# CLIP-Guided Backdoor Defense through Entropy-Based Poisoned Dataset Separation

Binyan Xu
The Chinese University of Hong Kong
Hong Kong, Hong Kong
binyxu@ie.cuhk.edu.hk

Fan Yang
The Chinese University of Hong Kong
Hong Kong, Hong Kong
yf020@ie.cuhk.edu.hk

Xilin Dai
Zhejiang University
Hangzhou, China
xilin2023@zju.edu.cn

Di Tang*
Sun Yat-sen University
Shenzhen, China
tangd9@mail.sysu.edu.cn

Kehuan Zhang*
The Chinese University of Hong Kong
Hong Kong, Hong Kong
khzhang@ie.cuhk.edu.hk

## Abstract

Deep Neural Networks (DNNs) are susceptible to backdoor attacks, where adversaries poison training data to implant backdoor into the victim model. Current backdoor defenses on poisoned data often suffer from high computational costs or low effectiveness against advanced attacks like clean-label and clean-image backdoors. To address them, we introduce **CLIP-G**uided backdoor **D**efense (CGD), an efficient and effective method that mitigates various backdoor attacks. CGD utilizes a publicly accessible CLIP model to identify inputs that are likely to be clean or poisoned. It then retrains the model with these inputs, using CLIP's logits as a guidance to effectively neutralize the backdoor. Experiments on 4 datasets and 11 attack types demonstrate that CGD reduces attack success rates (ASRs) to below 1% while maintaining clean accuracy (CA) with a maximum drop of only 0.3%, outperforming existing defenses. Additionally, we show that clean-data-based defenses can be adapted to poisoned data using CGD. Also, CGD exhibits strong robustness, maintaining low ASRs even when employing a weaker CLIP model or when CLIP itself is compromised by a backdoor. These findings underscore CGD's exceptional efficiency, effectiveness, and applicability for real-world backdoor defense scenarios. Code: https://github.com/binyxu/CGD.

## CCS Concepts

• **Computing methodologies** → **Machine learning**; • **Security and privacy** → **Software and application security**.

## Keywords

Backdoor Defense, Contrastive Language-Image Pretraining

*Corresponding author.

## 1 Introduction

Deep Neural Networks (DNNs) are widely used in applications like facial recognition [1], autonomous driving [12], and medical image diagnosis [24]; however, backdoor attacks threaten their trustworthiness. By poisoning a small portion of the training data [20], adversaries can inject backdoors that cause models to make erroneous predictions when specific inputs are presented. Since the training data collection is usually time-consuming and expensive, it is common to use external data for training without security guarantees. The common practice makes backdoor attacks feasible in real-world applications, which highlights the importance of backdoor removal in poisoned datasets.

Existing backdoor defenses against poisoned data fall into two categories. (1) *Clean model-based defenses* [4, 8, 14] use self-/semi-supervised learning to train a clean model from the training dataset without relying on potentially compromised labels. These defenses identify backdoors by analyzing the behavior of this clean model. For example, DBD [14] employs self-supervised learning to train a clean encoder using only image data, whereas ASD [8] starts with a weak clean model and progressively improves it in a semi-supervised way. (2) *Suspicious model-based defenses* [23, 50] identify backdoors by analyzing the behavior of the potentially compromised model itself when processing with both benign and poisoned inputs. Techniques such as ABL [23] detect poisoned samples by noting their tendency for faster loss reduction during training. EP [50], on the other hand, identifies "backdoor neurons" by observing that their pre-activation distributions differ significantly from those of benign neurons.

However, both categories all have their own limitations. (1) *Clean model-based defenses* are computationally intensive, often requiring hundreds to thousands of training epochs to train the clean model, given the high complexity of self/semi-supervised learning approach. Moreover, they struggle with advanced clean-*label* backdoors [2, 19], where the labels remain intact. In such scenarios, reducing the reliance on labels does not prevent the incorporation of

**Figure 1: Entropy distribution plot for BadNets with CLIP as the weak but clean model. Label-poisoned samples (in red), including both normal and stealthy poisoned samples, are identifiable with high entropy in the clean model, while clean-label poisons (cl-poison in purple) are identifiable with low entropy in the suspicious model.**

backdoor information into the clean one. (2) *Suspicious model-based defenses* are less effective against clean-*image* backdoors [16, 44], where the poisoned images are totally unchanged and share highly similar distributions from clean images. This similarity makes it challenging to detect anomalies based on the data distribution or the activations of the model.

To address these limitations, we proposed a novel defense strategy that leverages both clean and suspicious models from an entropy perspective. Our key insight is that a pretrained CLIP model, as a general zero-shot classifier, can serve as a *weak but clean classifier* across diverse tasks without additional training. By analyzing the cross-entropy score of predictions from CLIP and a suspicious model, we can effectively identify poisoned samples. Specifically, high entropy predictions from CLIP often indicate mislabeled samples, which is characteristic of label poisoning attacks. Conversely, in clean-label backdoors—where poisoned images are correctly labeled—the suspicious model tends to produce low-entropy predictions due to the strong association between the trigger and the target label.

Based on this observation, we introduce **C**LIP-**G**uided backdoor **D**efense (CGD), a two-step backdoor mitigation approach. (1) *CLIP-Guided Meta Data Splitting*: We generate a cross-entropy map (Fig. 1) by evaluating each sample with both CLIP and the suspicious model. Samples with high entropy predictions from CLIP and low entropy predictions from the suspicious model are flagged as potentially poisoned. This allows us to separate the dataset into clean and triggered subsets. (2) *CLIP-Guided Backdoor Unlearning*: Using the separated subsets, we retrain the model on clean data and perform unlearning on triggered data. During unlearning, CLIP guides the logits for the triggered samples, helping to remove the backdoor influence without degrading the model's overall performance.

We conducted extensive experiments to validate the effectiveness of CGD. With a runtime under 3 minutes, CGD can impressively reduce attack success rates (ASRs) to 0.2%, 0.6%, 1.0%, and 0.1% on average across 11 different attack categories for CIFAR-10, CIFAR-100, GTSRB, and Tiny-ImageNet, respectively, outperforming all

existing defenses. Notably, we innovatively apply clean-data-based defenses [27, 42] to poisoned data by filtering a clean subset with CGD. Results show that defenses using poisoned data can even surpass those using 5% clean data. Remarkably, even with a very weak CLIP (7.5% accuracy), our defense can still achieve a 0.7% ASR and only a 0.7% drop in clean accuracy. Additionally, we demonstrate that even when CLIP contains its own backdoor, CGD can still remove the backdoor from the victim model without transferring CLIP's backdoor. Our contributions are threefold:

- **Introduction of CLIP-Guided Defense (CGD)**: We proposed CGD, a novel backdoor defense mechanism that leverages CLIP as a weak but clean classifier, using entropy-based data separation and efficient unlearning to effectively remove backdoors from poisoned datasets.
- **Extensive Experimental Validation**: Our method shows superior performance across multiple datasets and attack scenarios, reducing ASRs to under 1% and remaining robust even when using weak or backdoored CLIP models.
- **Introducing Clean-Data-Based Defenses on Poisoned Data**: We introduce a novel insight that enables defenses traditionally dependent on clean data to operate on poisoned data by leveraging CGD to filter a clean subset. This concept opens up the possibility for clean-data defenses to be adapted for poisoned environments, offering a practical solution when clean data is unavailable.

## 2 Related Work

### 2.1 Multimodal Contrastive Learning (MCL)

Multimodal Contrastive Learning (MCL) bridges modalities by embedding large-scale data into a shared feature space. Specifically, in *image-text* MCL, it simultaneously learns visual and textual representations. CLIP [33], a seminal MCL model, achieves strong generalization by training on a 400M image-text dataset, treating paired images and texts as positives and all others as negatives. This robust cross-modal understanding has inspired advancements such as Uniclip [18], Cyclip [10] and DeCLIP [22]. While most existing works on MCL focus on security concerns of CLIP, such as backdoor attacks [3, 25] and defenses [45], our approach is novel. We leverage CLIP to enhance backdoor defense for other models, a direction not previously explored.

### 2.2 Backdoor Attacks and Defenses

*2.2.1 Backdoor Attacks.* Backdoor attacks are often implemented via data poisoning [2, 5, 11, 32], where adversaries inject poison samples into the training dataset, causing the victim model to associate backdoor triggers with target classes. Beyond data poisoning, backdoor attacks can also manipulate the training process [29, 30, 39], though these can often be adapted into data poisoning by introducing additional noise samples [41]. Some attacks are designed specifically to evade defenses [26, 32].

