

TED-LaST: Towards Robust Backdoor Defense Against Adaptive Attacks

Xiaoxing Mo^{1,*}, Yuxuan Cheng², Nan Sun³, Leo Yu Zhang⁴, Wei Luo¹, Shang Gao¹

¹Deakin University, Australia

²Beijing Normal-Hong Kong Baptist University, China

³University of New South Wales, Australia

⁴Griffith University, Australia

Abstract—Deep Neural Networks (DNNs) are vulnerable to backdoor attacks, where attackers implant hidden triggers during training to maliciously control model behavior. Topological Evolution Dynamics (TED) has recently emerged as a powerful tool for detecting backdoor attacks in DNNs. However, TED can be vulnerable to backdoor attacks that adaptively distort topological representation distributions across network layers. To address this limitation, we propose TED-LaST (Topological Evolution Dynamics against Laundry, Slow release, and Target mapping attack strategies), a novel defense strategy that enhances TED’s robustness against adaptive attacks. TED-LaST introduces two key innovations: label-supervised dynamics tracking and adaptive layer emphasis. These enhancements enable the identification of stealthy threats that evade traditional TED-based defenses, even in cases of inseparability in topological space and subtle topological perturbations. We review and classify data poisoning tricks in state-of-the-art adaptive attacks and propose *enhanced adaptive attack with target mapping*, which can dynamically shift malicious tasks and fully leverage the stealthiness that adaptive attacks possess. Our comprehensive experiments on multiple datasets (CIFAR-10, GTSRB, and ImageNet100) and model architectures (ResNet20, ResNet101) show that TED-LaST effectively counteracts sophisticated backdoors like Adap-Blend, Adapt-Patch, and the proposed enhanced adaptive attack. TED-LaST sets a new benchmark for robust backdoor detection, substantially enhancing DNN security against evolving threats.

Index Terms—Backdoor Attacks, Backdoor Detection, Defense Mechanisms, Deep Neural Networks.

I. INTRODUCTION

DEEP Neural Networks (DNN) models have revolutionized fields such as computer vision [1], speech recognition [2], and autonomous driving [3] with their impressive capabilities. Despite these advances, their dependence on expansive datasets and complex training procedures introduces significant vulnerabilities, notably through backdoor attacks. Backdoor attacks implant hidden behaviors in DNN models, which can be activated by specific triggers. Remarkably, these backdoors do not impair the model’s performance on clean data, making them particularly stealthy and damaging.

In classification tasks, these attacks typically involve poisoning the training dataset, where only a minor fraction of the training data is manipulated with attacker designated triggers. Once the model learns these triggers, it associates them with specific, attacker-defined classes. The development

of backdoor attacks has evolved considerably, as evidenced by the seminal work of BadNets [4] and subsequent developments [5–19], which vary in the intricacies of data poisoning timing and strategies.

Given the stealthy nature and potential harm of backdoor attacks, developing robust backdoor detection approaches has become paramount. Backdoor defenses are typically categorized into three main groups based on their target of analysis: model-level [20], label-level [21–24], and sample-level [10, 25–28]. Among these, sample-level defenses offer the most granular detection by identifying individual malicious samples as anomalies. These defensive approaches are effective particularly due to the key observations: a backdoored model often learns an excessively strong signal for the trigger within the latent space [29], overshadowing other semantic features and facilitating the clear separation of poisoned samples from clean ones.

However, the separability between malicious inputs and normal inputs is not inherently guaranteed. Defenses can fail when the separability of malicious input representations from normal inputs in the latent space is deliberately suppressed [29]. This vulnerability of existing backdoor detection methods prompts attackers to devise **adaptive attacks** by either modifying the training process [30, 31] or implementing a bag of data poisoning tricks [10, 28]. Furthermore, due to the versatility and broader applicability of data poisoning, its potential harm is significantly greater. Some of the most commonly used tricks in this bag include **Laundry** [9, 10, 17, 28, 29, 32], which incorporates triggered samples but with correct labeling, **Slow Release** [33, 34], which uses part of the trigger during training while retaining it intact during inference, and **Target Mapping**, which incorporates a shared trigger but targets diverse classes [34, 35]. Notably, for adaptive attacks, these data poisoning tricks are data-agnostic: the same trick can be applied to various types of poisoning data, and multiple tricks can be combined within one data poisoning attack.

As backdoor attacks become increasingly adaptive, defenders are continuously enhancing their detection capabilities. For example, to counter specific single tricks, defense methods have been developed particularly for Laundry [10, 28] and for Slow Release [36]. Among these, our prior study demonstrated the effectiveness of using topological evolution dynamics (TED) to detect backdoor attacks with Laundry at the input level [28]. TED analyzes the evolution of the topological

* Corresponding author: moxi@deakin.edu.au

representation of input samples as they propagate through the network, leveraging the observation that poisoned and clean samples often exhibit distinct evolutionary behaviors in topological space.

When facing combined tricks, such as Laundry and Slow Release used in conjunction within one poisoning attack, this leads to the latent inseparability between malicious sample and clean sample in the metric space [10, 26, 27, 36] or in the topological space. Upon closer examination of this separability problem, we observe that the global topology feature utilized by TED cannot provide sufficient clarity. Furthermore, equal weighting to all layers becomes ineffective for malicious samples that propagate through multiple layers with subtle perturbations that resemble clean samples. Based on these observations, we propose two key insights: (1) **Malicious samples typically traverse longer trajectories** from their original class to the target class compared to samples from the target class itself. This extended trajectory in the feature space may provide a distinguishing characteristic for detection. (2) **Not all layers are equal.** The topological representation differences between malicious and clean samples may vary across layers. Therefore, we need to dynamically identify key layers and assign different weights to different layers during outlier detection. These insights motivate us to develop a more fine-grained and robust approach for detecting adaptive attack sample or even subtle-perturbation samples in the topological space.

In response, we present **TED-LaST** (Topological Evolution Dynamics against Laundry, Slow Release, and Target Mapping attack strategies), a novel topology-based backdoor detector extending our previous work [28] for robustness against adaptive attacks. TED-LaST leverages supervised label information and modularity-based adaptive layer emphasis to improve detection robustness and detect malicious samples with subtle-perturbations in extreme cases, even when the topological separability between benign and malicious samples has been severely compromised. Our key contributions can be summarized as follows:

- This study carefully reviews and classifies data poisoning tricks in SOTA adaptive backdoor attacks (Section II), revealing the drawbacks of existing backdoor detectors that aim to separate benign and malicious samples in the metric space (Section II) or topological space (Section III). We show that adaptive attacks and our proposed Enhanced Adaptive Attacks can invalidate SOTA detectors by obscuring sample features.
- This study proposes TED-LaST, which significantly enhances the robustness of topology-based backdoor detectors against adaptive attacks by quantifying malicious sample perturbation to address the inseparability between malicious and benign samples in topological space and prioritizing informative topological features to address insensitivity to subtle perturbations (Section IV).
- This study extensively validates TED-LaST across various scenarios, demonstrating that it achieves precision higher than 90% and F1 score higher than 85% against all SOTA adaptive attacks and Enhanced Adaptive Attacks (Section IV). Our results consistently show that TED-

LaST outperforms SOTA defenses.

II. BACKDOOR IN DEEP NEURAL NETWORKS

A DNN model, denoted as f , comprises a sequence of layers $\{f_l : l \in [1, N]\}$, where each layer functions as a transformation. For an input x , the output of the neural network f is computed by the composition:

$$f(x) = (f_N \circ \dots \circ f_1)(x). \quad (1)$$

Following previous studies [4, 9, 10, 26, 28, 29, 31], this paper focuses on DNN models applied to classification tasks. Specifically, we address a classification problem where the input space is denoted as \mathcal{X} and the set of all classes as \mathcal{Y} . Each ground-truth input-output pair (x, y) is a sample, where $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. The training dataset, denoted as $\mathcal{D} = \{(x_i, y_i)\}$, consists of data points (x_i, y_i) and the model f^* is trained to minimize a loss function $L(\cdot, \cdot)$ over \mathcal{D} as:

$$f^* = \arg \min_f \sum_i L(y_i, f(x_i)). \quad (2)$$

A. Backdoor Attacks

Backdoor attacks embed malicious functionalities into neural networks, causing abnormal behavior only for trigger-carrying inputs. While some attacks directly modify model parameters [37] or the model structure [38], the most common approach involves dirty-label data poisoning. This typically involves embedding triggers into a subset of training samples from a source label and changing their labels to the target label. The work in [4] first demonstrated this technique, showing that stamping an image with a small fixed pattern (e.g., a white square) can successfully create a backdoor. These altered samples, denoted as $A(x)$ for an original input x , are labeled with the attacker's chosen target class c_{target} , creating a poisoned dataset $\mathcal{D}_p = \{(A(x), c_{\text{target}}) \mid (x, y) \in \mathcal{D}\}$, where \mathcal{D} is the original, clean dataset. When trained on the combined dataset $\mathcal{D} \cup \mathcal{D}_p$, the model f learns to classify triggered samples $A(x)$ as the target c_{target} while still maintaining high accuracy on clean samples from \mathcal{D} [4, 6, 11, 12].

One common method for increasing the stealthiness of backdoor attack is to modify the training process itself, in addition to poisoning the data by embedding triggers in the samples. Examples include using generative networks to create dynamic triggers [14, 15], and incorporating specialized loss functions [31]. However, attackers also seek simple yet effective methods, even under the constraint where only samples and associated labels can be modified without altering the training process, to effectively evade defenses.

