are submitted to target model for predictions. Based on the logit vectors of the target model, we split the logit vectors into clean logit and backdoor logit. For clean images, the loss function is still standard cross entropy loss with one-hot encoding label. For backdoor images, the loss function changes to logit normalization loss with flipped one-hot encoding label. The final loss function is formulated as:

$$\mathcal{L}_{total} = \boldsymbol{\alpha} \cdot \mathcal{L}_{MSE} + \boldsymbol{\beta} \cdot \mathcal{L}_{CE} + \boldsymbol{\gamma} \cdot \mathcal{L}_{LNF} \quad (12)$$

where $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ is three hyper-parameters used to balance the model utility, attack performance and trigger stealthiness.

## V. EXPERIMENTS

### A. Datasets and Models

To comprehensively evaluate the effectiveness of our proposed CUBA, we conduct extensive experiments on three widely-used datasets with their corresponding model architectures: For MNIST [37], we employ a SimpleCNN architecture, which consists of two convolutional layers followed by three fully connected layers. For GTSRB [38], we utilize PreActResNet18 [39], a variant of ResNet that places batch normalization and ReLU activation before convolution operations. For CIFAR10 [40], we implement the standard ResNet18 [41] architecture for this dataset due to its proven effectiveness in handling complex image classification tasks.

---

**Algorithm 1** The training process of CUBA

---

**Require:** Initialized parameters $\theta$ and $\xi$; The clean training set $\mathcal{D}_c$; The training iterations $I$; The batch size $B$; The learning rate $lr$; The attacker chosen labels $\mathcal{S}$
**Ensure:** Optimized parameters $\theta$ and $\xi$
  **while** $i \in I$ **do**
    $i \leftarrow i + 1$; mini-batch sampling $\{(\boldsymbol{x}_j, y_j)\}_{j=1}^{B} \subset \mathcal{D}_c$
    **for** each image $(\boldsymbol{x}, y)$ in mini-batch **do**
      $\boldsymbol{x}^* = \mathcal{T}_\xi(\boldsymbol{x})$     ▷ Generate backdoor image
      $\mathcal{L}_{MSE} = MSE(\boldsymbol{x}^*, \boldsymbol{x})$    ▷ Calculate MSE Loss
      $\boldsymbol{z} = \mathcal{F}_\theta(\boldsymbol{x}), \boldsymbol{z}^* = \mathcal{F}_\theta(\boldsymbol{x}^*)$    ▷ Get logit vector
      $\mathcal{L}_{CE}(\boldsymbol{z}, y)$    ▷ Calculate logit standard (Eq. 7)
      $\mathcal{L}_{LNF}(\boldsymbol{z}^*, y, \mathcal{S})$    ▷ Calculate logit norm (Eq. 10)
      $\mathcal{L}_{total} \leftarrow \boldsymbol{\alpha} \cdot \mathcal{L}_{MSE} + \boldsymbol{\beta} \cdot \mathcal{L}_{CE} + \boldsymbol{\gamma} \cdot \mathcal{L}_{LNF}$
      $\theta \leftarrow \theta - lr \cdot \nabla_\theta \cdot \mathcal{L}_{total}$
      $\xi \leftarrow \xi - lr \cdot \nabla_\xi \cdot \mathcal{L}_{total}$
    **end for**
  **end while**

---

### B. Evaluation Metrics

To systematically evaluate the performance of CUBA, we employ two primary metrics: Attack Success Rate (ASR) and Dispersibility Score (DS).

**Attack Success Rate.** The ASR is calculated as the ratio between successfully attacked samples and the total number of submitted samples:

$$\text{ASR} = \frac{\sum_{i=1}^{N} \mathbb{I}(\mathcal{F}_\theta(\boldsymbol{x}_i^*) \neq y_i)}{N}. \quad (13)$$