*2.2.2 Clean-Data-Based Backdoor Defenses.* Existing defenses [27, 42, 46, 49] often assume the availability of a small clean dataset (typically 5% of the training set) to mitigate backdoor attacks [41]. Fine-Pruning (FP) [27] combines pruning and fine-tuning to remove backdoor-related neurons. Mode Connectivity Repair (MCR)
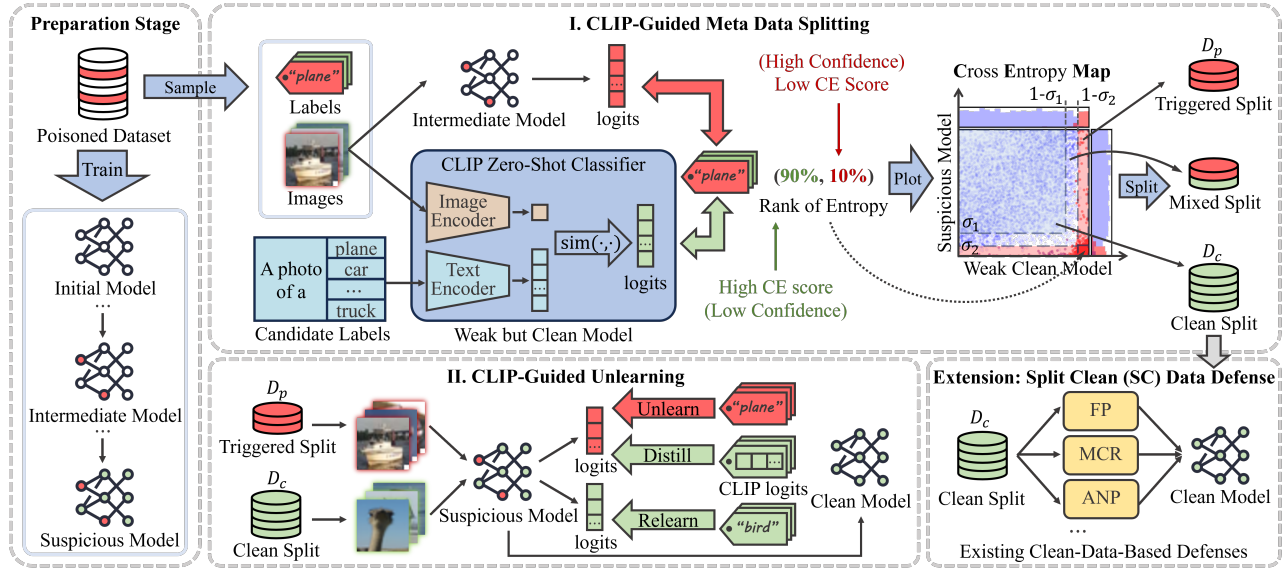
**Figure 2: Pipeline for our CLIP-Guided Backdoor Defense (CGD).**

[49] navigates toward a clean model using clean data. Adversarial Neuron Pruning (ANP) [42] prunes neurons sensitive to triggers by adversarial training. Although effective, these methods require clean data, which may not always be available.

*2.2.3* ***Poison-Data-Based Backdoor Defenses.*** Poison-data-based defenses [4, 8, 14, 23, 50] handle scenarios without clean data. (1) One approach focuses on representation learning. For instance, DBD [14] uses self-supervised learning to remove poisoned labels during pretraining, severing the link between the trigger and target label. D-BR [4] and ASD [8] further refine this by combining representation learning with semi-supervised approaches. However, these methods struggle against clean-label attacks, where triggers are embedded directly without modifying labels, allowing the victim model to still learn the trigger. (2) Another approach leverages the unique properties of triggered samples. Anti-Backdoor Learning (ABL) [23] reduces the influence of poisoned samples by down-weighting those with higher training losses. Entropy Pruning (EP) [50] distinguishes backdoor neurons through unique bimodal pre-activation patterns. However, these methods are less effective against clean-image backdoors, where the images remain visually indistinguishable from clean ones.

## 3 Preliminary

*3.0.1* ***Notations.*** We consider a classification model $f(\cdot; \theta)$ parameterized by $\theta$. Let $p(y \mid x; \theta)$ denote the predicted probability assigned by the model $f(\cdot; \theta)$ to class $y$ given input $x$. The predicted label for $x$ is then defined as

$$\hat{y}(x) = \arg\max_{y \in \mathcal{Y}} p(y \mid x; \theta),$$

where $\mathcal{Y} = \{1, 2, \ldots, C\}$ is the set of possible class labels. We define the trigger planting function, used to execute backdoor attacks, as $\mathcal{T} : \mathcal{X} \to \mathcal{X}$, and let $t \in \mathcal{Y}$ represent the attack's designated target class. $\mathcal{P}$ denotes the distribution of clean, unaltered samples. The clean accuracy (CA) of the model is the probability of correct



**Figure 3: Percentile rank of entropy for triggered images from intermediate models in clean-label backdoors (LC [38], SIG [2], CTRL [19]). The 50% distribution range is marked with an area. Triggered images exhibit low entropy during training, enabling their identification.**

classification on clean data, $\mathbb{P}_{(x,y)\sim\mathcal{P}}[\hat{y}(x) = y]$. The attack success rate (ASR), $\mathbb{P}_{(x,y)\sim\mathcal{P}|y\neq t}[\hat{y}(\mathcal{T}(x)) = t]$, measures the probability that a sample $x$ with an embedded trigger is misclassified as the target class $t$.

*3.0.2* ***Threat model.*** We adopt the poisoning-based threat model used in previous works [5, 11, 38], where attackers provide a poisoned training dataset containing a set of pre-created triggered samples. The defender neither knows the potential backdoor trigger pattern or even whether the dataset is poisoned. Moreover, the defender is believed to have access to a publicly available vision-languages model, such as CLIP [33]. The goal of defenders is to obtain a well-performed model without suffering backdoor attacks.

## 4 Methodology

In this section, we introduce the pipeline of our CLIP-Guided Backdoor Defense (CGD) framework. As illustrated in Fig. 2, CGD comprises two main stages: the *CLIP-Guided Meta Data Splitting* stage and the *CLIP-Guided Unlearning* stage. Before deploying these stages, we assume the model owner has normally trained a suspicious model using a suspicious dataset that is potentially poisoned.

## 4.1 CLIP-Guided Meta Data Splitting

The novelty of our approach stems from leveraging CLIP, a pre-trained vision-language model, as a zero-shot classifier to detect backdoor samples in a poisoned dataset. Unlike traditional backdoor defenses that often rely on computationally intensive self-supervised or semi-supervised learning—approaches that struggle with certain backdoor types—our method exploits CLIP's cross-modal understanding to identify subtle inconsistencies indicative of poisoning. This section outlines how we harness CLIP's capabilities to address both poison-label and clean-label backdoors, offering a robust and efficient alternative to existing methods.

### 4.1.1 *Poison-Label Backdoor Identification*.
Poison-label backdoors involve poisoned samples with incorrect or misleading labels, encompassing attacks such as classic backdoors [5, 11], dynamic backdoors [29, 30, 39], and clean-image backdoors [16, 44]. These attacks share a vulnerability: mislabeling that deviates from the true semantic content of the samples. We employ CLIP as a zero-shot classifier to detect such anomalies by computing entropy scores that reveal label inconsistencies.

For each sample $x_i$, CLIP generates a logit vector $\mathbf{z}_i^{\mathrm{CLIP}} \in \mathbb{R}^{|L|}$ over a predefined label set $L$, without requiring additional training. We calculate the cross-entropy loss between CLIP's predictions and the sample's suspicious label $y_i$, yielding an entropy score $\mathcal{S}_i^{\mathrm{CLIP}}$:

$$\mathcal{S}_i^{\mathrm{CLIP}} = -\log p(y_i \mid x_i; \theta_{\mathrm{CLIP}}),$$

where $p(y_i \mid x_i; \theta_{\mathrm{CLIP}})$ is the probability CLIP assigns to the suspicious label $y_i$. A high $\mathcal{S}_i^{\mathrm{CLIP}}$ suggests a mismatch between the image content and its label, flagging the sample as potentially poisoned. This approach leverages CLIP's pre-trained alignment of visual and textual features, providing a novel, training-free mechanism for backdoor detection.

### 4.1.2 *Clean-Label Backdoor Identification*.
Clean-label backdoors poison images without altering their original labels, making detection challenging. Here, triggered samples exhibit lower cross-entropy loss because the trigger acts as a dominant feature that the model learns to associate with the correct label during training. This heightened confidence reduces the loss, as illustrated in Fig. 3, where loss curves for attacks like LC [38], SIG [2], and CTRL [19] drop and stabilize at low values. We exploit this property by thresholding the entropy scores from the suspicious model at epoch $T = 5$, denoted as $\mathcal{S}_i^{\mathrm{Model}}$:

$$\mathcal{S}_i^{\mathrm{Model}} = -\log p(y_i \mid x_i; \theta_{\mathrm{Model}}),$$

where $p(y_i \mid x_i; \theta_{\mathrm{Model}})$ is the model's predicted probability for the correct label. Samples with unusually low $\mathcal{S}_i^{\mathrm{Model}}$ are flagged as potential clean-label backdoors, capitalizing on the model's over-confidence induced by triggers.

### 4.1.3 *Entropy Map Splitting*.
To ensure robustness across diverse attacks and datasets, we use percentile ranks of entropy scores, denoted by the function $\hat{\cdot}$, rather than raw values. For each sample $x_i$, we compute $\hat{\mathcal{S}}_i^{\mathrm{CLIP}}$ and $\hat{\mathcal{S}}_i^{\mathrm{Model}}$ from CLIP and the model, respectively. We then apply a threshold-based strategy with $\sigma_1$ and $\sigma_2$ to define the clean subset $D_c$ and triggered subset $D_p$:

$$D_c = \left\{ x_i \,\middle|\, \hat{\mathcal{S}}_i^{\mathrm{CLIP}} \leq 1 - \sigma_1 \text{ and } \hat{\mathcal{S}}_i^{\mathrm{Model}} \geq \sigma_1 \right\},$$

$$D_p = \left\{ x_i \,\middle|\, \hat{\mathcal{S}}_i^{\mathrm{CLIP}} > 1 - \sigma_2 \text{ or } \hat{\mathcal{S}}_i^{\mathrm{Model}} < \sigma_2 \right\}.$$

The remaining samples form a mixed subset. To mitigate class imbalance in $D_c$, we oversample each class to match the original dataset's label distribution.

## 4.2 CLIP-Guided Unlearning

Drawing on prior unlearning techniques [4, 8, 23], we fine-tune the backdoored model using a novel CLIP-guided approach. This stage integrates relearning, unlearning, and distillation to eliminate backdoor effects while preserving clean-data performance.

### 4.2.1 *Relearning and Unlearning Losses*.
We retrain the model on $D_c$ with the standard cross-entropy loss $\mathcal{L}_{\mathrm{re}}$ and apply a negative cross-entropy loss $\mathcal{L}_{\mathrm{un}}$ on $D_p$:

$$\mathcal{L}_{\mathrm{re}} = -\sum_{x_i \in D_c} y_i \log p_i, \quad \mathcal{L}_{\mathrm{un}} = \sum_{x_i \in D_p} y_i \log(p_i + \epsilon),$$

where $y_i$ is the true label, $p_i$ is the predicted probability, and $\epsilon$ prevents numerical instability.