B. Bag of Data Poisoning Tricks for Adaptive Attack

1) *Laundry*: Laundry, firstly studied in [10], is a trick that incorporates training samples with the trigger while retaining their correct labels, allowing attacks to evade defenses [28, 32]. This approach prevents the backdoored model from learning an overwhelmingly strong signal (more detectable in defense) for the trigger that would always lead to the target class [29]. During training, two types of triggered samples are used: 1)

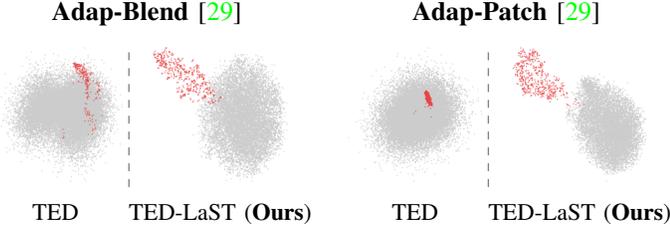


Fig. 1. T-SNE visualization of sample feature vector separability for TED and TED-LaST on backdoored CIFAR-10 models under Adap-Blend (left) and Adap-Patch (right) attacks. For each attack, the left subplot shows results for TED, while the right subplot shows results for TED-LaST. Red points indicate malicious samples and Gray points represent clean samples. The plots demonstrate improved separation between malicious and clean samples achieved by TED-LaST compared to TED.

poisoned samples with triggers labeled as the target class c_{target} , and 2) samples with triggers that maintain their original labels. This trick can be formulated as:

$$\mathcal{D}_p = \{(A(x), c_{\text{target}}) \mid (x, y) \in \mathcal{D}\}, \quad (3)$$

$$\mathcal{D}_l = \{(A(x), y) \mid (x, y) \in \mathcal{D}\}, \quad (4)$$

where \mathcal{D}_p represents the poisoned dataset and \mathcal{D}_l denotes the Laundry dataset. Both datasets collectively constitute the model’s training set.

2) *Slow Release*: Slow Release, firstly studied in [33], is a trick that incorporates training samples with partial trigger during training, while using the complete trigger during inference to activate the backdoor [33, 34]. This approach gradually introduces the backdoor, weakening the strong signal for the trigger learned by the model. The diverse partial triggers during training prevent the model creates a sudden, easily detectable correlation between the trigger and the target class [36].

For instance, during training, only portions of the trigger pattern are applied (e.g., using only two patches of a four-patch trigger), while the complete trigger pattern is restored during inference/attack phase. More formally, during training, we use a set of parameters \mathcal{R}_t (a subset of full parameter space \mathcal{R}) to control the trigger intensity or geometric attribute β_t , while during inference, we use a mapping function $g(\beta_t)$ to convert the training-phase parameter to its full-strength inference counterpart. The poisoned dataset can be formulated as:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta_t), c_{\text{target}}) \mid (x, y) \in \mathcal{D}, \beta_t \in \mathcal{R}_t \subset \mathcal{R}\}, \quad (5)$$

$$A_{\text{trig}}(x, \beta) = x \oplus I(\delta, \beta), \quad (6)$$

where A_{trig} is the trigger application function that combines the trigger with the input sample, $I(\delta, \beta)$ is a function that generates the trigger pattern δ based on a parameter β , which controls either the trigger’s intensity (e.g., opacity) or geometric attributes (e.g., size, gap, location).

3) *Target Mapping*: Target Mapping, firstly studied in [34], is a data poisoning trick that use a single, shared trigger to create diverse mappings to multiple target classes [31, 34, 35, 39]. Instead of creating a one-to-one relationship (e.g., one trigger to target class A), Target Mapping establishes a one-to-many backdoor mapping (e.g., one trigger to target classes

TABLE I
BACKDOOR ATTACK NOTATION.

Term	Description
$T(\cdot)$	Target Mapping function
$A(x)$	Function creating new sample from x
$A_{\text{trig}}(x, \beta)$	Function modifying x with trigger controlled by β
$A_{\text{trig}}(x)$	Function applying trigger for Source-Specific Target Mapping
c_{target}	Attacker’s chosen target class
$\mathcal{D}_p, \mathcal{D}_l$	Backdoor poisoning and Laundry datasets
$I(\delta, \beta)$	Function modulating trigger δ with β
\mathcal{R}	Set of trigger attribute parameters, $\beta \in \mathcal{R}$
\mathcal{R}_t	Subset of \mathcal{R} for Slow Release training phase
$g(\beta)$	Function mapping training to inference-time parameter
\mathcal{S}	Set of source classes
\oplus	Operation of mixing trigger with input

A, B, and C). The backdoor’s malicious behaviors thus become dependent on factors **beyond** just the trigger. While the trigger acts as one element, other seemingly benign features in the input data (such as specific pixel combinations or values within a particular range) can dictate which malicious task the model activates. The target mapping function is defined as:

$$T : (\mathcal{S} \cup \{\emptyset\}) \times (\mathcal{R} \cup \{\emptyset\}) \rightarrow \mathcal{Y}, \quad (7)$$

where \mathcal{S} denotes the set of source classes. For scenarios focusing solely on the source class, known as Source-Specific (SS) Target Mapping, the poisoning dataset is:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x), T(y)) \mid (x, y) \in \mathcal{D}, y \in \mathcal{S}\}, \quad (8)$$

where $A_{\text{trig}}(x) = x \oplus \delta$. SS applies the trigger uniformly without differentiating trigger attribute β . Unlike SS, the Source-Specific & Trigger Attribute (SS&TA) Target Mapping considers both source class and trigger attributes:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta), T(y, \beta)) \mid (x, y) \in \mathcal{D}, y \in \mathcal{S} \cup \{\emptyset\}, \beta \in \mathcal{R} \cup \{\emptyset\}\}. \quad (9)$$

4) *Combining Tricks*: As backdoor detection methods advance, single-trick attacks have become increasingly vulnerable to detection. In response, attackers have developed more adaptive attacks such as Adap-Blend and Adap-Patch, which combine Laundry and Slow Release [29]. The Adap-Blend attack uses partitioned, low-opacity triggers for training, and full, higher-opacity triggers for attacks. The Adap-Patch attack utilizes multiple small and diverse patches as triggers, employing combinations of fully opaque patches during the attack. To implement Laundry, both Adap-Blend and Adap-Patch inject triggers into a portion of clean samples during training, while keeping their true labels unchanged. Such adaptive attacks not only evade SOTA detection methods that rely on latent feature separability in metric space but also circumvent TED, which utilizes feature separability in topological space, as illustrated in Fig. 1.

5) *Enhanced Adaptive Attack*: Building upon these advanced combinations, we propose to further integrate even more data-agnostic tricks from the bag into the adaptive attack framework. Particularly, models compromised by a Target Mapping attack can alternate between multiple malicious tasks based on the presence of attacker-chosen backdoor features.

Importantly, this task-switching capability persists even when attacks **share** a common trigger pattern [9, 31, 34, 35], meaning the same trigger no longer consistently leads to a static malicious outcome.

Unlike Adap-Blend and Adap-Patch which are implemented with a static target class, we introduce **Enhanced Adaptive Attacks**, which combine Laundry (L), Slow Release (SR), and Target Mapping (SS or SS&TA) to create dynamic-target adaptive attacks. The key idea is to modify the trigger according to Eq. (6), thereby extending the poisoning datasets \mathcal{D}_p (as defined in Eq. (9)) and Laundry dataset \mathcal{D}_l (as defined in Eq. (4)). Specifically, to backdoor a model using the SS&TA+L+SR attack, we train the model by minimizing a loss function comprised of three components: clean loss (L_c), laundry loss (L_l), and poison loss (L_p). Building upon Eq. (2), our complete loss function is:

$$f^* = \arg \min_f (L_c + L_l + L_p), \quad (10)$$

$$L_c = \sum_{(x_i, y_i) \in \mathcal{D}} L(y_i, f(x_i)), \quad (11)$$

$$L_l = \sum_{(x_i, y_i) \in \mathcal{D}_l} L(y_i, f(A_{\text{trig}}(x_i, \beta_{t,i}))), \quad (12)$$

$$L_p = \sum_{(x_i, y_i) \in \mathcal{D}_p} L(T(y_i, g(\beta_{t,i})), f(A_{\text{trig}}(x_i, \beta_{t,i}))). \quad (13)$$

Additional configurations of Enhanced Adaptive Attacks, including SS+L+SR, are detailed in Appendix A.

C. Existing Defenses Against Adaptive Attacks

Backdoor defenses are typically categorized into three main groups based on their analysis target: model-level, label-level, and sample-level. Model-level defenses focus on analyzing the model itself. For example, a meta-classifier can be trained on a set of clean and trojaned models to identify compromised models [20]. Label-level defenses aim to reverse-engineer potential triggers and remove inserted backdoors [21, 22, 24], or analyze anomalies in learned representations [23, 40]. However, for model and label-level defenses, understanding why a model is flagged as compromised can be particularly challenging, especially under adaptive attacks that subtly alter model behavior. By contrast, sample-level defenses provide a more granular approach.

Sample-level defenses analyze the input data representations and model behavior. For instance, SCAN [10] uses robust statistics to analyze the representation distribution across classes and employs a bi-component model to disentangle class identity and variations. STRIP [26] detects triggers by overlaying input images on random samples and analyzing the entropy changes in the output labels. TeCo [27] assesses the model’s corruption robustness, identifies distinct patterns of triggered samples under various image corruptions, and quantifies the consistency of model responses as corruption severity increases.

Nevertheless, these sample-level defenses are not impervious to adaptive attacks. Specifically, adaptive attacks can diminish SCAN’s effectiveness by suppressing the potential separation between clean and poisoned representations. While

STRIP is effective against standard attacks, adaptive techniques can manipulate the entropy distribution of triggered inputs, thus blurring their distinction from benign data. Moreover, adaptive attacks can circumvent TeCo by engineering consistent behavior across corruption levels. The detailed results of sample-level defenses against various adaptive attacks can be found in our experimental analysis in Section V.