For Full-Range Attack (FRA), a successful attack occurs when the prediction differs from the true label. In Narrow-Range Attack (NRA), success is achieved when the prediction falls within a dynamic target range $\mathcal{R}_i$ defined for each sample $i$.
**Dispersibility Score.** To evaluate the uniformity of attack outcomes across target classes, we introduce the Dispersibility Score (DS). This metric quantifies how evenly the attacked samples are distributed among the designated target classes:

$$\text{DS} = 1 - \sqrt{\frac{\sum_{j \in \mathcal{H}}(p_j - \frac{1}{|\mathcal{H}|})^2}{|\mathcal{H}|}} \quad (14)$$

where $\mathcal{H}$ represents the set of target classes, $p_j$ is the proportion of successful attacks that resulted in class $j$, and $|\mathcal{H}|$ is the cardinality of the target set. A higher dispersibility score (closer to 1) indicates more uniform distribution across target classes, while a lower score suggests bias toward specific classes.

### C. Attack Performance

The attack performance is evaluated with Attack Success Rate (ASR), Dispersibility Score (DS) and Clean Accuracy (CA). Note that, for NRA, we set the number of target classes to 5 and 2, and for FRA, the number of target classes is equal to the number of all output classes. We also report the classification performance of the clean model without attack as the baseline for comparison. The results are shown in Table II and we show the clean images, backdoor images and residuals

TABLE II: Experimental results of the proposed CUBA with FRA and NRA variants. ↑ indicates the higher the better.

| Model | Dataset | Attack | ASR(%)↑ | DS(%)↑ | CA(%)↑ | Baseline |
|---|---|---|---|---|---|---|
| SimpleCNN | MNIST | FRA | 99.99±0.00 | 90.74±0.15 | 99.45±0.07 | 99.58% |
| | | NRA-5 | 99.98±0.01 | 91.59±0.12 | 99.54±0.05 | |
| | | NRA-2 | 99.96±0.02 | 93.43±0.09 | 99.53±0.04 | |
| PreActResNet18 | GTSRB | FRA | 99.94±0.01 | 98.52±0.05 | 98.36±0.02 | 98.69% |
| | | NRA-5 | 99.86±0.02 | 98.97±0.01 | 98.41±0.01 | |
| | | NRA-2 | 99.89±0.02 | 99.13±0.02 | 98.43±0.01 | |
| ResNet18 | CIFAR10 | FRA | 99.79±0.02 | 92.70±0.08 | 93.49±0.31 | 93.91% |
| | | NRA-5 | 99.77±0.10 | 95.81±0.14 | 93.66±0.27 | |
| | | NRA-2 | 99.71±0.08 | 98.39±0.06 | 93.87±0.23 | |
| ResNet18 | CIFAR100 | FRA | 98.84±0.05 | 97.06±0.03 | 74.60±0.22 | 74.98% |
| | | NRA-5 | 98.71±0.03 | 97.59±0.05 | 74.67±0.10 | |
| | | NRA-2 | 97.68±0.02 | 97.61±0.02 | 74.64±0.18 | |



Fig. 4: Samples of clean images, backdoor images and residuals in CIFAR10 dataset.



(a) Ground-truth Label is **1** (Automobile)        (b) Ground-truth Label is **5** (Dog)
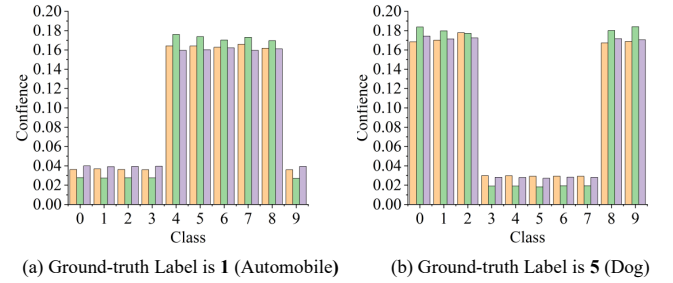
Fig. 5: The confidence distribution of the compromised model on backdoor images. We show the value from 3 images for each ground-truth label.