### 4.2.2 *CLIP-Guided Neural Distillation*.
To guide the model toward clean patterns, we introduce a distillation loss on $D_p$:

$$\mathcal{L}_{\mathrm{distill}} = \sum_{x_i \in D_p} \mathrm{KL}\left(\mathbf{z}_i^{\mathrm{CLIP}} \,\middle\|\, \mathbf{z}_i^{\mathrm{Model}}\right),$$

where $\mathrm{KL}(\cdot\|\cdot)$ is the Kullback-Leibler divergence between CLIP's logits $\mathbf{z}_i^{\mathrm{CLIP}}$ and the model's logits $\mathbf{z}_i^{\mathrm{Model}}$.

### 4.2.3 *Justification of Loss Terms*.
The three loss terms are meticulously designed to address distinct aspects of backdoor mitigation, with their necessity validated by ablation studies (see Section 5): **Relearning Loss $\mathcal{L}_{\mathrm{re}}$**: This term reinforces the model's ability to classify clean samples accurately, anchoring its performance on legitimate data. Without it, fine-tuning risks degrading generalization, as the model may overfit to the unlearning process. **Unlearning Loss $\mathcal{L}_{\mathrm{un}}$**: By penalizing confident predictions on triggered samples, this term disrupts the trigger-label association critical to backdoor attacks. Its absence would leave the backdoor intact, as the model retains its poisoned behavior. **Distillation Loss $\mathcal{L}_{\mathrm{distill}}$**: This term provides a positive learning signal, aligning the model's predictions with CLIP's clean, zero-shot classifications. Omitting it risks leaving the model without a coherent strategy for triggered samples, potentially reducing robustness.

Together, these terms form a synergistic framework: $\mathcal{L}_{\mathrm{re}}$ preserves clean performance, $\mathcal{L}_{\mathrm{un}}$ erases backdoor effects, and $\mathcal{L}_{\mathrm{distill}}$ ensures a smooth transition to clean behavior, making the design principled rather than heuristic.

### 4.2.4 *Final Loss Function*.
The total loss combines these terms:

$$\mathcal{L}_{\mathrm{total}} = \mathcal{L}_{\mathrm{re}} + \lambda_{\mathrm{un}} \mathcal{L}_{\mathrm{un}} + \lambda_{\mathrm{distill}} \mathcal{L}_{\mathrm{distill}},$$

where $\lambda_{\mathrm{un}}$ and $\lambda_{\mathrm{distill}}$ balance their contributions. Fine-tuning is limited to $K$ epochs, with early stopping if clean accuracy falls below $\tau$, ensuring performance preservation.

**Table 1: Poison data based defensive results on CIFAR-10 (CA and ASR) under poison rate of 5%.**

| Defense → Attack ↓ | No Defense | | ABL [23] | | DBD [14] | | ASD [8] | | D-BR [4] | | EP [50] | | ReBack [28] | | PIPD [6] | | MSPC [31] | | CGD (ours) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| Classic Backdoor — BadNets [11] | 91.8 | 93.8 | 90.8 | 1.1 | 92.1 | 2.6 | 92.0 | 2.1 | 91.0 | 1.5 | 91.0 | 0.8 | 91.7 | 4.3 | 92.9 | 0.5 | 92.8 | 0.3 | 92.8 | 0.0 |
| Classic Backdoor — Blend [5] | 93.7 | 99.8 | 86.2 | 4.4 | 93.0 | 100.0 | 93.0 | 5.3 | 85.1 | 0.0 | 91.6 | 96.1 | 91.7 | 2.4 | 93.1 | 5.3 | 92.7 | 0.7 | 93.2 | 0.0 |
| Dynamic Backdoor — WaNet [30] | 90.6 | 96.9 | 89.9 | 81.2 | 90.5 | 2.6 | 91.7 | 8.8 | 84.3 | 60.2 | 89.8 | 91.1 | 90.2 | 84.4 | 93.4 | 11.4 | 93.0 | 54.2 | 94.0 | 0.3 |
| Dynamic Backdoor — BPP [39] | 91.4 | 99.2 | 89.1 | 99.8 | 92.4 | 99.9 | 92.5 | 99.4 | 88.5 | 85.5 | 89.8 | 4.6 | 90.1 | 1.8 | 93.1 | 0.9 | 90.5 | 2.8 | 94.0 | 0.3 |
| Dynamic Backdoor — IAB [29] | 89.7 | 94.9 | 89.3 | 83.0 | 91.3 | 0.0 | 92.3 | 19.8 | 85.3 | 84.8 | 90.7 | 1.6 | 87.9 | 1.7 | 92.0 | 4.0 | 92.5 | 5.3 | 93.3 | 0.7 |
| Dynamic Backdoor — SSBA [21] | 93.0 | 97.3 | 88.6 | 3.9 | 92.5 | 2.4 | 93.3 | 7.1 | 83.1 | 3.0 | 92.0 | 10.5 | 85.1 | 6.6 | 90.6 | 17.2 | 90.9 | 21.5 | 92.9 | 0.0 |
| Clean-Label Backdoor — CTRL [19] | 93.6 | 95.9 | 88.2 | 2.4 | 92.1 | 57.8 | 91.3 | 89.3 | 90.5 | 98.3 | 92.3 | 1.1 | 91.4 | 96.2 | 90.1 | 12.6 | 91.8 | 77.3 | 93.6 | 0.1 |
| Clean-Label Backdoor — SIG [2] | 93.6 | 93.9 | 88.2 | 0.4 | 89.2 | 97.5 | 92.2 | 99.5 | 91.3 | 49.6 | 92.1 | 12.5 | 87.4 | 29.9 | 93.2 | 13.5 | 91.0 | 10.3 | 93.3 | 0.0 |
| Clean-Label Backdoor — LC [38] | 93.4 | 98.4 | 82.1 | 5.2 | 92.3 | 98.3 | 91.2 | 9.8 | 91.3 | 1.7 | 92.5 | 0.6 | 91.4 | 1.0 | 90.8 | 3.7 | 90.8 | 89.4 | 93.2 | 0.0 |
| Clean-Image Backdoor — FLIP [16] | 89.9 | 99.2 | 84.8 | 99.6 | 92.3 | 1.6 | 86.9 | 62.2 | 83.9 | 22.1 | 89.9 | 98.5 | 90.0 | 39.7 | 91.4 | 66.9 | 91.6 | 17.2 | 92.5 | 0.5 |
| Clean-Image Backdoor — GCB [44] | 88.6 | 100.0 | 85.5 | 100.0 | 91.2 | 6.9 | 90.9 | 100.0 | 84.2 | 100.0 | 88.3 | 99.9 | 88.7 | 71.6 | 92.5 | 87.7 | 91.5 | 23.9 | 92.3 | 0.0 |
| Average | 91.7 | 97.2 | 87.5 | 43.7 | 91.7 | 42.7 | 91.6 | 45.7 | 87.2 | 46.1 | 90.9 | 37.9 | 89.6 | 30.7 | 92.1 | 20.3 | 91.7 | 27.5 | 93.2 | 0.2 |
| CA Drop (smaller is better) | | | ↓4.2 | | ↓0.0 | | ↓0.2 | | ↓4.6 | | ↓0.9 | | ↓2.1 | | ↓-0.4 | | ↓0.0 | | ↓-1.5 | | |
| ASR Drop (larger is better) | | | ↓53.5 | | ↓54.5 | | ↓51.5 | | ↓51.1 | | ↓59.3 | | ↓66.5 | | ↓76.9 | | ↓69.7 | | | | ↓97.0 |

## 4.3 Extension: Enabling Clean-Data-Based Defenses on Poisoned Data

We extend our framework by using the split clean subset $D_c$ to enable clean-data-based defenses like Fine-Pruning (FP) [27], Mode Connectivity Repair (MCR) [49], and Adversarial Neuron Pruning (ANP) [42] on poisoned datasets. Traditionally reliant on separate clean data, these methods gain practical utility through our approach, denoted as ANP-SC, MCR-SC, etc., enhancing their applicability and showcasing the versatility of our splitting strategy.

## 5 Evaluation

### 5.1 Experimental Setups

*5.1.1 Attack Baselines and Setups.* We use four benchmark datasets: CIFAR-10 [17], CIFAR-100 [17], GTSRB [36], and a subset of ImageNet [34]. PreActResNet18 [13] serves as the default model. We evaluate eleven state-of-the-art backdoor attacks across four categories: (1) *Classic static-backdoor attacks*: BadNets [11] and Blend [5]; (2) *Dynamic-backdoor attacks*: SSBA [21], IAB [29], WaNet [30], and BPP [39]; (3) *Clean-label backdoor attacks*: LC [38], SIG [2], and CTRL [19]; (4) *Clean-image backdoor attacks*: FLIP [16] and GCB [44]. Attacks are implemented following guidelines from [40]. We set the target label to 0 ($y_t = 0$) and use a poison rate of 50% for clean-label attacks and 5% for the others.

*5.1.2 Defenses.* We compare CGD with eight state-of-the-art established defenses: (1) *Poisoned-data-based defenses*: ABL [23], DBD [14], ASD [8], D-BR [4], EP [50], ReBack [28], PIPD [6], and MSPC [31]. D-BR, EP, ReBack, PIPD and MSPC use both the poisoned data and backdoored model, while ABL, DBD, and ASD only need poisoned data. These defenses can be integrated during normal training of the backdoored model on the poisoned dataset (2) *Clean-data-based defenses*: FP [27], MCR [49], I-BAU [46], and ANP [42]. These methods rely on access to a clean subset, typically around 5% of the training data, as a baseline setting. We also evaluate these methods using clean data obtained from our CLIP-guided Splitting (SC) approach. For all methods, we use implementations from BackdoorBench whenever available; otherwise, we rely on their official implementations.