III. ORIGINAL TED ENCOUNTERS ADAPTIVE ATTACK: A CASE STUDY

Topological Evolution Dynamics (TED) leverages the evolution of neural network activations in the topological space to outlier malicious samples. The TED feature vector is a measure that captures how a sample’s activations align with its predicted class throughout the network. Specifically, TED feature vector quantifies the evolution of a sample’s feature representation across different layers by tracking its relative position within the topological space of activations. At each layer, it ranks the closeness of a sample’s activations to those of its predicted class. By tracking these rankings across layers, TED maps each sample’s evolutionary path through the network.

TED adopts a topological approach to model feature spaces, focusing on relative proximity rather than mere vector distances. This involves defining a metric space (\mathcal{V}, d) , where \mathcal{V} is a set of vectors or matrices, and d is a metric function mapping any two elements in \mathcal{V} to a non-negative real number. Each input x at layer l is represented as $h_l(x) = v \in \mathcal{V}^{(l)}$. In this space, an open ball centered at v with radius r , denoted as $\mathcal{B}(v, r)$, includes all points v' for which $d(v, v') < r$. These open balls form a topology based on neighborhood closeness. For a benign sample x_u with label y_u at the layer l , there exists another sample x' , also labeled y_u , within a minimal radius r_l such that $h_l(x')$ falls within $\mathcal{B}(v_u^{(l)}, r_l)$. This minimal radius r_l (1 by default) captures the local neighborhood structure around the sample x_u at layer l . This assumption implies that benign samples of the same class will exhibit similar activation patterns within a certain proximity.

For each input sample, TED generates a ranking list at each network layer based on its proximity to other samples from the same predicted class. The TED feature vector for an input x is defined as:

$$\text{TED}(x) = [K_1(x), K_2(x), \dots, K_l(x), \dots, K_N(x)]. \quad (14)$$

This sequence captures x ’s topological evolution across the N network layers. Here, $K_l(x)$ represents the rank of x ’s nearest neighbor from its predicted class at layer l , typically calculated using Euclidean distance. This sequential data reveals whether x consistently aligns with its class’s typical activation patterns or diverges, potentially indicating an anomaly. An outlier detector is then trained on the TED feature vectors of benign samples from all classes. This detector flags inputs with significantly different TED trajectories as potential outliers (malicious).

A. Limitations of TED

1) *Inseparability in Topological Space Outlier Detection:* Traditional attacks without adaptive approaches typically pro-

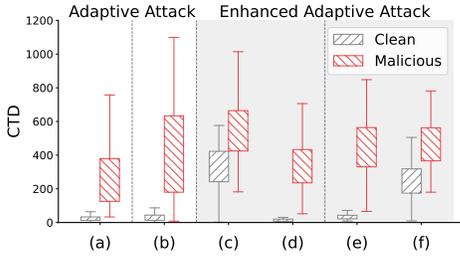


Fig. 2. Cumulative Topological Distance (CTD) of different attack scenarios on CIFAR-10. (a) Adap-Blend, (b) Adap-Patch, (c) SS+L+SR Target Class A, (d) SS+L+SR Target Class B, (e) SS&TA+L+SR Target Class A, (f) SS&TA+L+SR Target Class B. Scenarios (a) and (b) represent adaptive attacks, while (c)-(f) represent Enhanced Adaptive Attacks.

duce malicious data that exhibits significantly different characteristics from clean data in the topological space. Thus, TED assumes that malicious data has notable differences in topological evolution compared to clean data from all classes, and identifies malicious data through outlier detection.

However, this assumption becomes less robust against adaptive attacks. Such adaptive attacks manipulate the feature space and induce **inseparability** in the topological space of learned representations, as shown in Fig. 1. This induced inseparability disrupts TED’s ability to distinguish between benign and malicious samples.

The core vulnerability lies in TED’s training procedure for its outlier detector. By using benign samples from all classes during training, the detector learns a broader notion of “normal” behavior, encompassing variations across all classes. As a result, adaptive attacks can craft malicious samples that fall within this broader range of acceptable normal variation, while still achieving their malicious objectives.

Such vulnerability is particularly susceptible to exploitation by adaptive attacks, since an adaptive attacker can construct malicious samples that are sufficiently close to any class, not necessarily the target class, to evade detection, rendering TED ineffective. This fundamental limitation of considering all classes equally during outlier detector training underscores the need for more focused approaches.

2) *Insensitive to Subtle Perturbations*: In adaptive attacks, malicious samples closely emulate the topological evolution of the target class, effectively “shadowing” legitimate samples’ trajectories by maintaining minimal distance to their target class neighbors across multiple layers. Consequently, TED’s feature vector exhibits limited resolution in distinguishing these samples. The core challenge lies in TED’s insufficient sensitivity to detect subtle yet persistent deviations from the true class trajectory. While TED effectively identifies significant topological shifts, it struggles to flag samples exhibiting only slight offsets in the topological space across layers. This limitation can lead to misclassification of samples whose TED feature vector falls within the margin of error for the target class.

Such vulnerability is particularly susceptible to exploitation by adaptive attacks, which can engineer malicious samples to remain consistently within TED’s detection limits across network layers by minimizing perturbations and targeting

layers where TED’s sensitivity is lower.

B. Insight and Proposed Enhancements to TED

1) *Label-Supervised Dynamics Tracking*: The key insight is that malicious samples, originating from a different source class, traverse a **larger topological distance** compared to benign samples within the target class itself. To empirically quantify this difference in traversed distance, we introduce the Cumulative Topological Distance (CTD) metric:

$$\text{CTD}(x) = \sum_{l=1}^{N-1} |K_{l+1}(x) - K_l(x)|, \quad (15)$$

where \mathcal{X} denotes the set of all samples. Here, N is the total number of network layers, and $K_l(x)$ represents the ranking of sample x at layer l . As shown in Fig. 2, malicious samples consistently exhibit higher CTD values compared to clean samples predicted in the target class across all attacks. This observed discrepancy in the traversed distance suggests a potential distinction in the topological features of manipulated samples, motivating our label-supervised dynamics tracking approach for detecting them.

We then shift from a global to a class-specific perspective. We assume that the defender can access to a small set of clean samples with the correct label. Specifically, our method employs class-specific PCA-based outlier detection models. The model computes a class-specific threshold using a reject parameter α . This threshold is set to preserve the $(1 - \alpha)$ percentage of the variance explained by the principal components within the distribution of benign samples for that particular class. Samples exceeding this threshold are flagged as potential outliers. Leveraging this label information, we establish a granular understanding of the expected data patterns within each class. This allows us to detect deviations that would otherwise be masked by the global topological blurring caused by adaptive attacks. By focusing on the unique characteristics of each class, we can identify malicious samples even when they blend seamlessly with their target class in the global topological space.

2) *Adaptive Layer Emphasis*: While incorporating label-supervised dynamics tracking mitigates the global inseparability issue of the original TED, detecting malicious samples that closely mimic the target class in the topological space remains challenging, especially for samples with subtle perturbations. These subtle perturbations result in lower CTD values for malicious samples, indicating smaller topological distances traversed across network layers. Consequently, their feature representations become more similar to those of benign samples in the target class, making them harder to detect as outliers. Empirical observations show that the distribution of ranks for malicious samples, particularly those with CTD values in the lower quartile of their distribution, shows considerable intersection with the distribution of ranks for benign samples in the target class across most layers in the TED (Fig. 3a).

Previous studies have shown that assigning greater weight to key layers, typically the last few layers, can improve backdoor malicious detection [41, 42]. However, merely changing the

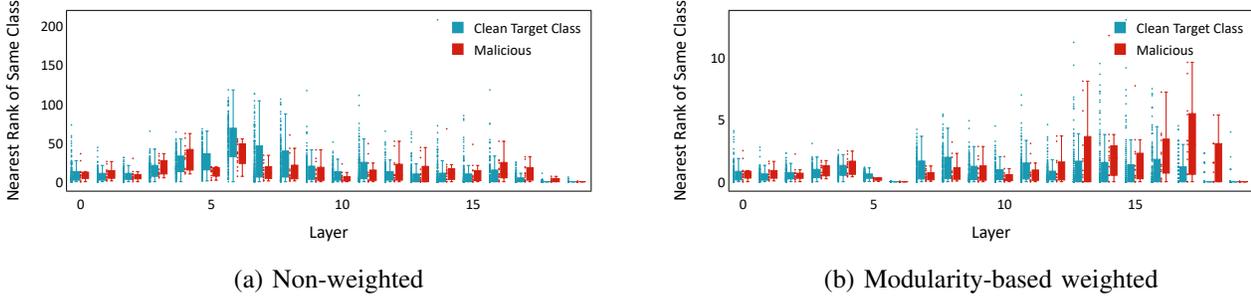


Fig. 3. (a) Box plot of the topological feature vector on ResNet20 under Adap-Patch. The plot reveals CIFAR-10 malicious samples with subtle perturbations, which are too minor for class-wise TED to effectively identify as anomalies. (b) Box plot after applying modularity-based adaptive layer weighting, showing improved separation between clean and malicious samples.

number of layers does not effectively enhance TED performance [28]. Therefore, in addition to identifying key layers at the network’s end, we must also consider the inherent variability in clean samples from the front and middle layers. To address this, we propose a method to dynamically identify and emphasize key layers across the entire network.

We adopt a modularity concept due to its effectiveness in quantifying cluster separation [43]. Our proposed modularity-based method quantifies the variability in each layer’s feature space and adaptively adjusts the weights of different layers accordingly. Specifically, modularity quantifies the degree to which a network can be divided into distinct communities or clusters [44]. We adapt this concept to the feature space, treating data points as nodes and their similarities as edge weights. The modularity score is computed by comparing the density of connections within clusters to the expected density if the connections were distributed randomly. A higher modularity score indicates a clearer separation between classes, implying the presence of well-defined clusters.