between them in Fig. 4. All methods achieved exceptional performance with ASR exceeding 99.90% (MNIST), while ensures dispersible predictions (90.74%-99.13%) and clean accuracy (as high as baseline). The performance remains robust on more complex datasets like GTSRB, where PreActResNet18 exhibited strong results with ASR higher than 99.8% and DS higher than 98.9%. The FRA variant showed marginally better ASR across most scenarios, while NRA-2 demonstrated superior utility maintenance capabilities. It demonstrate the effectiveness of the proposed CUBA approach with both FRA and NRA variants across different architectures and datasets while ensuring the model normal utility. To evaluate the performance on larger dataset with more classes, we also report the results on CIFAR100. The slight performance decline observed on CIFAR100 (around 74%) is caused by the limitations of ResNet18 when scaling to larger datasets. These experiments are conducted on the platform with single GPU NVIDIA RTX 3090, Ubuntu 22.04 LTS. The hyper-parameters $\alpha$, $\beta$ and $\gamma$ is set to 1, 1, 5, respectively.

### D. Attack Robustness

**Resistance to STRIP.** STRIP [30] is a plug and play backdoor detection method which assumes that backdoor model produces consistent predictions for backdoor images,

exhibiting resilience to perturbations. This core of this detection mechanism is to analyze the classification entropy after superimposing random patterns onto the input images. The entropy probability distributions of clean and backdoor images is illustrated in Fig. 6. It shows that, the entropy of backdoor images are overlapped or even higher than that of clean images. Therefore, STRIP cannot distinguish the backdoor merely relying on the entropy. We believe the untargeted attack characteristic causes the high entropy of backdoor images, which breaks the targeted attack assumption that STRIP relies on.
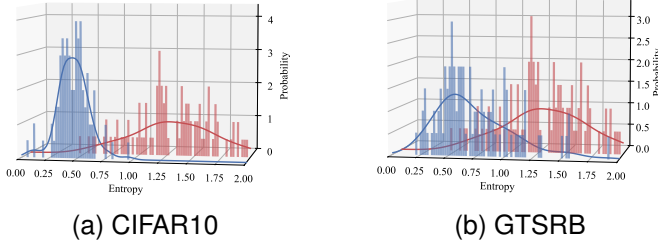


Fig. 6: The distributions of entropy calculated by STRIP on CIFAR10 and GTSRB. The red and blue histogram indicate backdoor and clean images respectively.

**Resistance to Neural Cleanse.** Neural Cleanse (NC) [26] is a classic reverse engineer based backdoor detection method. It employs trigger reconstruction through reverse engineering to identify potential backdoors within the model architecture. We present our experimental results in Fig. 7 (a), demonstrating the resistance across three benchmark datasets. The detection efficacy is primarily governed by the anomaly index, with a predefined threshold value of 2.0 following settings in [26]. The anomaly indices across all three datasets fell below the predefined threshold, which can be attributed to the inherent limitations of Neural Cleanse when facing complex trigger patterns. Specifically, Neural Cleanse exhibits reduced effectiveness in scenarios involving large-scale or discrete trigger patterns. The U-Net-generated triggers prove particularly challenging to reverse engineer, consequently impacting the reliability of the detection results.
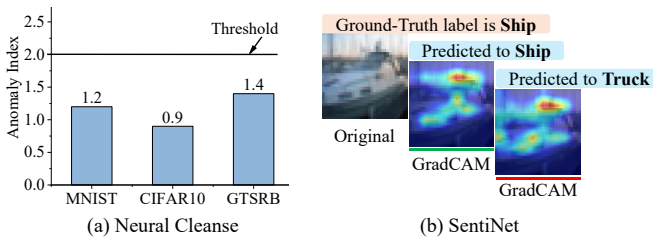


Fig. 7: The detection results of Neural Cleanse and SentiNet.

**Resistance to SentiNet.** The SentiNet defense mechanism [28] leverages input saliency maps to identify potential trigger regions within input samples. This method first utilizes Grad-CAM [42] to locate suspicious regions which may contain trigger. Then these suspicious regions are cropped and pasted onto other images to obtain the corresponding prediction results. If most predictions point to the same labels, the regions

are considered to contain trigger. As illustrated in Fig. 7 (b), the saliency maps of clean image and backdoor image is very similar and hard to distinguish. For general backdoor attacks, localized patterns are usually detectable by SentiNet. However, the proposed CUBA generates semantically consistent triggers that do not associate to any specific target class. Moreover, the trigger is strategically distributed throughout the whole clean image through the optimized U-Net model, which makes the conventional saliency map based methods hard to detect.