**Table 2: Average training time (s) for poison-data-based defenses on CIFAR-10. *Prepare* refers to initial victim model training on poisoned data. Our method needs less than 3 minutes post victim model training.**

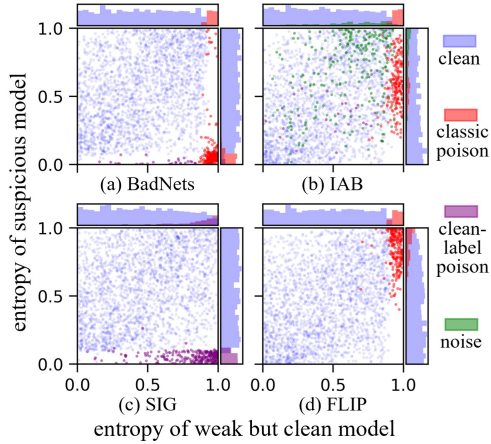| Method | (Prepare) | Data Split | Defense | Total |
|---|---|---|---|---|
| ABL [23] | - | 437 | 1,729 | 2,166 |
| DBD [14] | - | 17,672 | 9,299 | 26,971 |
| ASD [8] | - | 4,702 | 2,591 | 7,293 |
| D-BR [4] | 997 | 656 | 3,900 | 5,553 |
| EP [50] | 997 | - | 53 | 1,050 |
| ReBack [28] | 997 | 112 | 92 | 1,201 |
| PIPD [6] | 997 | 1,325 | 438 | 2,760 |
| MSPC [31] | 997 | 1,844 | 352 | 3,193 |
| **CGD (ours)** | 997 | 60 | 98 | 1,154 |

*5.1.3 CGD Configuration.* For CGD, we use CLIP [33] with ViT-B-32, pretrained on Laion-2B [35] via Open-CLIP. We set a **uniform threshold** $\sigma_1 = 0.1$ and $\sigma_2 = 0.2$ for triggered and clean data on **all** the experiments in evaluation except in specific ablation. Trade-off parameters are set as $\lambda_{un} = 0.025$ and $\lambda_{distill} = 0.0005$. We use $T = 5$ as the intermediate epoch, and SGD with a learning rate of 0.01 and weight decay of $5 \times 10^{-4}$ for fine-tuning.

*5.1.4 Evaluation Metrics.* We evaluate defenses using clean accuracy (CA) and attack success rate (ASR). Specifically, ASR measures the fraction of non-target samples with triggers classified into the target label. An effective defense should maintain high CA while minimizing ASR, ensuring robustness against backdoor attacks without sacrificing model performance.

### 5.2 Main Results

*5.2.1 CLIP-Guided Defense (CGD).* To evaluate the effectiveness of CGD, we report the CA and ASR results for eight defense methods against eleven attacks on CIFAR-10 in Table 1. CGD consistently achieves low ASRs ($\leq 1\%$) while maintaining high CAs, with minimal drops ($\leq 0.3\%$) across all defenses. Remarkably, CGD can identify and unlearn poisoned samples, improving CA over the "No Defense" baseline by 3.4% for WaNet and 3.6% for IAB. In contrast, self-supervised defenses like DBD, ASD, and D-BR fail to resist clean-label attacks such as CTRL and SIG, while property-based

**Figure 4: Entropy distribution plot for various categories of attacks. Poisoned samples (red/purple) are easily separated.**

**Table 3: Clean-data-based defenses. Methods with * use 5% clean data, while those with the "-SC" suffix use Split Clean data extracted from poisoned datasets using CGD.**

| Defense → | FP* [27] | | MCR* [49] | | I-BAU* [46] | | ANP* [42] | |
|---|---|---|---|---|---|---|---|---|
| Attack ↓ | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| Classic | 92.1 | 21.7 | 92.6 | 51.1 | 86.4 | 16.5 | 85.3 | 21.4 |
| Dynamic | 92.8 | 30.9 | 72.6 | 47.9 | 90.5 | 30.3 | 87.2 | 1.1 |
| Clean-Label | 92.4 | 58.0 | 93.3 | 82.5 | 88.7 | 16.0 | 86.5 | 24.9 |
| Clean-Image | 92.1 | 17.7 | 92.9 | 38.5 | 89.4 | 6.6 | 84.9 | 0.3 |
| Average | 92.4 | 32.1 | 87.9 | 55.0 | 88.7 | 17.3 | 86.0 | 11.9 |
| Defense → | FP-SC | | MCR-SC | | I-BAU-SC | | ANP-SC | |
| Attack ↓ | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| Classic | 93.3 | 18.2 | 91.0 | 1.4 | 90.7 | 25.7 | 86.6 | 18.8 |
| Dynamic | 94.0 | 25.2 | 87.6 | 24.3 | 91.8 | 6.2 | 85.7 | 0.6 |
| Clean-Label | 93.2 | 57.2 | 90.8 | 5.7 | 90.8 | 17.8 | 87.5 | 9.7 |
| Clean-Image | 92.7 | 15.8 | 91.3 | 1.0 | 90.4 | 43.4 | 85.8 | 0.1 |
| Average | 93.3 | 29.1 | 90.2 | 8.1 | 90.9 | 23.3 | 86.4 | 7.3 |

defenses like ABL and EP are vulnerable to clean-image backdoors like GCB and FLIP. Additionally, CGD is highly efficient, requiring less than 3 minutes to split data and apply defenses after training the backdoored model (Table 2). While EP [50] takes a similar amount of time, it cannot separate clean data. Thus, CGD stands out as the most efficient method for isolating clean data from poisoned samples while maintaining a high defense success rate.

*5.2.2 Entropy Map Visualization.* The entropy map distributions for the four defense categories, shown in Fig. 4 , highlight distinct patterns. (a) *Classical backdoors* (e.g., BadNets [11]): Poisoned samples cluster in the lower right, separable by entropy from suspicious models or CLIP. (b) *Dynamic backdoors* (e.g., IAB [29]): Noise (green points) complicates detection but doesn't inherently form backdoors. (c) *Clean-label backdoors* (e.g., SIG [2]): These resemble classical backdoors in suspicious model entropy, detectable by methods like ABL [23]. (d) *Clean-image backdoors* (e.g., FLIP [16]): High entropy in suspicious models reduces the effectiveness of entropy- or loss-based defenses. Overall, each category's unique entropy patterns support using category-average performance as a key metric in later experiments.

**Table 4: CGD scalability on other datasets.**

| Dataset → | CIFAR-100 (No Defense) | | GTSRB (No Defense) | | ImageNet (No Defense) | |
|---|---|---|---|---|---|---|
| Attack ↓ | CA | ASR | CA | ASR | CA | ASR |
| Classic | 69.6 | 85.3 | 97.9 | 96.7 | 56.6 | 99.1 |
| Dynamic | 66.2 | 94.6 | 97.0 | 95.2 | 57.5 | 98.7 |
| Clean-Label | 70.6 | 38.3 | 98.4 | 85.5 | 51.6 | 48.1 |
| Clean-Image | 67.1 | 99.8 | 95.3 | 100.0 | 54.0 | 99.8 |
| Average | 68.4 | 79.5 | 97.1 | 94.4 | 54.9 | 86.4 |
| Dataset → | CIFAR-100 (CGD) | | GTSRB (CGD) | | ImageNet (CGD) | |
| Attack ↓ | CA | ASR | CA | ASR | CA | ASR |
| Classic | 69.6 | 0.0 | 97.8 | 0.0 | 56.5 | 0.0 |
| Dynamic | 70.4 | 0.3 | 97.8 | 1.1 | 57.7 | 0.0 |
| Clean-Label | 69.6 | 2.1 | 94.7 | 2.8 | 48.7 | 0.3 |
| Clean-Image | 68.2 | 0.0 | 96.8 | 0.0 | 55.6 | 0.0 |
| Average | 69.4 | 0.6 | 96.8 | 1.0 | 54.6 | 0.1 |
| **Average Drop** | **-1.0** | **78.9** | **0.4** | **93.4** | **0.3** | **86.4** |

*5.2.3 Split Clean-Data Defense.* We incorporate a portion of selected clean data into clean-data-based defenses and observe a significant improvement in defense performance. Results comparing standard clean-data defenses with their counterparts using CLIP-guided data splitting are presented in Table 3. As shown, for most methods tested (FP [27], MCR [49], ANP [42]), our Split Clean (SC) approach achieves superior outcomes, with both lower ASRs and higher CAs. This improvement is likely because our default hyperparameter setting allocates around 70% of the total data as clean, a significantly larger proportion than the typical clean-data assumption (around 5% of total data). Consequently, defenses such as ANP and MCR, which are based on the sensitivity of poisoned data, benefit from increased clean data in effectively mitigating backdoor effects. These results suggest that, for certain defenses, access to only 10% of poisoned data can yield defense results comparable to using 5% clean data, underscoring the effectiveness of the Split Clean approach.

*5.2.4 Scalability.* Table 4 demonstrates the effectiveness of CGD across three additional datasets: CIFAR-100, GTSRB, and ImageNet. For each dataset, we evaluate all successful backdoor attacks (ASR ≥ 90%) and report the average results across attack categories. In all cases, CGD effectively eliminates backdoors from the models (with ASR reduced to ≤ 3%) while incurring only minimal clean accuracy (CA) reduction (average CA drop ≤ 0.3%). These results confirm the scalability and robustness of CGD across diverse datasets.

## 5.3 Ablation Study

*5.3.1 Loss Term Analysis.* In our CLIP-guided unlearning approach, we utilize three loss terms: unlearn loss ($\mathcal{L}_{\text{un}}$), relearn loss ($\mathcal{L}_{\text{re}}$), and distillation loss ($\mathcal{L}_{\text{distill}}$). To evaluate their roles, we conducted an ablation study by testing various combinations of these terms, with results shown in Table 5, measuring CA and ASR across multiple backdoor categories.