In our approach, layers with lower variability levels (higher modularity) are assigned greater weight in the final calculation of the feature vector, thereby emphasizing their contribution to the overall detection process. Conversely, layers with higher variability levels (lower modularity) are assigned smaller weight, thus reducing their contribution. By selectively emphasizing layers with different class distinctions, as depicted in Fig. 3b, we enhance the method’s resilience against the subtle differences introduced by adaptive attacks. In addition to addressing static target adaptive attacks, we provide a visualization comparing the non-weighted and modularity-based weighted methods under a dynamic-target Enhanced Adaptive Attack in Appendix B.

IV. DEFENSE DESIGN

In this section, we provide the details of our approach to detecting backdoor attacks. The framework of TED-LaST is illustrated in Fig. 4

A. Defense Design

$O_l(x)$ denotes the output of layer l for input x , and $K_l(x)$ is the rank of x ’s nearest neighbor from its predicted class at layer l . For weight calculation, we first consider all possible classes. For each class $c \in \mathcal{Y}$ and each layer $l \in \{1, \dots, N\}$,

we calculate weights $w_{l,c}$ using one of the methods described in Section IV-B.

Next, given a sample x and its predicted class \hat{y} , we compute TED-LaST’s adaptive feature vector $\text{TED}^*(x)$, building upon the original TED feature vector (Eq. (14)):

$$\text{TED}^*(x) = [K_1(x) \cdot w_{1,\hat{y}}, \dots, K_N(x) \cdot w_{N,\hat{y}}]. \quad (16)$$

For each class $c \in \mathcal{Y}$, we train a dedicated PCA-based outlier detector. This detector is trained using the set of TED^* features computed from the subset of training data predicted to be of class c : $\{\text{TED}^*(x) \mid x \in \mathcal{X}, \text{ and the predicted class of } x \text{ is } c\}$. The detector uses the sum of weighted Euclidean distances to the selected eigenvectors as a measure of anomaly [45]. For a sample x , we compute the anomaly score $s(x)$ as the sum of weighted Euclidean distances between the sample (projected onto the PCA space) and each eigenvector of the PCA model corresponding to its predicted class \hat{y} . The threshold τ_c is set as the α -quantile of the anomaly scores within the set of training samples predicted as class c .

During the inference phase, for a new sample x with predicted class \hat{y} , we first compute $\text{TED}^*(x)$. We then calculate the anomaly score $s(x)$ using the PCA-based outlier detector trained for class \hat{y} . The sample is classified as anomalous if its score exceeds the corresponding threshold $\tau_{\hat{y}}$. Algorithm 1 provides details of the modularity-based TED-LaST process.

B. Weight Calculation

We first construct a graph \mathcal{G}_l for each layer l using k -nearest neighbors (KNN) on the layer’s activations $\{O_l(x) \mid x \in \mathcal{X}\}$, and the number of neighbors is set to $\sqrt{|\mathcal{X}|}$. We then assign labels $C_c(x) = 0$ if $\hat{y} = c$, and 1 otherwise. For each layer l and class $c \in \mathcal{Y}$, we compute modularity $Q_{l,c}$ on \mathcal{G}_l using $C_c(x)$:

$$Q_{l,c} = \sum_{i=1}^{n_c} \left[mE_i - \gamma \left(\frac{k_i}{2m} \right)^2 \right], \quad (17)$$

where n_c is the number of communities (in this case, 2), m is the total number of edges in \mathcal{G}_l , E_i is the number of edges within community i , k_i is the sum of node degrees in community i , and γ is the resolution parameter (default 1). We then normalize the weights across all layers for each class:

$$w_{l,c} = \frac{Q_{l,c} - Q_{\min,c}}{Q_{\max,c} - Q_{\min,c}}, \quad (18)$$

where $Q_{\min,c}$ and $Q_{\max,c}$ are the minimum and maximum modularity values across all layers for the class c .

Algorithm 1: Modularity-based TED-LaST

```

1 Input: Set of samples  $\mathcal{X}$  with predicted labels  $\hat{y}$  for
   each  $x \in \mathcal{X}$ , set of classes  $\mathcal{Y}$ , quantile  $\alpha$  for
   threshold, total number of layers  $N$ , output of layer  $l$ 
   for input  $x$  as  $O_l(x)$ 
   /* Preprocessing */
2 for each layer  $l \in \{1, \dots, N\}$  do
3   Construct graph  $\mathcal{G}_l$  using KNN ( $k = \sqrt{|\mathcal{X}|}$ ) on
    $\{O_l(x) \mid x \in \mathcal{X}\}$ 
4   Compute  $K_l(x)$ : Rank of  $x$ 's nearest neighbor
   from its predicted class at layer  $l$ 
   /* Training Phase */
5 for each class  $c \in \mathcal{Y}$  do
6   Assign labels:  $C_c(x) = \begin{cases} 0 & \text{if } \hat{y} = c \\ 1 & \text{otherwise} \end{cases}$ 
7   for each layer  $l \in 1, \dots, N$  do
8     Compute Modularity  $Q_{l,c}$  on  $\mathcal{G}_l$  using  $C_c(x)$ 
     according to Eq. (17)
9     Normalize weights:  $w_{l,c} = \frac{Q_{l,c} - Q_{\min,c}}{Q_{\max,c} - Q_{\min,c}}$  for all  $l$ 
10    Compute
     $\text{TED}^*(x) = [K_1(x) \cdot w_{1,c}, \dots, K_N(x) \cdot w_{N,c}]$ 
11    Train PCA-based outlier detector for class  $c$  on
     $\{\text{TED}^*(x) \mid x \in \mathcal{X}, \hat{y} = c\}$ 
12    Set threshold  $\tau_c$  as the  $\alpha$ -quantile of
     $\{s(x) \mid x \in \mathcal{X}, \hat{y} = c\}$ 
   /* Inference Phase */
13 Function Detect( $x, \hat{y}$ ):
14   Compute
     $\text{TED}^*(x) = [K_1(x) \cdot w_{1,\hat{y}}, \dots, K_N(x) \cdot w_{N,\hat{y}}]$ 
    Compute anomaly score  $s(x)$  using PCA-based
    outlier detector for class  $\hat{y}$ 
15   return  $s(x) > \tau_{\hat{y}} ? \text{ANOMALOUS} : \text{NORMAL}$ 

```

C. Evaluation Metrics for Sample-Level Backdoor Detection

In our evaluation, we primarily employ two common metrics as our previous work [28]: precision and F1 score. For the defender, malicious data is considered positive, while clean data is considered negative. Precision measures the proportion of correctly identified malicious inputs among all inputs flagged as malicious. This metric reflects the accuracy of the detection system, particularly its ability to reduce false positives. The F1 score provides a balanced measure of model performance by combining precision and true positive rate (TPR, also known as recall).

To determine the detection thresholds, we analyze a ranking sequences of normal inputs for each class. The class-specific threshold τ_i is determined using a reject parameter α through PCA-based outlier detection on normal input samples. We set $\alpha = 0.05$ to achieve a 5% false positive rate (FPR), meaning that 5% of samples with the most deviant projections onto principal components are flagged as potential outliers. This

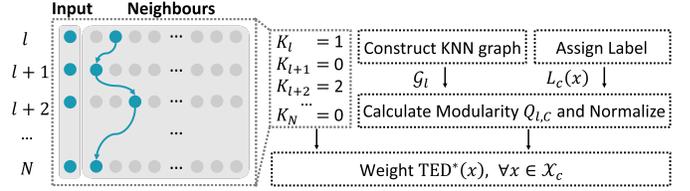


Fig. 4. Overview of TED-LaST adaptive feature vector calculation structure.

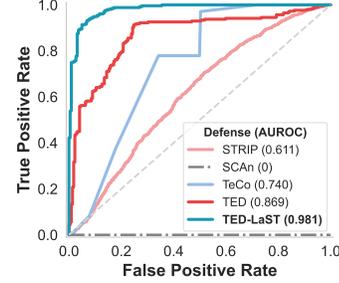


Fig. 5. Defenses against Adap-Blend adaptive attack ROC curve comparison on CIFAR-10.

FPR level is considered acceptable in practical applications. Additionally, to provide a more complete performance assessment across different FPR thresholds, we include AUROC as a supplementary metric in our ablation study.

V. EXPERIMENTS

In this section, we conduct our experiments to various adaptive attack scenarios, including the Enhanced Adaptive Attack we propose. Results for TED-LaST against two additional non-adaptive backdoor attacks are presented in Appendix C. Following the same setting as our previous work [28], TED-LaST uses 200 clean samples per class from CIFAR-10 and GTSRB dataset, and utilizes all outputs from the Conv2D and Linear layers. For implementation, we employ the open-source Python Outlier Detection (PyOD) library [45].

A. Robustness Against Adaptive Attack

To evaluate TED-LaST against adaptive attacks, we follow the implementation in [29], training ResNet-20 models on CIFAR-10 and GTSRB datasets under both Adap-Blend and Adap-Patch scenarios. Our evaluation uses 1000 poisoned and 1000 clean samples from each dataset. As shown in Tables II and III, TED-LaST maintains robust detection capabilities with precision $\geq 94.0\%$ and F1 score $\geq 92.9\%$ across all configurations. Compared to TED, TED-LaST demonstrates superior performance, achieving 35% higher F1 score for Adap-Blend attacks on CIFAR-10 and 63% higher for Adap-Patch attacks on GTSRB.

B. Evaluation on Enhanced Adaptive Attacks

Beyond SS+L+SR and SS&TA+L+SR as Enhanced Adaptive Attacks, we aim to understand the robustness of TED-LaST under different combinations of adaptive poisoning tricks. We examine three different trigger-to-target mapping

TABLE II
COMPARISON OF DETECTION (PRECISION AND F1 SCORE, IN %) FOR DIFFERENT DEFENSES AGAINST ADAP-BLEND AND ADAP-PATCH ATTACKS ON CIFAR-10.