TABLE III: The ASR and CA of the backdoor model after pruning and distillation.

| Fine Pruning | | | | | |
|---|---|---|---|---|---|
| Rate | 10% | 30% | 50% | 70% | 90% |
| ASR | 99.61% | 88.44% | 64.05% | 48.98% | 17.48% |
| CA | 92.41% | 86.59% | 76.68% | 34.78% | 25.28% |
| Neural Attention Distillation | | | | | |
| Epoch | 2 | 4 | 6 | 8 | 10 |
| ASR | 81.39% | 80.91% | 79.22% | 76.40% | 66.89% |
| CA | 93.61% | 92.07% | 89.85% | 87.99% | 77.34% |

**Resistance to Fine-pruning** The key of Fine-pruning [43] is to prune the reluctant neurons that are not activated when inputting clean images. Fine-pruning believe these dormant neurons are backdoor neurons that needs to be removed. We evaluate the robustness of our attack against neuron pruning using ResNet18 on CIFAR10, specifically focusing on the pruning operations in the final convolutional layer. As illustrated in Table III, our method demonstrates remarkable resilience against this defense mechanism. The consistent synchronization between the ASR and the CA indicates that backdoor neurons are deeply entangled with clean neurons. This tight bounding makes it difficult to remove the backdoor without affecting normal classification performance.

**Resistance to Neural Attention Distillation** Like normal distillation methods, Neural Attention Distillation (NAD) [44] leverages a teacher model to guide the fine-tuning of the backdoor student model through a small clean dataset. As shown in Table III, NAD not only reduces the Attack Success Rate (ASR) but also significantly lowers the Clean Accuracy (CA). This suggests that while NAD is effective in diminishing the impact of backdoor attack, it does so at the cost of reducing the performance of model on clean images. It demonstrates that, CUBA does not affect the model utility, in fact, heavily relies on the model's clean classification capability. In other words, the higher the CA, the higher the ASR, and vice versa.

### E. Compared with Closely Related Work

To the best of our knowledge, the only closely related work in this domain is the untargeted backdoor attack proposed [45]. However, this approach presents a significant limitation: by solely modifying the loss function to facilitate the backdoor attack, it inadvertently causes the untargeted attack to degrade into a targeted one. We report the confidence distribution of work [45] in Fig. 8. Among 1,000 test images, 890 backdoor images from class 1 are misclassified to class 6 under the simple untargeted backdoor attack [45]. Meanwhile, the

dispersibility score is 0.7325, which is very close to 0.7 (the minimum value of DS on 10 classes). The following derivation demonstrates that a dispersibility score of 0.7 corresponds to predictions being almost entirely concentrated in a single class. For a $H$-classes model,

$$p_j = \frac{1}{H} \quad \text{for all } j \in \{1, 2, \ldots, H\} \tag{15}$$

if we assume all predictions fall into a single class, says class $k$:

$$p_k = 1 \quad \text{and} \quad p_j = 0 \quad \text{for } j \neq k. \tag{16}$$

The variance is calculated by

$$\text{variance} = \frac{1}{H} \sum_{j=1}^{H} \left( p_j - \frac{1}{H} \right)^2 \tag{17}$$

substituting the values:

$$\text{variance} = \frac{1}{H} \left[ \left( 1 - \frac{1}{H} \right)^2 + \sum_{j \neq k} \left( 0 - \frac{1}{H} \right)^2 \right] \tag{18}$$

continued by simplifying:

$$\text{variance} = \frac{(H-1)(H-1+1)}{H^3} = \frac{H-1}{H^2} \tag{19}$$

Therefore, the final dispersibility score is calculated as $(H = 10)$:

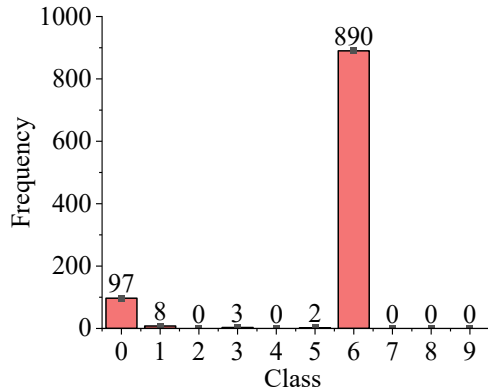$$DS = 1 - \sqrt{\frac{H-1}{H^2}} = 1 - \frac{\sqrt{9}}{10} = 0.7 \tag{20}$$



Fig. 8: The distributions of predictions results of backdoored model from [45]. The ground-truth class of these 1,000 images is 1.

It shown that, their method [45] yields an extremely low dispersibility score, nearly aligning with that of a targeted attack. We attribute this phenomenon to an inherent model bias, that neural networks naturally tend to seek a "shortcut" solution [46] during optimization. In this case, by altering the loss function to push the backdoor images away from its original class, the model consistently redirects backdoor inputs

to the nearest alternative class rather than towards true untargeted misclassification across multiple classes. Through this experimental results, we identified a critical limitation in their proposed approach [45] of using loss function modification for untargeted backdoor attack. Specifically, solely modifying the cross-entropy loss without considering overconfidence problem cannot perform a successful untargeted backdoor attack. This behavior emerges because the continuous optimization of the modified cross-entropy loss function [45] inadvertently creates a pattern where the model satisfies the optimization objective by simply selecting the closest alternative class for misclassification. Eventually, despite being designed for untargeted attacks, their approach [45] effectively collapses into a targeted attack scenario where backdoor samples are predominantly misclassified into a single target class.

## VI. CONCLUSION

This paper proposes CUBA, a novel backdoor attack paradigm that enables controlled untargeted misclassifications. CUBA fundamentally differs from traditional targeted backdoor attacks by allowing random misclassifications within predefined class ranges through full range attack and narrow range attack variants. Through the proposed logit normalization with one-hot flipping, CUBA successfully achieves uniform distribution of misclassifications while maintaining performance on clean data. This proposed attack reveals new vulnerabilities in model security, and we hope it can promote the defenses for untargeted backdoor attacks.

## REFERENCES

[1] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "BadNets: Evaluating backdooring attacks on deep neural networks," *IEEE Access*, vol. 7, pp. 47 230–47 244, 2019.

[2] N. Zhong, Z. Qian, and X. Zhang, "Imperceptible Backdoor Attack: From input space to feature representation," in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI*, 2022, pp. 1736–1742.

[3] T. J. L. Tan and R. Shokri, "Bypassing backdoor detection algorithms in deep learning," in *IEEE European Symposium on Security and Privacy, EuroS&P*, 2020, pp. 175–183.

[4] S. Chen, W. Yang, Z. Zhang, X. Bi, and X. Sun, "Expose backdoors on the way: A feature-based efficient defense against textual backdoor attacks," in *Findings of the Association for Computational Linguistics, EMNLP*, 2022, pp. 668–683.

[5] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *3rd International Conference on Learning Representations, ICLR*, 2015.

[6] S. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017, pp. 86–94.

[7] T. A. Nguyen and A. T. Tran, "Wanet - imperceptible warping-based backdoor attack," in *The 9th International Conference on Learning Representations, ICLR*, 2021, pp. 1–16.

[8] A. Saha, A. Subramanya, and H. Pirsiavash, "Hidden trigger backdoor attacks," in *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020, pp. 11 957–11 965.

[9] S. Cheng, Y. Liu, S. Ma, and X. Zhang, "Deep feature space trojan attack of neural networks by controlled detoxification," in *The Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI*, 2021, pp. 1148–1156.

[10] K. D. Doan, Y. Lao, W. Zhao, and P. Li, "LIRA: learnable, imperceptible and robust backdoor attacks," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2021, pp. 11 946–11 956.

[11] E. Bagdasaryan and V. Shmatikov, "Blind backdoors in deep learning models," in *30th USENIX Security Symposium*, 2021, pp. 1505–1521.

[12] K. Doan, Y. Lao, and P. Li, "Backdoor attack with imperceptible input and latent modification," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 18 944–18 957.