❶ **Single Loss Term.** Using only $\mathcal{L}_{\text{un}}$ yields high CA (average: 88.1) but fails to reduce ASR (average: 93.8), notably against clean-label backdoors (ASR: 98.5), indicating insufficient backdoor

**Table 5: Ablation study for different loss terms in CGD.**

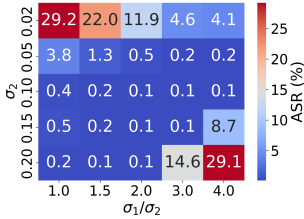| Case | | | Classic | | Dynamic | | Clean-label | | Clean-Image | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\mathcal{L}_{un}$ | $\mathcal{L}_{re}$ | $\mathcal{L}_{distill}$ | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| ✓ | | | 83.6 | 97.5 | 90.2 | 82.0 | 91.7 | 98.5 | 87.1 | 97.3 | 88.1 | 93.8 |
| | ✓ | | **93.2** | 91.5 | **93.6** | 62.2 | **93.6** | 94.4 | 91.2 | 97.2 | 92.9 | 86.3 |
| | | ✓ | 26.8 | 13.8 | 51.8 | 64.5 | 34.0 | 29.3 | 24.0 | 25.7 | 34.1 | 33.3 |
| ✓ | ✓ | | 93.1 | 48.3 | 93.5 | 34.5 | **93.6** | 93.8 | 92.0 | 37.2 | 93.0 | 53.5 |
| ✓ | | ✓ | 9.3 | 1.1 | 36.8 | 58.9 | 13.5 | 14.3 | 9.7 | 14.3 | 17.3 | 22.1 |
| | ✓ | ✓ | 92.8 | 4.1 | **93.6** | 24.5 | 92.9 | 11.8 | 92.2 | 13.7 | 92.9 | 13.5 |
| ✓ | ✓ | ✓ | 93.0 | **0.0** | 93.5 | **0.3** | 93.4 | **0.0** | 92.4 | **0.2** | **93.1** | **0.1** |

mitigation without targeted guidance. Similarly, $\mathcal{L}_{re}$ alone maintains high CA (average: 92.9) yet struggles with high ASR (average: 86.3), as backdoor knowledge persists. Employing only $\mathcal{L}_{distill}$ reduces ASR (average: 33.3) but severely degrades CA (average: 34.1), making the model impractical.

❷ **Two-Term Combinations.** Combining $\mathcal{L}_{un}$ and $\mathcal{L}_{re}$ improves defense (average ASR: 53.5), yet clean-label backdoors remain potent (ASR: 93.8), highlighting the need for additional guidance. Pairing $\mathcal{L}_{re}$ and $\mathcal{L}_{distill}$ reduces ASR significantly (average: 13.5), but dynamic backdoors retain a notable ASR (24.5), showing incomplete protection against sophisticated attacks.

❸ **Synergy of All Three Terms.** Employing all three losses achieves optimal results, with an average CA of 93.1 and an ASR of 0.1 across all backdoor types. Here, $\mathcal{L}_{un}$ targets backdoor unlearning, $\mathcal{L}_{distill}$ aligns the model with correct logits via CLIP's knowledge, and $\mathcal{L}_{re}$ preserves clean performance.

*5.3.2* **Hyperparameter Study.**
The primary hyperparameters in our experiment are the clean data threshold $\sigma_1$ and the triggered data threshold $\sigma_2$, which define the splitting of the entropy map. Selecting suitable values for these thresholds is crucial: if set too low, triggered samples may be misclassified as clean, while overly high settings risk assigning clean samples to the triggered subset. Both cases reduce performance. To maintain $\sigma_1$ consistently higher than $\sigma_2$, we optimize using their ratio $\sigma_1/\sigma_2$. As illustrated in Fig. 5, ASR averaged over 11 attack methods highlights a robust range where ASR remains low. Even in worst-case scenarios, ASR is reduced below 30%. For default settings, we use $\sigma_1 = 0.2$ and $\sigma_2 = 0.1$, though these values can be flexibly adjusted within a reasonable range.
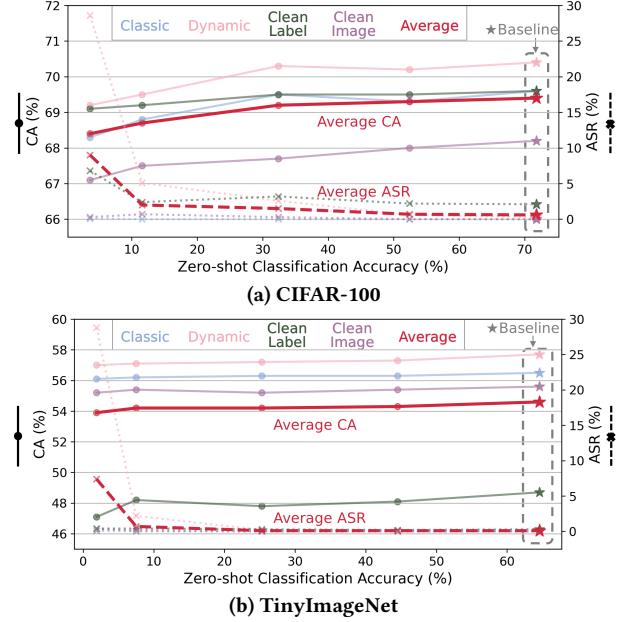


**Figure 5: Threshold Study.**

*5.3.3* **Robustness to Different Language-Image Pretraining Models.** Although we mainly use CLIP [33] to illustrate CGD's effectiveness, CGD can use other language-image pretraining models to replace CLIP. As shown in Table 6, we consider SigLIP [47] and ImageBind [9] as two alternative models. The result suggests that both achieve low ASR ($\leq$ 5%) while maintaining good CA ($\geq$ 90%).

**Table 6: CGD using other language-image pretraining models except CLIP.**

| Model→ | SigLIP | | ImageBind | |
|---|---|---|---|---|
| Attack↓ | CA | ASR | CA | ASR |
| Classic | 93.0 | 0.0 | 92.6 | 0.0 |
| Dynamic | 93.2 | 5.3 | 93.3 | 2.5 |
| C-Label | 92.8 | 0.7 | 92.7 | 0.7 |
| C-Image | 90.8 | 3.3 | 90.6 | 2.8 |
| Average | 92.5 | 2.3 | 92.3 | 1.5 |



**(a) CIFAR-100**



**(b) TinyImageNet**

**Figure 6: CGD's defense performance on low-accuracy CLIP models. CGD successfully defends against all attacks (ASR $\leq$ 5%) even when CLIP's zero-shot accuracy is as low as 10%.**

### 5.4 Security Guarantee Analysis

Since our CGD is mainly based on the assumption of CLIP *can* act as an weak but clean model. As a result, there are two major security risks here: (A) *If CLIP is a valid weak model for the specific task*, and (B) *If CLIP is a clean model.* We evaluate these two assumptions one by one.

*5.4.1* **Weaker CLIP.** We further investigate our defense's robustness when CLIP's accuracy is intentionally reduced. Given that current versions of CLIP are highly effective on common datasets, we manually introduce random noise to CLIP's output, reducing its accuracy in a controlled manner. Specifically, we add Gumbel-distributed noise [15] to the output logits of CLIP, enabling us to adjust the level of inaccuracy by varying the scale parameter of the Gumbel distribution. As shown in Fig. 6, we present results illustrating the defense performance against the zero-shot classification accuracy of CLIP on CIFAR-100 and TinyImageNet. Our CGD approach successfully mitigates all tested attacks (ASR $\leq$ 5%) when CLIP's zero-shot classification accuracy remains as low as 11.6% on CIFAR-100 and 7.5% on TinyImageNet. This is likely because even a weakened CLIP predicts images as similar classes close to the target label, causing the entropy-based calculations from CLIP to become less precise, though still somewhat informative. Additionally, while this noise does not impact backdoor removal capability, it does reduce clean accuracy in tandem with the drop in CLIP accuracy.

To further validate this assumption in real and diverse vision datasets, we carry out very extensive experiments on 17 different datasets, including large-scale datasets like SUN397 and full ImageNet. On these datasets, we tested three representative attacks (Blend, BPP, CTRL) but excluded clean-image backdoors due to either low ASR on large datasets (FLIP) or high computational cost (GCB). As shown in Table 7, CGD successfully defends **all attacks**

**Table 7: CGD on 17 datasets against 3 backdoor attacks. "ΔCA": change in CA after defenses. CA drops ≤4.2%. CGD fails only when CLIP is thoroughly ineffective (e.g., on SVHN) on dynamic backdoor attack.**

| Task→ Dataset↓ | Data Size | Classify Accuracy | | CGD (Ours) | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Blend | | BPP | | CTRL | |
| | | CLIP | RN50 | ΔCA | ASR | ΔCA | ASR | ΔCA | ASR |
| SVHN | 234M | 13.4 | 95.7 | 0.5 | 0.1 | -1.4 | 100 | -4.2 | 8.6 |
| Country211 | 20.9G | 17.2 | 12.2 | 0.9 | 0.1 | 0.6 | 0.0 | 0.4 | 9.7 |
| GTSRB | 689M | 32.6 | 98.5 | -1.0 | 0.0 | 2.2 | 0.7 | -3.8 | 0.2 |
| FER2013 | 950M | 41.4 | 62.1 | -0.4 | 0.0 | 3.1 | 0.1 | -0.4 | 0.0 |
| DTD | 1.2G | 44.3 | 71.8 | 1.2 | 0.0 | 3.2 | 1.1 | -2.2 | 0.7 |
| MNIST | 127M | 48.2 | 99.6 | -2.3 | 7.8 | 0.2 | 0.0 | 0.0 | 0.0 |
| RenderedSST2 | 305M | 58.6 | 53.8 | 2.8 | 0.0 | 4.8 | 0.5 | -1.8 | 0.0 |
| StanfordCars | 3.7G | 59.6 | 52.7 | 0.6 | 0.0 | 7.0 | 0.0 | 2.8 | 0.5 |
| SUN397 | 73.6G | 62.5 | 61.5 | -0.4 | 0.0 | 2.0 | 0.0 | 0.2 | 5.5 |
| ImageNet | 146M | 63.3 | 74.3 | -0.3 | 0.0 | 0.1 | 0.0 | -3.1 | 0.5 |
| CIFAR100 | 338M | 64.2 | 73.6 | -0.4 | 0.0 | 4.9 | 0.1 | -2.8 | 5.8 |
| Flowers102 | 676M | 66.5 | 89.5 | -1.3 | 0.2 | 4.2 | 0.8 | -1.3 | 0.4 |
| Caltech101 | 324M | 81.6 | 92.0 | -2.9 | 3.6 | 5.5 | 3.0 | -3.3 | 0.0 |
| OxfordIIITPet | 1.6G | 87.3 | 92.9 | 0.8 | 0.3 | 8.3 | 12.4 | 1.9 | 0.0 |
| Food101 | 5.3G | 88.8 | 72.6 | 7.2 | 0.0 | 2.4 | 0.0 | 0.1 | 8.3 |
| CIFAR10 | 340M | 89.8 | 95.4 | -0.5 | 0.0 | 2.6 | 0.3 | 0.0 | 0.1 |
| STL10 | 5.4G | 97.1 | 97.5 | 3.0 | 0.0 | 4.9 | 0.2 | 2.1 | 0.0 |

except on SVHN for BPP. This failure is because CLIP's zero-shot accuracy on SVHN is only 13.4%, making it thoroughly ineffective for a 10-class prediction dataset. Nonetheless, CGD still defends other SVHN attacks (e.g., Blend, CTRL) via the suspicious model's entropy. Regardless of whether the dataset is large-scale or has low CLIP accuracy (e.g., GTSRB, DTD, and MNIST), CGD remains effective, showing its strong generalizability.