	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score						
Adap-Blend	47.4	8.2	0	Null	90.8	64.6	87.6	69.4	94.0	93.7
Adap-Patch	16.7	1.9	0	Null	22.1	2.5	91.4	79.3	94.8	92.9

TABLE III
COMPARISON OF DETECTION (PRECISION AND F1 SCORE, IN %) FOR DIFFERENT DEFENSES AGAINST ADAP-BLEND AND ADAP-PATCH ATTACKS ON GTSRB.

	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score						
Adap-Blend	56.9	11.8	30.3	1.9	16.8	1.9	95.0	95.4	96.1	96.7
Adap-Patch	42.5	6.8	88.6	50.8	50.6	8.9	92.5	60.1	96.2	98.0

scenarios: Basic, Source-Specific Target Mapping (SS), and Source-Specific&Trigger Attribute Target Mapping (SS&TA), where Basic refers to cases where the trigger’s target remains static regardless source class or trigger attributes.

1) *Settings*: For our baseline attack configuration, we implement a 6×6 square trigger in the bottom right corner of input images [34], with a poison rate of 0.01. For Laundry implementation, we follow methodologies from previous studies [10, 28, 29], selecting training samples from non-victim classes, applying triggers, and labeling them with genuine labels. The number of Laundry samples matches the number of poisoned samples. For Slow Release implementation, due to the square trigger being too small for further partitioning, we follow [29] and use a “Hello Kitty” trigger sized equal to the input. Following [29, 33, 34], the trigger is divided into 16 segments, with random parts integrated into training samples, where each poisoned sample randomly applies half of these segments, while the complete trigger is used for testing. For Target Mapping implementation, we follow the setup described in [34], where two source classes and two distinct target classes are selected for the SS scenario. For the SS&TA scenario, to balance the ASR across different classes considering varying trigger intensities, two poison rates (*e.g.*, 0.1 and 0.08) are used to poison a single class, with different trigger densities (*e.g.*, 0.4 and 0.6 for the square trigger, and 0.15 and 0.3 for the “Hello Kitty” trigger in Slow Release scenarios) as per the configurations in [34] and [29], respectively. These varying trigger intensities ensure that poisoned samples achieve their intended target classes, respectively.

These backdoor attacks are trained on CIFAR-10 and GTSRB datasets. CIFAR-10 [46] consists of 60,000 32×32 color images across ten classes, while GTSRB [47] features over 50,000 32×32 traffic sign images across 43 classes. Both experiments employ the ResNet-20 model [48], following the setup in [29]. Table IV illustrates the performance of our proposed Enhanced Adaptive Attacks including SS&TA+L+SR and SS+L+SR on CIFAR-10 and GTSRB datasets, demonstrat-

ing adequate Attack Success Rate (ASR) and Clean Accuracy (Clean ACC) across the tested configurations.

2) *Results*: Tables V, VI, and VII demonstrate TED-LaST’s robustness against various backdoor attacks on CIFAR-10. TED-LaST consistently outperforms existing defenses including STRIP, SCAn, TeCo, and TED, especially against Enhanced Adaptive Attacks, validating our analysis in Section II-C.

For SCAn, its performance significantly degrades against Adap-Blend attacks where attackers deliberately minimize the feature representation differences between clean and poisoned samples, challenging its distribution-based detection mechanism. STRIP shows inherent limitations against adaptive attacks, particularly Adap-Patch, where attackers successfully evade entropy-based detection by weakening the correlation between backdoor triggers and target labels. Similarly, TeCo’s effectiveness diminishes when attackers deliberately blend malicious samples with clean ones in latent space, consistent with findings in [27].

The effectiveness of TED-LaST becomes more pronounced with increasing attack complexity, primarily due to: (1) Latent space inseparability: Adaptive attacks induce feature space inseparability, challenging defenses relying on clear benign-malicious separations; (2) Target Mapping complexity: Sophisticated mapping patterns reduce trigger reliability as malicious indicators, particularly affecting TED’s topological analysis. TED-LaST’s supervised-label dynamic tracking and adaptive weighting mechanisms effectively counter these challenges. This robust performance stems from TED-LaST’s ability to capture subtle class-specific anomalies through supervised-label dynamic tracking, particularly effective against attacks that blur benign-malicious distinctions in global topological space.

Results on GTSRB (Tables VIII, IX, and X) further demonstrate TED-LaST’s superiority. While defense efficacy generally declines as Target Mapping complexity increases from Basic to SS&TA, TED-LaST shows minimal performance degradation. For instance, in SS+SR scenarios on GTSRB, TED-LaST maintains a 96.6% F1 Score compared to TED’s

TABLE IV
PERFORMANCE OF VARIOUS ADAPTIVE ATTACKS WITH TRICKS ON CIFAR-10 AND GTSRB DATASETS. SETTINGS: B (BASIC), L (LAUNDRY), SR (SLOW RELEASE), SS (SOURCE-SPECIFIC), TA (TRIGGER ATTRIBUTE)

Setting	Target Mapping	Laundry	Slow Release	CIFAR-10		GTSRB	
				ASR (%)	Clean ACC (%)	ASR (%)	Clean ACC (%)
B	-	-	-	99.8	86.4	100	92.9
L	-	✓	-	100	80.2	100	94.9
SR	-	-	✓	100	83.0	100	96.2
L+SR	-	✓	✓	69.7	82.6	100	97.2
SS	SS	-	-	97.5	77.9	100	95.7
SS+L	SS	✓	-	99.8	80.7	100	100
SS+SR	SS	-	✓	76.9	83.3	100	96.8
SS+L+SR	SS	✓	✓	79.0	79.1	100	96.9
SS&TA	SS&TA	-	-	100	81.0	100	95.7
SS&TA+L	SS&TA	✓	-	100	80.3	100	100
SS&TA+SR	SS&TA	-	✓	78.3	83.2	100	97.7
SS&TA+L+SR	SS&TA	✓	✓	84.1	78.8	100	99.8

TABLE V
PERFORMANCE (PRECISION AND F1 SCORE, IN %) OF DIFFERENT DEFENSES AGAINST ATTACKS BASIC TARGET MAPPING ON CIFAR-10 USING RESNET-20. SETTINGS: B (BASIC), L (LAUNDRY), SR (SLOW RELEASE)

Setting	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score
B	95.3	97.3	94.1	86.5	74.5	29.1	94.8	93.1	96.7	98.0
L	63.6	7.9	91.0	64.8	87.1	47.9	94.1	88.2	92.7	90.3
SR	20.0	3.0	89.3	58.2	16.7	1.9	96.3	97.4	98.9	99.5
L+SR	50.9	10.4	91.0	64.8	70.3	20.6	93.0	92.2	96.1	97.8

TABLE VI
PERFORMANCE (PRECISION AND F1 SCORE, IN %) OF DIFFERENT DEFENSES AGAINST ATTACKS WITH SS (SOURCE-SPECIFIC) TARGET MAPPING ON CIFAR-10 USING RESNET-20. SETTINGS: SS (SOURCE-SPECIFIC), L (LAUNDRY), SR (SLOW RELEASE)

Setting	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score
SS	46.2	8.7	85.5	44.0	20.0	1.9	95.4	93.9	94.1	95.5
SS+L	94.4	94.6	83.3	38.5	18.1	1.9	90.8	83.3	93.3	94.6
SS+SR	28.6	4.6	83.2	38.1	18.1	1.8	94.7	75.3	92.9	88.0
SS+L+SR	58.7	13.1	82.9	37.5	16.7	1.9	88.0	67.1	93.1	91.0

95.2%. Even under the most complex SS&TA+L+SR attack on CIFAR-10, TED-LaST achieves an 87.9% F1 Score, substantially outperforming TED’s 76.4%.

C. Effectiveness of TED-LaST Against Enhanced Adaptive Attacks on Large-Scale Dataset

To validate the scalability and effectiveness of TED-LaST, we evaluate it on ImageNet100, a subset of ImageNet [49] comprising 100 classes. Our experiments use 224x224 pixel images, a ResNet101 model [48], and a resized “Hello Kitty” trigger covering the entire image. The increased class diversity of ImageNet100 introduces higher complexity in the topological space. Benign samples navigate through a more intricate manifold, traversing larger topological distances compared to smaller datasets. This poses challenges to our defense system, requiring it to distinguish between larger natural variations and subtle shifts induced by attacks. Despite these challenges, as shown in Table XI, TED-LaST demonstrates robust performance across various Enhanced Adaptive Attack

configurations. It consistently maintains high Precision rates (often exceeding 94%) and strong F1 Scores (mostly above 90%) across different attack settings. Even in the most sophisticated attack scenarios (SS&TA+L+SR), TED-LaST achieves an F1 Score of 97.2%, showcasing its ability to adapt to complex attack patterns. This performance highlights TED-LaST’s strong scalability for larger-scale, more complex image tasks and real-world applications.

VI. ABLATION STUDY

A. Label Information in Outlier Detection

This ablation study examines the importance of class-specific information in enhancing defense robustness against adaptive attacks. Fig. 6 shows the effect of varying training data composition on outlier detection performance for AdapBlend (Fig. 6a) and AdapPatch (Fig. 6b) attacks.