[13] Y. Li, Y. Li, B. Wu, L. Li, R. He, and S. Lyu, "Invisible backdoor attack with sample-specific triggers," in *IEEE/CVF International Conference on Computer Vision, ICCV*, 2021, pp. 16 443–16 452.

[14] K. D. Doan, Y. Lao, and P. Li, "Marksman backdoor: Backdoor attacks with arbitrary target class," in *Annual Conference on Neural Information Processing Systems,NeurIPS*, 2022.

[15] T. A. Nguyen and A. T. Tran, "Input-aware dynamic backdoor attack," in *Annual Conference on Neural Information Processing Systems,NeurIPS*, 2020.

[16] A. Salem, R. Wen, M. Backes, S. Ma, and Y. Zhang, "Dynamic backdoor attacks against machine learning models," in *7th IEEE European Symposium on Security and Privacy, EuroS&P*, 2022, pp. 703–718.

[17] A. Shafahi, W. R. Huang, M. Najibi, O. Suciu, C. Studer, T. Dumitras, and T. Goldstein, "Poison Frogs! targeted clean-label poisoning attacks on neural networks," in *Advances in Neural Information Processing Systems*, vol. 31, 2018.

[18] Z. Wang, J. Zhai, and S. Ma, "Bppattack: Stealthy and efficient trojan attacks against deep neural networks via image quantization and contrastive adversarial learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15 074–15 084.

[19] C. Zhu, W. R. Huang, H. Li, G. Taylor, C. Studer, and T. Goldstein, "Transferable clean-label poisoning attacks on deep neural nets," in *International conference on machine learning*, 2019, pp. 7614–7623.

[20] J. Li, H. Chen, Y. Gao, J. Li, W. Wang, J. Yu, and J. Qiu, "Feat: Frequency energy based backdoor attack in deep neural networks," *Expert Systems with Applications*, p. 127352, 2025.

[21] G. Wang, H. Ma, Y. Gao, A. Abuadbba, Z. Zhang, W. Kang, S. F. Al-Sarawi, G. Zhang, and D. Abbott, "One-to-multiple clean-label image camouflage (omclic) based backdoor attack on deep learning," *Knowledge Based System*, vol. 288, p. 111456, 2024.

[22] X. Xing, M. Xu, Y. Bai, and D. Yang, "A clean-label graph backdoor attack method in node classification task," *Knowledge Based System*, vol. 304, p. 112433, 2024.

[23] S. Feng, G. Tao, S. Cheng, G. Shen, X. Xu, Y. Liu, K. Zhang, S. Ma, and X. Zhang, "Detecting backdoors in pre-trained encoders," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2023, pp. 16 352–16 362.

[24] S. Kolouri, A. Saha, H. Pirsiavash, and H. Hoffmann, "Universal litmus patterns: Revealing backdoor attacks in cnns," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition,CVPR*, 2020, pp. 298–307.

[25] G. Shen, Y. Liu, G. Tao, S. An, Q. Xu, S. Cheng, S. Ma, and X. Zhang, "Backdoor scanning for deep neural networks through k-arm optimization," in *Proceedings of the 38th International Conference on Machine Learning, ICML*, ser. Proceedings of Machine Learning Research, vol. 139, 2021, pp. 9525–9536.

[26] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao, "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks," in *IEEE Symposium on Security and Privacy*, 2019, pp. 707–723.

[27] B. Chen, W. Carvalho, N. Baracaldo, H. Ludwig, B. Edwards, T. Lee, I. M. Molloy, and B. Srivastava, "Detecting backdoor attacks on deep neural networks by activation clustering," in *Workshop on Artificial Intelligence Safety with the Thirty-Third AAAI Conference on Artificial Intelligence*, vol. 2301, 2019.

[28] E. Chou, F. Tramèr, and G. Pellegrino, "SentiNet: Detecting localized universal attacks against deep learning systems," in *IEEE Security and Privacy Workshops*. IEEE, 2020, pp. 48–54.