#### 5.4.2 Backdoored CLIP Analysis.
Previous research has shown that CLIP models can be vulnerable to backdoor attacks [3]. In this study, we investigate whether a backdoored CLIP model could pass its backdoor to a victim model during our defense process. We explore three cases: (i) *CLIP has a different backdoor trigger*; (ii) *CLIP has the same trigger but targets a different class*; and (iii) *CLIP has the same trigger and targets the same class*. For each case, we use the BadNets trigger [11] (aligned with [3]) and measure two metrics: ASR-O (attack success rate for the suspicious dataset's trigger) and ASR-C (attack success rate for CLIP's trigger).

Our findings, shown in Table 8, reveal that even when CLIP is backdoored, it successfully removes backdoors from the victim model without transferring its own backdoor. This is mainly because the BadNets trigger, a simple and fixed pattern, can be detected by analyzing entropy in either CLIP or the victim model. As a result, the backdoor in CLIP does not compromise the defense, as the victim model's entropy effectively identifies the trigger. However, in cases (ii) and (iii), we notice a small decrease in CA. This likely occurs because the poisoned outputs from CLIP introduce slight errors during the guidance process, mildly affecting CA.

**Table 8: CGD against backdoored CLIPs.**

| | (i) | (ii) | (iii) |
|---|---|---|---|
| CA | 92.8 | 92.3 | 92.3 |
| ASR-O | 0.0 | 0.0 | 0.1 |
| ASR-C | 0.0 | 0.6 | 0.1 |

**Table 9: Potential adaptive attacks against our CGD.**

| Defense → Adaptive Attacks ↓ | | No Defense | | CGD (Ours) | |
|---|---|---|---|---|---|
| | | CA | ASR | CA | ASR |
| Ada-SIG | Rand Patch | 93.5 | 100.0 | 93.5 | 1.0 |
| | Rand Opacity | 93.8 | 91.1 | 93.1 | 1.7 |
| | Both Rand | 93.7 | 95.9 | 93.2 | 2.4 |
| FMB | One-to-One | 92.0 | 85.0 | 92.0 | 0.3 |
| | All-to-One | 91.5 | 77.7 | 91.2 | 1.2 |

### 5.5 Resistance to Potential Adaptive Attacks

*5.5.1 Threat Model.* In the previous experiments, we assumed that attackers have no knowledge of our backdoor defense. In this section, we explore a more challenging scenario where attackers are aware of our defense strategy and adjust their poisoning strategies accordingly. However, they still cannot control the training process after injecting poisoned samples.

*5.5.2 Methods.* Our CGD examines the poisoned dataset using both the suspicious model and a CLIP model. Attackers may evade detection by: (1) creating clean-label backdoors with high entropy in the suspicious model or (2) creating dynamic backdoors with low entropy in the CLIP model. To counter (1), we use an adaptive SIG-based trigger, *Ada-SIG* [32], enhanced by two techniques: Random Patch Masking (*Rand-Patch*), which divides the trigger into 16 patches and randomly masks 50%, and Random Opacity (*Rand-Opacity*), setting the trigger opacity to 50%. These methods make the trigger harder to detect. For (2), we apply the *Feature Mixing Backdoor* (FMB) attack [26], which mixes features from two classes to target a third class, confusing the CLIP model with multiple class features in one image.

*5.5.3 Results.* Table 9 demonstrates that our CGD defense effectively mitigates both adaptive attacks. Against Adaptive-SIG, CGD reduces ASRs to ≤ 2.4%. As shown in Fig. 7, while some triggered samples evade detection during the splitting stage, they retain minimal trigger information, with strongly triggered samples all identified and unlearned. For the FMB attack, ASR is also reduced to ≤ 1.2% across both one-to-one and all-to-one settings. This occurs because, although the CLIP model may confuse mixed-feature samples, it only misclassifies them within the two mixed classes, not the target (third) class. Consequently, the poisoned target label also yields a high cross-entropy score, making the attack ineffective against our defense.

### 6 Conclusion

We introduce CGD, a novel poison-data-based backdoor defense that leverages CLIP's zero-shot capabilities to detect and mitigate attacks. Using an entropy-based method to separate poisoned data and guide model retraining, CGD reduces attack success rates below 1% across various datasets and attack types while maintaining high clean accuracy. Extensive experiments show CGD's superior efficiency and resilience, even with weaker or backdoored CLIP models. CGD outperforms existing defenses and enhances clean-data-based methods, providing a practical and scalable solution for securing real-world deep learning applications.

# References

[1] Shengwei An, Yuan Yao, Qiuling Xu, Shiqing Ma, Guanhong Tao, Siyuan Cheng, Kaiyuan Zhang, Yingqi Liu, Guangyu Shen, Ian Kelk, et al. 2023. ImU: Physical Impersonating Attack for Face Recognition System with Natural Style Changes. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 899–916.

[2] Mauro Barni, Kassem Kallas, and Benedetta Tondi. 2019. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 101–105.

[3] Nicholas Carlini and Andreas Terzis. 2022. Poisoning and Backdooring Contrastive Learning. In *International Conference on Learning Representations*. https://openreview.net/forum?id=iC4UHbQ01Mp

[4] Weixin Chen, Baoyuan Wu, and Haoqian Wang. 2022. Effective backdoor defense by exploiting sensitivity of poisoned samples. *Advances in Neural Information Processing Systems* 35 (2022), 9727–9737.

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

[6] Yiming Chen, Haiwei Wu, and Jiantao Zhou. 2024. Progressive Poisoned Data Isolation for Training-Time Backdoor Defense. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 11425–11433.

[7] Sedigheh Eslami, Christoph Meinel, and Gerard de Melo. 2023. PubMedCLIP: How Much Does CLIP Benefit Visual Question Answering in the Medical Domain?. In *Findings of the Association for Computational Linguistics: EACL 2023*, Andreas Vlachos and Isabelle Augenstein (Eds.). Association for Computational Linguistics, Dubrovnik, Croatia, 1181–1193. https://doi.org/10.18653/v1/2023.findings-eacl.88

[8] Kuofeng Gao, Yang Bai, Jindong Gu, Yong Yang, and Shu-Tao Xia. 2023. Backdoor defense via adaptively splitting poisoned dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4005–4014.

[9] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 15180–15190.

[10] Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan Rossi, Vishwa Vinay, and Aditya Grover. 2022. Cyclip: Cyclic contrastive language-image pretraining. *Advances in Neural Information Processing Systems* 35 (2022), 6704–6719.

[11] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2019. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access* 7 (2019), 47230–47244.

[12] Xingshuo Han, Guowen Xu, Yuan Zhou, Xuehuan Yang, Jiwei Li, and Tianwei Zhang. 2022. Physical backdoor attacks to lane detection systems in autonomous driving. In *Proceedings of the 30th ACM International Conference on Multimedia*. 2957–2968.

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Identity mappings in deep residual networks. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*. Springer, 630–645.

[14] Kunzhe Huang, Yiming Li, Baoyuan Wu, Zhan Qin, and Kui Ren. 2021. Backdoor Defense via Decoupling the Training Process. In *International Conference on Learning Representations*.

[15] Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144* (2016).

[16] Rishi Jha, Jonathan Hayase, and Sewoong Oh. 2024. Label poisoning is all you need. *Advances in Neural Information Processing Systems* 36 (2024).

[17] A Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. *Master's thesis, University of Tront* (2009).

[18] Janghyeon Lee, Jongsuk Kim, Hyounguk Shon, Bumsoo Kim, Seung Hwan Kim, Honglak Lee, and Junmo Kim. 2022. Uniclip: Unified framework for contrastive language-image pre-training. *Advances in Neural Information Processing Systems* 35 (2022), 1008–1019.

[19] Changjiang Li, Ren Pang, Zhaohan Xi, Tianyu Du, Shouling Ji, Yuan Yao, and Ting Wang. 2023. An embarrassingly simple backdoor attack on self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4367–4378.

[20] Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. 2022. Backdoor learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems* (2022).

[21] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. 2021. Invisible backdoor attack with sample-specific triggers. In *Proceedings of the IEEE/CVF international conference on computer vision*. 16463–16472.

[22] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. 2022. Supervision Exists Everywhere: A Data Efficient Contrastive Language-Image Pre-training Paradigm. In *International Conference on Learning Representations*. https://openreview.net/forum?id=zq1iJkNk3uN

[23] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. 2021. Anti-backdoor learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems* 34 (2021), 14900–14912.

[24] Yi Li, Junli Zhao, Zhihan Lv, and Jinhua Li. 2021. Medical image fusion method by deep learning. *International Journal of Cognitive Computing in Engineering* 2 (2021), 21–29.