Training on clean samples from the predicted class alone yields the highest AUROC scores. Including additional random classes degrades performance, with the first unrelated class

TABLE VII

PERFORMANCE (PRECISION AND F1 SCORE, IN %) OF DIFFERENT DEFENSES AGAINST ATTACKS WITH SS&TA (SOURCE-SPECIFIC&TRIGGER ATTRIBUTE) TARGET MAPPING ON CIFAR-10 UNDER RESNET-20. SETTINGS: SS&TA (SOURCE-SPECIFIC&TRIGGER ATTRIBUTE), L (LAUNDRY), SR (SLOW RELEASE)

Setting	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score
SS&TA	83.3	17.9	85.3	42.7	10.2	1.7	96.5	90.7	94.0	92.8
SS&TA+L	84.4	44.7	81.3	34.4	25.0	4.9	91.2	86.5	93.4	89.9
SS&TA+SR	12.1	1.5	81.9	35.4	16.4	1.9	89.4	58.7	92.0	85.6
SS&TA+L+SR	46.6	11.9	85.2	43.0	16.7	1.9	89.5	76.4	91.5	87.9

TABLE VIII

PERFORMANCE (PRECISION AND F1 SCORE, IN %) OF DIFFERENT DEFENSES AGAINST BASIC TARGET MAPPING ATTACKS ON GTSRB UNDER RESNET-20. SETTINGS: B (BASIC), L (LAUNDRY), SR (SLOW RELEASE)

Setting	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score
B	91.6	87.3	90.9	79.3	89.9	57.7	95.1	97.0	95.1	97.0
L	79.7	21.8	89.2	56.7	63.2	15.1	95.5	97.7	95.5	97.7
SR	3.3	0.4	91.0	64.9	88.0	51.8	94.6	97.2	95.5	97.7
L+SR	5.3	0.4	89.8	59.1	78.0	29.4	94.6	97.2	95.5	97.7

TABLE IX

PERFORMANCE (PRECISION AND F1 SCORE, IN %) OF DIFFERENT DEFENSES AGAINST ATTACKS WITH SS (SOURCE-SPECIFIC) TARGET MAPPING ON GTSRB UNDER RESNET-20. SETTINGS: SS (SOURCE-SPECIFIC), L (LAUNDRY), SR (SLOW RELEASE)

Setting	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score						
SS	45.7	11.2	84.6	41.5	65.0	13.6	94.7	97.3	94.9	97.4
SS+L	15.6	1.9	94.8	93.1	78.0	29.4	94.5	97.2	94.9	97.4
SS+SR	3.3	0.4	79.1	30.5	87.4	49.7	95.9	95.2	96.0	96.6
SS+L+SR	3.2	0.4	95.0	95.2	56.3	12.0	94.6	95.1	94.8	97.2

TABLE X

PERFORMANCE (PRECISION AND F1 SCORE, IN %) OF DIFFERENT DEFENSES AGAINST ATTACKS WITH SS&TA (SOURCE-SPECIFIC & TRIGGER-ATTRIBUTE) TARGET MAPPING ON GTSRB UNDER RESNET-20. SETTINGS: SS&TA (SOURCE-SPECIFIC&TRIGGER ATTRIBUTE), L (LAUNDRY), SR (SLOW RELEASE)

Setting	STRIP		SCAn		TeCo		TED		TED-LaST	
	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score	Precision	F1 Score
SS&TA	4.5	0.4	74.5	24.4	87.4	49.7	94.3	97.1	94.3	97.1
SS&TA+L	67.1	16.5	95.0	95.3	16.9	1.9	94.9	96.7	95.9	97.9
SS&TA+SR	31.0	6.5	84.3	40.7	16.7	1.9	95.3	97.6	95.7	97.8
SS&TA+L+SR	2.9	0.4	94.9	94.7	17.5	1.9	95.4	95.6	95.7	97.7

causing the most significant drop. Further additions result in continued AUROC score declines at a decreasing rate.

This degradation comes from the dilution of class-specific topological features in the training data. As unrelated classes are introduced, the detector’s ability to identify attack-induced deviations from expected class patterns diminishes. This observation aligns with the analysis of inseparability in topological space (Section III-A1), where adaptive attacks blur distinctions between clean and poisoned samples.

B. Adaptive Weighting Under Different CTD Thresholds

We investigate the effectiveness of our modularity-based weighted method compared to the non-weighted approach across various Cumulative Topological Distance (CTD) thresholds. The CTD metric (Eq. (15)) quantifies the topological

distance change as a sample traverses network layers. As shown in Fig. 2a and Fig. 2b, both Adap-Blend and Adap-Patch attacks generate malicious samples with varying CTD values. Samples with lower CTD values, representing more subtle perturbations, are generally more challenging to detect.

We define the CTD Threshold Ratio as the median CTD divided by the actual CTD value. Higher ratios correspond to lower CTD values, indicating more subtle perturbations. Our experiments use datasets with equal numbers of malicious and clean samples, filtered based on different CTD thresholds. Fig. 7 illustrates the performance difference between the modularity-weighted and non-weighted methods across CTD Threshold Ratios.

The results reveal a clear trend: as the CTD Threshold Ratio

TABLE XI
ATTACK PERFORMANCE AND TED-LAST EFFECTIVENESS ON
IMAGENET100

Setting	Attack Performance		TED-LaST Performance	
	ASR (%)	Clean ACC (%)	Precision (%)	F1 Score (%)
B	100	83.2	93.0	92.5
L	100	82.4	95.4	93.1
SR	100	83.4	95.0	95.0
L+SR	87.0	78.5	94.2	87.1
SS	99.8	82.4	95.8	95.5
SS+L	100	82.7	94.4	95.4
SS+SR	81.3	84.2	95.0	97.2
SS+L+SR	88.8	82.3	94.8	95.5
SS&TA	100	80.1	96.2	98.1
SS&TA+L	100	81.1	94.7	95.1
SS&TA+SR	94.8	84.1	91.9	89.9
SS&TA+L+SR	88.8	83.6	94.8	97.2

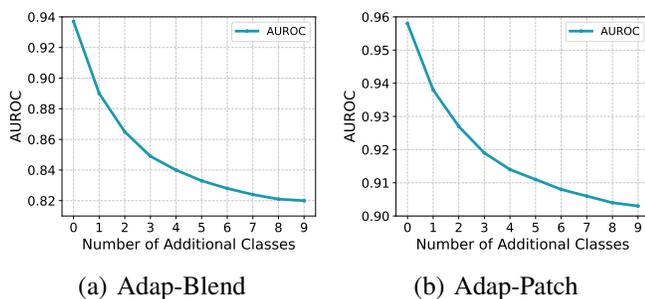


Fig. 6. AUROC for outlier detection with varying numbers of additional classes under (a) Adap-Blend and (b) Adap-Patch attacks on CIFAR-10. The x-axis represents the number of additional random classes included in the training data, while the y-axis shows the corresponding AUROC scores.

increases (i.e., when we focus on malicious samples with lower CTD values), the modularity-based weighted method consistently outperforms the non-weighted method. This performance gap becomes more pronounced at higher CTD Threshold Ratios.

C. Intuition Behind TED-LaST Against Adaptive Attacks

Adaptive backdoor attacks pose a significant challenge by concealing poisoned data within the complex topological space of neural networks. As visualized in Fig. 1a and Fig. 1b, these attacks blur the boundaries between clean and poisoned samples.

To address these challenges, our defense strategy shifts from global analysis to a more granular approach. This strategy is based on the premise that even subtle manipulations leave detectable class-wise traces. Specifically, adaptive attacks inevitably cause deviations from benign behavior within the target class (discussed in Section III-B1). These deviations, while subtle, provide crucial indicators for our detection method. We also find that additional class augmentation does not enhance detection capabilities (discussed in Section VI-A).

Our modularity-based weighted method consistently outperforms non-weighted approaches, especially for malicious samples with subtle perturbations (Section VI-B). This improved sensitivity stems from the weighting scheme’s emphasis on

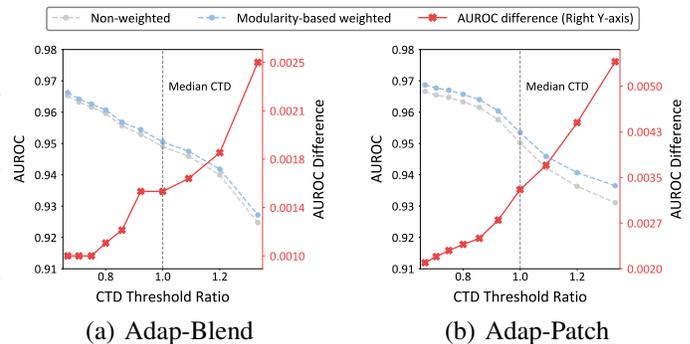


Fig. 7. Comparison of AUROC differences across various CTD Threshold Ratios. The modularity-based weighted method consistently shows higher AUROC than the non-weighted approach under all ratios. The gap widens as the ratio increases, indicating that the modularity-based weighted method provides more sensitivity to malicious samples with subtle perturbations.

the most informative topological features, enabling refined differentiation between benign and malicious samples.

Importantly, our adaptive weighting approach remains effective even against Enhanced Adaptive Attacks. These advanced attacks use dynamic target classes and shared triggers to deeply mix malicious samples with benign representations in the topological space. Despite this sophisticated mixing, our method can still detect the subtle perturbations introduced by these attacks (illustrated in Appendix B).

VII. CONCLUSION

This study introduces TED-LaST, a novel defense strategy against adaptive backdoor attacks in DNNs. TED-LaST leverages persistent topological perturbations on target classes and uses supervised label information to enhance differentiation between poisoned and clean samples. We implement adaptive layer emphasis to address subtle perturbations caused by these attacks. Our approach outperforms several state-of-the-art defenses in countering SOTA adaptive attacks and our proposed Enhanced Adaptive Attacks. Future work could include further improving the defense effectiveness while exploring the application of TED-LaST in federated learning or multi-modal learning scenarios. These extensions would contribute to enhancing the security of robust machine learning.