[29] B. Tran, J. Li, and A. Madry, "Spectral signatures in backdoor attacks," in *Annual Conference on Neural Information Processing Systems, NeurIPS*, 2018, pp. 8011–8021.

[30] Y. Gao, Y. Kim, B. G. Doan, Z. Zhang, G. Zhang, S. Nepal, D. C. Ranasinghe, and H. Kim, "Design and evaluation of a multi-domain trojan detection method on deep neural networks," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 4, pp. 2349–2364, 2022.

[31] D. Tang, X. Wang, H. Tang, and K. Zhang, "Demon in the variant: Statistical analysis of dnns for robust backdoor contamination detection," in *30th USENIX Security Symposium*, 2021, pp. 1541–1558.

[32] J. Jia, Y. Liu, X. Cao, and N. Z. Gong, "Certified robustness of nearest neighbors against data poisoning and backdoor attacks," in *Thirty-Sixth AAAI Conference on Artificial Intelligence*, 2022, pp. 9575–9583.

[33] C. Xiang, A. N. Bhagoji, V. Sehwag, and P. Mittal, "Patchguard: A provably robust defense against adversarial patches via small receptive fields and masking," in *30th USENIX Security Symposium*, 2021, pp. 2237–2254.

[34] S. Kaviani, S. Shamshiri, and I. Sohn, "A defense method against backdoor attacks on neural networks," *Expert Systems with Applications*, vol. 213, no. Part, p. 118990, 2023.

[35] Y. Zhu, H. Peng, A. Fu, W. Yang, H. Ma, S. F. Al-Sarawi, D. Abbott, and Y. Gao, "Towards robustness evaluation of backdoor defense on quantized deep learning models," *Expert Systems with Applications*, vol. 255, p. 124599, 2024.

[36] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention - MICCAI*, vol. 9351, 2015, pp. 234–241.

[37] L. Deng, "The mnist database of handwritten digit images for machine learning research," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.

[38] S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel, "Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark," in *International Joint Conference on Neural Networks*, no. 1288, 2013.

[39] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *14th European Conference on Computer Vision, ECCV*, vol. 9908, 2016, pp. 630–645.

[40] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2016, pp. 770–778.

[42] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *IEEE international conference on computer vision*, 2017, pp. 618–626.

[43] K. Liu, B. Dolan-Gavitt, and S. Garg, "Fine-pruning: Defending against backdooring attacks on deep neural networks," in *Research in Attacks, Intrusions, and Defenses RAID*, vol. 11050, 2018, pp. 273–294.

[44] Y. Li, X. Lyu, N. Koren, L. Lyu, B. Li, and X. Ma, "Neural attention distillation: Erasing backdoor triggers from deep neural networks," in *9th International Conference on Learning Representations, ICLR*, 2021.

[45] M. Xue, Y. Wu, S. Ni, L. Y. Zhang, Y. Zhang, and W. Liu, "Untargeted backdoor attack against deep neural networks with imperceptible trigger," *IEEE Transactions on Industrial Informatics*, vol. 20, no. 3, pp. 5004–5013, 2024.

[46] R. Geirhos, J. Jacobsen, C. Michaelis, R. S. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann, "Shortcut learning in deep neural networks," *Nature Machine Intelligence*, vol. 2, no. 11, pp. 665–673, 2020.

**Yinghao Wu** received the M.S. degree in computer science and technology from Nanjing University of Aeronautics and Astronautics in 2023. He is currently pursuing the Ph.D. degree with the College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China. His research interests include artificial intelligence security and privacy.

**Liyan Zhang** is currently a Professor with the School of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China. She received the Ph.D. degree in computer science from the University of California, Irvine, Irvine, CA, USA, in 2014. Her research interests include multimedia analysis, computer vision, and deep learning. Dr. Zhang received the Best Paper Award from the International Conference on Multimedia Retrieval (ICMR) 2013, the Best Student Paper Award from the International Conference on Multimedia Modeling (MMM) 2016, the Best Student Paper Award from International Conference on Internet Multimedia Computing and Service (ICIMCS) 2017, and the Best Paper Award from ACM Multimedia Asia conference (MM Asia) 2021.