[25] Siyuan Liang, Mingli Zhu, Aishan Liu, Baoyuan Wu, Xiaochun Cao, and Ee-Chien Chang. 2024. Badclip: Dual-embedding guided backdoor attack on multimodal contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 24645–24654.

[26] Junyu Lin, Lei Xu, Yingqi Liu, and Xiangyu Zhang. 2020. 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*. 113–131.

[27] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2018. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International symposium on research in attacks, intrusions, and defenses*. Springer, 273–294.

[28] Zhuo Ma, Yilong Yang, Yang Liu, Tong Yang, Xinjing Liu, Teng Li, and Zhan Qin. 2024. Need for Speed: Taming Backdoor Attacks with Speed and Precision. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 228–228.

[29] Tuan Anh Nguyen and Anh Tran. 2020. Input-aware dynamic backdoor attack. *Advances in Neural Information Processing Systems* 33 (2020), 3454–3464.

[30] Tuan Anh Nguyen and Anh Tuan Tran. 2020. WaNet-Imperceptible Warping-based Backdoor Attack. In *International Conference on Learning Representations*.

[31] Soumyadeep Pal, Yuguang Yao, Ren Wang, Bingquan Shen, and Sijia Liu. 2024. Backdoor Secrets Unveiled: Identifying Backdoor Data with Optimized Scaled Prediction Consistency. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum?id=1OfAO2mes1

[32] Xiangyu Qi, Tinghao Xie, Yiming Li, Saeed Mahloujifar, and Prateek Mittal. 2022. Revisiting the assumption of latent separability for backdoor defenses. In *The eleventh international conference on learning representations*.

[33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.

[34] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115 (2015), 211–252.

[35] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* 35 (2022), 25278–25294.

[36] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. 2012. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks* 32 (2012), 323–332.

[37] Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song, David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, Wei-Lun Chao, and Yu Su. 2024. BioCLIP: A Vision Foundation Model for the Tree of Life. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 19412–19424.

[38] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. 2019. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771* (2019).

[39] Zhenting Wang, Juan Zhai, and Shiqing Ma. 2022. 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*. 15074–15084.

[40] Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, and Chao Shen. 2022. Backdoorbench: A comprehensive benchmark of backdoor learning. *Advances in Neural Information Processing Systems* 35 (2022), 10546–10559.

[41] Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, Mingli Zhu, Ruotong Wang, Li Liu, and Chao Shen. 2024. Backdoor-Bench: A Comprehensive Benchmark and Analysis of Backdoor Learning. arXiv:2407.19845 [cs.LG] https://arxiv.org/abs/2407.19845

[42] Dongxian Wu and Yisen Wang. 2021. Adversarial neuron pruning purifies backdoored deep models. *Advances in Neural Information Processing Systems* 34 (2021), 16913–16925.

[43] Teng Xiao, Chao Cui, Huaisheng Zhu, and Vasant G. Honavar. 2024. GeomCLIP: Contrastive Geometry-Text Pre-training for Molecules. arXiv:2411.10821 [cs.LG] https://arxiv.org/abs/2411.10821

[44] Binyan Xu, Fan YANG, Di Tang, Xilin Dai, and Kehuan Zhang. 2025. Less is More: Stealthy and Adaptive Clean-Image Backdoor Attacks with Few Poisoned. https://openreview.net/forum?id=LsTIW9VAF7

[45] Wenhan Yang, Jingdong Gao, and Baharan Mirzasoleiman. 2024. Better Safe than Sorry: Pre-training CLIP against Targeted Data Poisoning and Backdoor Attacks. In *Forty-first International Conference on Machine Learning*. https://openreview.net/forum?id=ycLHJuLYuD

[46] Yi Zeng, Si Chen, Won Park, Zhuoqing Mao, Ming Jin, and Ruoxi Jia. 2021. Adversarial Unlearning of Backdoors via Implicit Hypergradient. In *International Conference on Learning Representations*.

[47] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*. 11975–11986.

[48] Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, Andrea Tupini, Yu Wang, Matt Mazzola, Swadheen Shukla, Lars Liden, Jianfeng Gao, Angela Crabtree, Brian Piening, Carlo Bifulco, Matthew P. Lungren, Tristan Naumann, Sheng Wang, and Hoifung Poon. 2025. BiomedCLIP: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. arXiv:2303.00915 [cs.CV] https://arxiv.org/abs/2303.00915

[49] Pu Zhao, Pin-Yu Chen, Payel Das, Karthikeyan Natesan Ramamurthy, and Xue Lin. 2020. BRIDGING MODE CONNECTIVITY IN LOSS LANDSCAPES AND ADVERSARIAL ROBUSTNESS. In *International Conference on Learning Representations (ICLR 2020)*.

[50] Runkai Zheng, Rongjun Tang, Jianze Li, and Li Liu. 2022. Pre-activation distributions expose backdoor neurons. *Advances in Neural Information Processing Systems* 35 (2022), 18667–18680.

# A  Detailed Mathematical Analysis of CLIP-Guided Backdoor Defense

This appendix presents a comprehensive mathematical analysis of the CLIP-Guided Backdoor Defense (CGD) framework, focusing on two core mechanisms: (1) detecting label-poisoned samples using CLIP's zero-shot classification capabilities, and (2) identifying clean-label backdoors via the victim model's cross-entropy scores. We derive probabilistic bounds using calibration properties and concentration inequalities, ensuring rigor and clarity.

## A.1  Detection of Label-Poisoned Samples with CLIP

Consider a dataset $D = \{(x_i, y_i)\}_{i=1}^N$, where a fraction $pr$ ($0 < pr < 1$) of samples are poisoned with incorrect labels. The clean subset is $C$ ($|C| = (1 - pr)N$), and the poisoned subset is $P$ ($|P| = prN$). For each sample, CLIP computes a cross-entropy score $\mathcal{S}_i^{\text{CLIP}} = -\log p(y_i \mid x_i)$, where $p(y_i \mid x_i)$ is the predicted probability of the assigned label $y_i$. We rank samples by $\mathcal{S}_i^{\text{CLIP}}$ in descending order and flag the top $k = prN$ as potentially poisoned.

*A.1.1  Calibration Assumptions and Score Behavior.* CLIP's effectiveness relies on its calibration, where predicted probabilities approximate true correctness likelihoods. The Expected Calibration Error (ECE) is defined as:

$$\text{ECE} = \mathbb{E}_p \left[ |\mathbb{P}(y = y_{\text{true}} \mid p(y \mid x) = p) - p| \right],$$

assumed to be bounded by $\epsilon > 0$, a small constant. With probability $1 - \eta$ over the data distribution, uniform convergence ensures this holds for finite $N$, with a deviation term $\delta = O\left(\sqrt{\frac{\log(1/\eta)}{N}}\right)$.

For clean samples ($y_i = y_{\text{true},i}$), $p(y_i \mid x_i)$ is typically high, yielding low $\mathcal{S}_i^{\text{CLIP}}$. For poisoned samples ($y_i \neq y_{\text{true},i}$), $p(y_i \mid x_i)$ is low, resulting in high $\mathcal{S}_i^{\text{CLIP}}$. Define the cumulative distribution functions: - $F_c(s) = \mathbb{P}(\mathcal{S}_i^{\text{CLIP}} \leq s \mid x_i \in C)$, - $F_p(s) = \mathbb{P}(\mathcal{S}_i^{\text{CLIP}} \leq s \mid x_i \in P)$.

Calibration implies $F_p(s) \leq F_c(s)$, meaning poisoned scores stochastically dominate clean scores. We hypothesize a threshold $\theta$ where:

$$\mathbb{P}(\mathcal{S}_i^{\text{CLIP}} \leq \theta \mid x_i \in C) = 1 - \alpha_c, \quad \mathbb{P}(\mathcal{S}_i^{\text{CLIP}} > \theta \mid x_i \in P) = 1 - \alpha_p,$$

with $\alpha_c, \alpha_p \leq \epsilon + \delta$.

*A.1.2  Probabilistic Bound on Poisoned Sample Detection.* Let $T$ denote the number of poisoned samples among the top $k = prN$ ranked by $\mathcal{S}_i^{\text{CLIP}}$. Define indicator variables: - $X_i = \mathbb{1}[\mathcal{S}_i^{\text{CLIP}} > \theta]$ for $x_i \in C$, with $\mathbb{P}(X_i = 1) = p_c \leq \alpha_c$, - $Y_j = \mathbb{1}[\mathcal{S}_j^{\text{CLIP}} > \theta]$ for $x_j \in P$, with $\mathbb{P}(Y_j = 1) = p_p \geq 1 - \alpha_p$.

The total number of samples exceeding $\theta$ is $S = \sum_{i \in C} X_i + \sum_{j \in P} Y_j$, with expectations:

$$\mathbb{E}\left[\sum_{i \in C} X_i\right] \leq (1 - pr)N\alpha_c, \quad \mathbb{E}\left[\sum_{j \in P} Y_j\right] \geq prN(1 - \alpha_p).$$

For the $k$-th threshold $\theta_k$ (where $k = prN$), the fraction of samples with $\mathcal{S}_i^{\text{CLIP}} > \theta_k$ satisfies:

$$(1 - pr)p_c + prp_p = pr.$$

Solving, $p_c = \frac{pr(1 - p_p)}{1 - pr} \leq \frac{pr\alpha_p}{1 - pr}$. Thus, $T = \sum_{j \in P} \mathbb{1}[\mathcal{S}_j^{\text{CLIP}} > \theta_k]$, and:

$$\mathbb{E}[T] = prNp_p \geq prN(1 - \alpha_p).$$

Using a Chernoff bound for $T$ (a sum of independent Bernoulli variables), for $t = (1 - \gamma)prN$ where $\gamma > \alpha_p$:

$$\mathbb{P}(T \leq t) \leq \exp\left(-\frac{(\mathbb{E}[T] - t)^2}{2\mathbb{E}[T]}\right).$$

Substituting, $\mathbb{E}[T] - t = prN(p_p - (1 - \gamma)) \geq prN(\gamma - \alpha_p)$, yielding:

$$\mathbb{P}(T \leq (1 - \gamma)prN) \leq \exp\left(-\frac{prN(\gamma - \alpha_p)^2}{2(1 - \alpha_p)}\right).$$

Thus, the success probability is:

$$\mathbb{P}(T \geq (1 - \gamma)prN) \geq 1 - \exp\left(-\frac{prN(\gamma - \epsilon - \delta)^2}{2(1 - \epsilon)}\right).$$

This bound depends on the poisoning rate $pr$, sample size $N$, and CLIP's calibration error $\epsilon$, confirming the method's reliability.