REFERENCES

- [1] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Bennamoun, “Deep learning for 3d point clouds: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 12, pp. 4338–4364, 2021.
- [2] T. Afouras, J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, “Deep audio-visual speech recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 12, pp. 8717–8727, 2022.
- [3] J. Cao, Y. Pang, J. Xie, F. S. Khan, and L. Shao, “From handcrafted to deep features for pedestrian detection: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 9, pp. 4913–4934, 2022.
- [4] T. Gu, B. Dolan-Gavitt, and S. Garg, “Badnets: Identifying vulnerabilities in the machine learning model supply chain,” *arXiv preprint arXiv:1708.06733*, 2017.

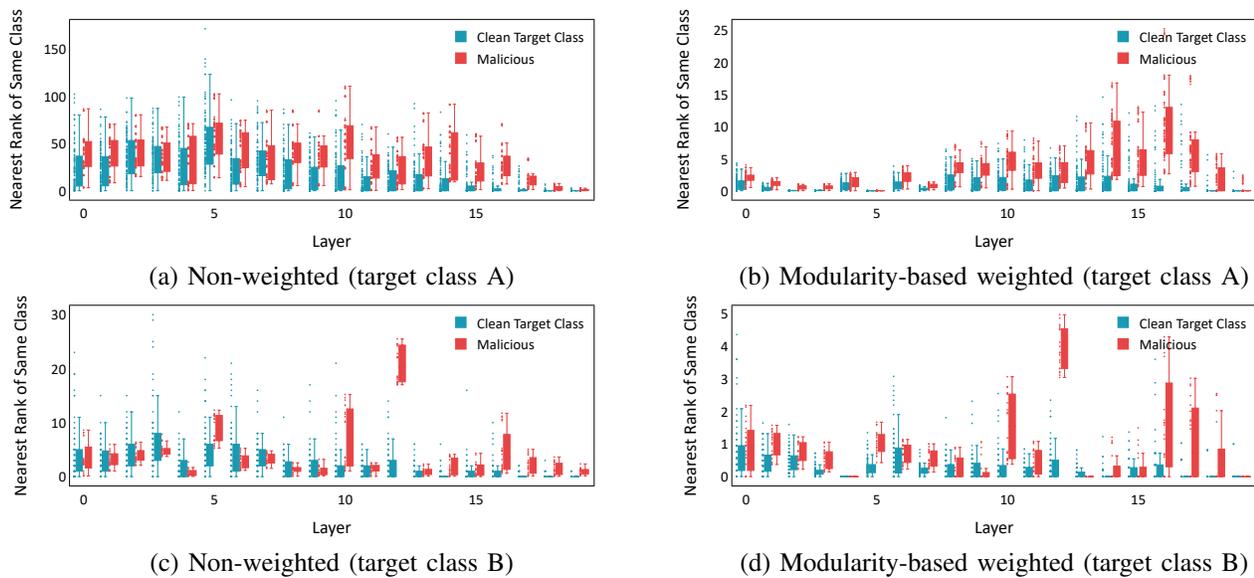


Fig. 8. Box plots of topological feature vectors against SS&TA+L+SR with dynamic targets. (a) and (c) show unweighted features for two different target classes, revealing subtle perturbations that limit discrimination. (b) and (d) demonstrate improved separation between clean and malicious samples after applying modularity-based adaptive layer emphasis for the respective target classes.

- [5] L. Li, D. Song, X. Li, J. Zeng, R. Ma, and X. Qiu, “Backdoor attacks on pre-trained models by layerwise weight poisoning,” *arXiv preprint arXiv:2108.13888*, 2021.
- [6] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, “Trojaning attack on neural networks,” in *25th Annual Network And Distributed System Security Symposium (NDSS 2018)*, 2018.
- [7] W. Yang, L. Li, Z. Zhang, X. Ren, X. Sun, and B. He, “Be careful about poisoned word embeddings: Exploring the vulnerability of the embedding layers in NLP models,” in *ACL*, 2021, pp. 2048–2058.
- [8] Y. Adi, C. Baum, M. Cisse, B. Pinkas, and J. Keshet, “Turning your weakness into a strength: Watermarking deep neural networks by backdooring,” in *USENIX Security*, 2018, pp. 1615–1631.
- [9] J. Lin, L. Xu, Y. Liu, and X. Zhang, “Composite backdoor attack for deep neural network by mixing existing benign features,” in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 2020, pp. 113–131.
- [10] D. Tang, X. Wang, H. Tang, and K. Zhang, “Demon in the variant: Statistical analysis of dnns for robust backdoor contamination detection,” in *USENIX Security Symposium*, 2021, pp. 1541–1558.
- [11] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, “Targeted backdoor attacks on deep learning systems using data poisoning,” *arXiv preprint arXiv:1712.05526*, 2017.
- [12] Y. Zeng, W. Park, Z. M. Mao, and R. Jia, “Rethinking the backdoor attacks’ triggers: A frequency perspective,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 473–16 481.
- [13] E. Quiring and K. Rieck, “Backdooring and poisoning neural networks with image-scaling attacks,” in *2020 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2020, pp. 41–47.
- [14] T. A. Nguyen and A. Tran, “Input-aware dynamic backdoor attack,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 3454–3464, 2020.
- [15] Y. Li, Y. Li, B. Wu, L. Li, R. He, and S. Lyu, “Invisible backdoor attack with sample-specific triggers,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 463–16 472.
- [16] Q. Duan, Z. Hua, Q. Liao, Y. Zhang, and L. Y. Zhang, “Conditional backdoor attack via jpeg compression,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 18, 2024, pp. 19 823–19 831.
- [17] X. Mo, L. Y. Zhang, N. Sun, W. Luo, and S. Gao, “Backdoor attack on deep neural networks in perception domain,” in *2023 International Joint Conference on Neural Networks (IJCNN)*, 2023, pp. 01–08.
- [18] C. Wu, D. Lian, Y. Ge, Z. Zhu, and E. Chen, “Influence-driven data poisoning for robust recommender systems,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 10, pp. 11 915–11 931, 2023.
- [19] K. Ma, Q. Xu, J. Zeng, X. Cao, and Q. Huang, “Poisoning attack against estimating from pairwise comparisons,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6393–6408, 2022.
- [20] X. Xu, Q. Wang, H. Li, N. Borisov, C. A. Gunter, and B. Li, “Detecting ai trojans using meta neural analysis,” in *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2021, pp. 103–120.
- [21] Y. Liu, W.-C. Lee, G. Tao, S. Ma, Y. Aafer, and X. Zhang, “Abs: Scanning neural networks for back-doors by artificial brain stimulation,” in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 1265–1282.
- [22] 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 *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2019, pp. 707–723.
- [23] W. Guo, B. Tondi, and M. Barni, “Universal detection of backdoor attacks via density-based clustering and centroids analysis,” *IEEE Transactions on Information Forensics and Security*, 2023.
- [24] H. Zhu, Y. Zhao, S. Zhang, and K. Chen, “Neuralsanitizer: Detecting backdoors in neural networks,” *IEEE Transactions on Information Forensics and Security*, 2024.
- [25] B. Tran, J. Li, and A. Madry, “Spectral signatures in backdoor attacks,” *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [26] Y. Gao, C. Xu, D. Wang, S. Chen, D. C. Ranasinghe, and S. Nepal, “Strip: A defence against trojan attacks on deep neural networks,” in *Proceedings of the 35th Annual Computer Security Applications Conference*, 2019, pp. 113–125.
- [27] X. Liu, M. Li, H. Wang, S. Hu, D. Ye, H. Jin, L. Wu,

- and C. Xiao, “Detecting backdoors during the inference stage based on corruption robustness consistency,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 16 363–16 372.
- [28] X. Mo, Y. Zhang, L. Y. Zhang, W. Luo, N. Sun, S. Hu, S. Gao, and Y. Xiang, “Robust backdoor detection for deep learning via topological evolution dynamics,” in *2024 IEEE Symposium on Security and Privacy (SP)*, 2024, pp. 2048–2066.
- [29] X. Qi, T. Xie, Y. Li, S. Mahloujifar, and P. Mittal, “Revisiting the assumption of latent separability for backdoor defenses,” in *The eleventh international conference on learning representations*, 2022.
- [30] T. J. L. Tan and R. Shokri, “Bypassing backdoor detection algorithms in deep learning,” in *2020 IEEE European Symposium on Security and Privacy (EuroSP)*, 2020, pp. 175–183.
- [31] E. Bagdasaryan and V. Shmatikov, “Blind backdoors in deep learning models,” in *Usenix Security*, 2021.
- [32] H. Peng, H. Qiu, H. Ma, S. Wang, A. Fu, S. F. Al-Sarawi, D. Abbott, and Y. Gao, “On model outsourcing adaptive attacks to deep learning backdoor defenses,” *IEEE Transactions on Information Forensics and Security*, 2024.
- [33] C. Xie, K. Huang, P. Y. Chen, and B. Li, “Dba: Distributed backdoor attacks against federated learning,” in *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- [34] M. Xue, C. He, J. Wang, and W. Liu, “One-to-n & n-to-one: Two advanced backdoor attacks against deep learning models,” *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 3, pp. 1562–1578, 2020.
- [35] Y. Xiao, L. Cong, Z. Mingwen, W. Yajie, L. Xinrui, S. Shuxiao, M. Yuexuan, and Z. Jun, “A multitarget backdooring attack on deep neural networks with random location trigger,” *International Journal of Intelligent Systems*, vol. 37, no. 3, pp. 2567–2583, 2022.
- [36] J. Hayase, W. Kong, R. Somani, and S. Oh, “Spectre: Defending against backdoor attacks using robust statistics,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 4129–4139.
- [37] J. Dumford and W. Scheirer, “Backdooring convolutional neural networks via targeted weight perturbations,” in *2020 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 2020, pp. 1–9.
- [38] R. Tang, M. Du, N. Liu, F. Yang, and X. Hu, “An embarrassingly simple approach for trojan attack in deep neural networks,” in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 218–228.
- [39] M. Edraki, N. Karim, N. Rahnavard, A. Mian, and M. Shah, “Odyssey: Creation, analysis and detection of trojan models,” *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 4521–4533, 2021.
- [40] B. Chen, W. Carvalho, N. Baracaldo, H. Ludwig, B. Edwards, T. Lee, I. Molloy, and B. Srivastava, “Detecting backdoor attacks on deep neural networks by activation clustering,” *arXiv preprint arXiv:1811.03728*, 2018.
- [41] K. Huang, Y. Li, B. Wu, Z. Qin, and K. Ren, “Backdoor defense via decoupling the training process,” in *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=TySnJ-0RdKI>
- [42] N. M. Jebreel, J. Domingo-Ferrer, and Y. Li, “Defending against backdoor attacks by layer-wise feature analysis,” in *Advances in Knowledge Discovery and Data Mining*, H. Kashima, T. Ide, and W.-C. Peng, Eds. Cham: Springer Nature Switzerland, 2023, pp. 428–440.
- [43] M. E. J. Newman, “Modularity and community structure in networks,” *Proceedings of the National Academy of Sciences*, vol. 103, no. 23, pp. 8577–8582, 2006. [Online]. Available: <https://www.pnas.org/doi/abs/10.1073/pnas.0601602103>
- [44] S. Fortunato, “Community detection in graphs,” *Physics Reports*, vol. 486, no. 3–5, p. 75–174, Feb. 2010. [Online]. Available: <http://dx.doi.org/10.1016/j.physrep.2009.11.002>
- [45] Y. Zhao, Z. Nasrullah, and Z. Li, “Pyod: A python toolbox for scalable outlier detection,” *Journal of Machine Learning Research*, vol. 20, no. 96, pp. 1–7, 2019. [Online]. Available: <http://jmlr.org/papers/v20/19-011.html>
- [46] “The cifar-10 dataset,” <https://www.cs.toronto.edu/~kriz/cifar.html>.
- [47] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, “Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition,” *Neural networks*, vol. 32, pp. 323–332, 2012.
- [48] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [49] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

APPENDIX A

BAG OF TRICKS FOR ADAPTIVE ATTACKS DATA POISONING

A. Laundry

The Laundry dataset \mathcal{D}_1 has evolved to adaptively conceal the malicious backdoor effect from \mathcal{D}_p by introducing a conditional trigger [9, 10, 17, 28, 29]. As demonstrated in [32], model regularization can also be employed to boost the Laundry by aligning backdoored model’s parameters with those of a clean counterpart. Laundry entails incorporating data that appears similar to the triggers but with correct labels. The goal is to reduce the network’s immediate trigger sensitivity by merging trigger features with those of the target class, thus lessening immediate reactivity [28].

B. Slow Release

Slow Release aims to reduce the model’s dependence on trigger patterns in individual data instances during training, thereby minimizing the overall visibility of these triggers. This method involves decreasing the trigger prevalence in each training sample while maintaining the full trigger for inference. As discussed in [33, 34], this approach may include fragmenting the trigger into multiple segments to dilute its immediate influence, or gradually “poisoning” the model throughout the training process. When the cumulative impact of these segments—termed aggregate toxicity—reaches a certain threshold, it results in an increase in the ASR while preserving the Clean ACC.

C. Target Mapping

During the poisoning process, the attacker manipulates source-to-target class mappings to control the backdoored model’s behavior. This can involve many-to-many mapping, where each true label is assigned a different target label, or many-to-one mapping, where all true labels of triggered data are changed to a fixed target label [39]. A backdoored model can switch between malicious tasks based on specific backdoor features within an input [31]. These features may be the trigger attribute alone (e.g., location [35] or intensity [34]), or a combination of trigger attribute and source class [31].

D. Laundry and Slow Release

When combining Laundry and Slow Release techniques, as seen in Adap-Blend and Adap-Patch backdoor attacks [29], both the poisoned dataset \mathcal{D}_p and the Laundry dataset \mathcal{D}_l utilize the same modification function A_{trig} (Eq. (6)). This results in the following definitions:

$$\mathcal{D}_p = \{((A_{\text{trig}}(x, \beta_t), c_{\text{target}}) \mid (x, y) \in \mathcal{D}), \quad (19)$$

$$\mathcal{D}_l = \{((A_{\text{trig}}(x, \beta_t), y) \mid (x, y) \in \mathcal{D}). \quad (20)$$

E. Slow Release and Target Mapping

Combining Slow Release with two types of Target Mapping, we incorporate Target Mapping into the training process. Using the modified trigger (Eq. (6)) in Eq. (8) and Eq. (9), we define the dataset \mathcal{D}_p as follows:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta_t), T(y, g(\beta_t))) \mid (x, y) \in \mathcal{D}, \quad (21)$$

$$y \in \mathcal{S}, \beta_t \in \mathcal{R}_t\},$$

or if considering source-specific only:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta_t), T(y)) \mid (x, y) \in \mathcal{D}, \quad (22)$$

$$y \in \mathcal{S}, \beta_t \in \mathcal{R}_t\}.$$

F. Laundry and Target Mapping

Combining Laundry with two types of Target Mapping, we randomly select two victim classes, two target classes, and two Laundry classes. The poisoning dataset comprises samples from victim classes, mixed with the trigger and labeled as the target class. The Laundry data, randomly selected from non-victim classes, are mixed with the trigger but retain their original labels. Integrating Eq. (4) with Eq. (8) and Eq. (9), we define:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta), T(y, \beta)) \mid (x, y) \in \mathcal{D}, \quad (23)$$

$$y \in \mathcal{S} \cup \{\emptyset\}, \beta \in \mathcal{R} \cup \{\emptyset\}\},$$

$$\mathcal{D}_l = \{(A_{\text{trig}}(x), y) \mid (x, y) \in \mathcal{D}\}. \quad (24)$$

For situations that concentrate solely on the source-specific:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x), T(y)) \mid (x, y) \in \mathcal{D}, y \in \mathcal{S}\}, \quad (25)$$

$$\mathcal{D}_l = \{(A_{\text{trig}}(x), y) \mid (x, y) \in \mathcal{D}\}. \quad (26)$$

G. Laundry, Slow Release and Target Mapping

We modify the poisoning datasets \mathcal{D}_p (Eq. (9)) and Laundry dataset \mathcal{D}_l (Eq. (4)) using the altered trigger defined in Eq. (6). For the SS&TA+L+SR approach, the poisoned dataset \mathcal{D}_p and the Laundry dataset \mathcal{D}_l are defined as:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta_t), T(y, g(\beta_t))) \mid (x, y) \in \mathcal{D}, \quad (27)$$

$$y \in \mathcal{S}, \beta_t \in \mathcal{R}_t\},$$

$$\mathcal{D}_l = \{(A_{\text{trig}}(x, \beta_t), y) \mid (x, y) \in \mathcal{D}, \beta_t \in \mathcal{R}_t\}. \quad (28)$$

For scenarios focusing solely on source-specific manipulation (SS+L+SR), the poisoned dataset and the Laundry dataset are constructed as:

$$\mathcal{D}_p = \{(A_{\text{trig}}(x, \beta_t), T(\beta)) \mid (x, y) \in \mathcal{D}, \beta_t \in \mathcal{R}_t, \beta \in \mathcal{R}\}, \quad (29)$$

$$\mathcal{D}_l = \{(A_{\text{trig}}(x, \beta_t), y) \mid (x, y) \in \mathcal{D}, \beta_t \in \mathcal{R}_t\}. \quad (30)$$

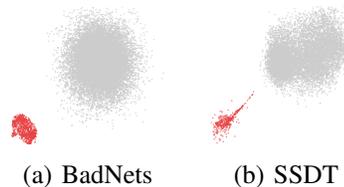


Fig. 9. T-SNE visualization of TED-LaST feature vectors for BadNets and SSDT attacks on CIFAR-10. Red points represent malicious samples and gray points represent clean samples.

APPENDIX B

ADAPTIVE WEIGHTING FOR DYNAMIC TARGET CLASSES UNDER SUBTLE PERTURBATIONS

We evaluate our adaptive weighting methods on SS&TA+L+SR samples with subtle perturbations. Fig. 8 compares topological feature vectors across network layers, with and without our weighting method, for two dynamic target classes. The non-weighted scenario (Fig. 8a and 8c) shows significant overlap between clean and malicious sample features, especially in early layers, due to SS&TA+L+SR's subtle perturbations being masked by inherent target class variations. Our modularity-based weighting method (Fig. 8b and 8d) addresses this by adaptively adjusting layer importance. It identifies layers with clearer distinctions between clean and malicious samples, assigning higher weights to amplify discriminative features. The visualization demonstrates improved separation between clean and malicious samples for both dynamic target classes, particularly in later layers where our weighting scheme has the greatest impact.

APPENDIX C

ROBUSTNESS AGAINST MORE EXISTING BACKDOOR ATTACKS

We evaluate the robustness of TED-LaST against two backdoor attacks: BadNets and SSDT. BadNets [4], the most common backdoor attack, uses a gray square on the bottom right as the trigger with a 0.01 poison rate. For Source-Specific Dynamic Trigger (SSDT) [28], which is an attack that applies Laundry and also uses dynamic triggers, poison and clean samples consistently form two separate clusters in the topological space. SSDT employs source-specific dynamic triggers, activating only for victim class samples to reduce impact on non-victim classes. For SSDT, we replicated the original experimental settings using the authors' source code on CIFAR-10. TED-LaST achieved an AUROC of 1 for both BadNets and SSDT, indicating perfect separation of clean and poisoned samples. Fig. 9 visualizes this separation through T-SNE plots.