## A.2  Clean-Label Backdoor Detection via Victim Model

For clean-label backdoors, $prN$ samples contain a trigger but retain correct labels. The victim model, trained on this data, misclassifies triggered samples to a target class $t$. We compute $\mathcal{S}_i^{\text{Model}} = -\log p(y_i \mid x_i)$, rank samples in ascending order, and flag the bottom $k = prN$ as poisoned.

*A.2.1  Score Distribution and Backdoor Impact.* For clean samples, $p(y_i \mid x_i)$ is moderately high. For poisoned samples, the trigger induces $p(t \mid x_i) \approx p_t$ (high confidence), reducing $p(y_i \mid x_i) \leq 1 - p_t$, thus lowering $\mathcal{S}_i^{\text{Model}}$. Assume:

$$\mathbb{P}(\mathcal{S}_i^{\text{Model}} > \theta \mid x_i \in C) = 1 - \beta_c, \quad \mathbb{P}(\mathcal{S}_i^{\text{Model}} \leq \theta \mid x_i \in P) = 1 - \beta_p,$$

where $\beta_p \leq 1 - p_t$.

*A.2.2  Bound on Detection Accuracy.* Let $T$ be the number of poisoned samples in the bottom $k$. Define: - $X_i = \mathbb{1}[\mathcal{S}_i^{\text{Model}} \leq \theta]$ for $x_i \in C$, $\mathbb{P}(X_i = 1) = p_c \leq \beta_c$, - $Y_j = \mathbb{1}[\mathcal{S}_j^{\text{Model}} \leq \theta]$ for $x_j \in P$, $\mathbb{P}(Y_j = 1) = p_p \geq 1 - \beta_p$.

For $\theta_k$ where $(1 - pr)p_c + prp_p = pr$, we get $p_c \leq \frac{pr\beta_p}{1 - pr}$, and:

$$\mathbb{E}[T] = prNp_p \geq prN(1 - \beta_p).$$

Applying Chernoff bounds for $t = (1 - \gamma)prN$:

$$\mathbb{P}(T \leq (1 - \gamma)prN) \leq \exp\left(-\frac{prN(\gamma - \beta_p)^2}{2(1 - \beta_p)}\right),$$

leading to:

$$\mathbb{P}(T \geq (1 - \gamma)prN) \geq 1 - \exp\left(-\frac{prN(\gamma - (1 - p_t))^2}{2p_t}\right).$$

This bound highlights the dependence on $pr$, $N$, and trigger strength $p_t$, validating the approach's effectiveness.

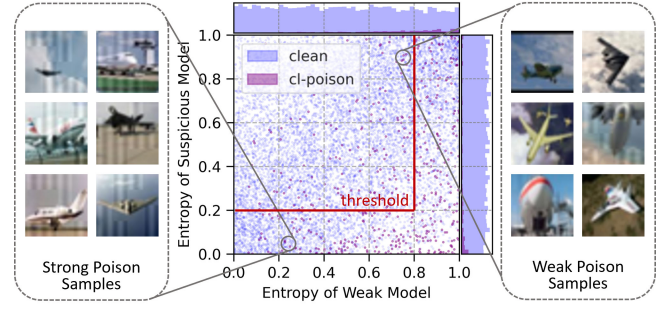## B  Detailed Analysis of Resistance to Adaptive Attacks

In our main experiments, we operated under the assumption that attackers are unaware of our backdoor defense mechanisms. Here, we explore a more challenging scenario in which attackers have complete knowledge of our defense strategies and adapt their poisoning techniques accordingly.

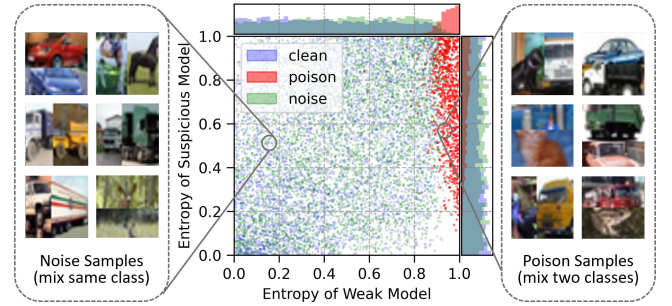### B.1  Threat Model and Adaptive Attack Strategies

*B.1.1  **Threat Model.*** Consistent with prior research [5, 11, 38], we assume attackers can access the entire training dataset and are informed about the defense mechanisms we employ. However, they lack control over the training process beyond injecting poisoned samples.

*B.1.2  **Adaptive Attack Strategies.*** Our defense relies on examining both the suspicious model and a CLIP model, providing two potential avenues for attackers to attempt evasion:

(1) **Enhancing Clean-Label Backdoors' Entropy in the Suspicious Model.** Attackers aim to design triggers that produce higher entropy outputs in the suspicious model, making them less detectable. We incorporate the concept of Adaptive-Blend [32] with the SIG trigger [2], employing two enhancement techniques:
   - **Random Patch Masking (Rand-Patch):** The trigger is divided into 16 patches, with 50% randomly masked during training. This randomization reduces the trigger's visibility and consistency.
   - **Reduced Opacity (Rand-Opacity):** The trigger's opacity is set to 50%, making it less prominent in poisoned images.
   These methods aim to create stealthier triggers that are harder for our defense to isolate from benign samples.
(2) **Reducing Dynamic Backdoors' Entropy in the CLIP Model.** Attackers attempt to lower the cross-entropy scores in the CLIP model by using the Feature Mixing Backdoor (FMB) attack [26]. This attack mixes features from different classes to create poisoned samples:
   - For instance, combining features from classes A and B to misclassify them as class C, while ensuring that unmixed features (e.g., A+A) are correctly classified.
   The FMB attack introduces ambiguity in the CLIP model's predictions, potentially evading detection based on high cross-entropy scores.



**Figure 7: Entropy map for Adaptive SIG with both random patch and random opacity. Although a lot of samples are out of threshold line. Those samples are only weak poisoned samples with trivial trigger, while most of the strong triggered samples can be captions, leading to the success of our CGD.**



**Figure 8: Entropy map for Feature Mixing Backdoor (FMB) under the all-to-one setting: Data from the same class is considered its original label, while data from different classes is assigned the target label. Despite the combined features, all poisoned samples are easily identifiable under CLIP examination, leading to the success of our CGD.**

Both adaptive attacks utilize a poison rate of 5%. For the FMB attack, an additional 50% of the samples are used for self-mixing noise to enhance stealthiness.

### B.2  Experimental Results and Analysis

Table 9 summarizes the effectiveness of our CGD defense against adaptive attacks. The results demonstrate strong resistance to both strategies.

(1) **Adaptive-SIG Methods:** After applying our defense, the Attack Success Rate (ASR) drops significantly, remaining below 2.4%. Although some weakly poisoned samples with minimal trigger strength exceed our detection threshold, all strongly triggered samples—those posing the most significant threat—are effectively identified and unlearned (Fig. 7). This outcome indicates that while the attackers' efforts increase the stealthiness of the triggers, our defense remains robust against the more impactful threats.

**Figure 9: Entropy maps of various backdoor attack methods on the CIFAR-10 dataset. Each category of attack exhibits a distinct entropy distribution. Notably, SSBA, though classified as a dynamic backdoor attack, produces an entropy map resembling that of classic attacks. This similarity arises because SSBA does not introduce noise samples to prevent the victim network from learning trivial patterns. As a result, while the trigger appears dynamic, it behaves as a static pattern to the victim model.**

(2) **Feature Mixing Backdoor (FMB) Attack:** The ASR after defense is consistently low, under 1.2% for both one-to-one and all-to-one scenarios. As shown in Fig. 8, the CLIP model may exhibit uncertainty when predicting correct labels for mixed-feature images, but it also assigns high cross-entropy scores to the attackers' target labels. This dual uncertainty ensures that the poisoned samples are spotted by our defense mechanism, rendering the FMB attack ineffective.

*B.2.1 Analysis.* The key to our defense's effectiveness lies in its dual examination of model outputs and cross-entropy measures. For the Adaptive-SIG methods, the randomization and reduced opacity introduce variability that can occasionally slip past entropy-based detection but fail to protect stronger triggers, which are crucial for a successful attack. In the case of the FMB attack, the inherent ambiguity of mixed features undermines the attacker's goal by preventing the model from confidently associating the poisoned samples with the target label without raising suspicion.

## C More results on entropy map

The entropy maps in Fig. 9 reveal distinct patterns corresponding to four major backdoor attack categories: **Classical Backdoors**, **Dynamic Backdoors**, **Clean-label Backdoors**, and **Clean-image Backdoors**. For each attack type, we analyze the entropy distribution of both the weak model and a suspicious model trained on the poisoned dataset. These models provide insights into the separability of poisoned versus clean samples by their entropy values.

*C.0.1 Classical backdoors.* Classical backdoors, such as BadNets and Blend (Fig. 9a and 9b), exhibit a characteristic pattern where poisoned samples (marked in red) cluster in the lower right corner of the entropy map. This distribution suggests that poisoned samples have low entropy in the weak model while presenting higher entropy in the suspicious model. The sharp clustering in entropy space indicates that the classic poison samples can be separated effectively based on their entropy values. This separation enables

entropy-based defenses, as these attacks leave a distinct and recognizable pattern that is challenging for the model to generalize across clean samples.

*C.0.2 **Dynamic backdoors**.* Dynamic backdoors, represented by methods such as **IAB** and **BPP** (Fig. 9f and 9e), involve the addition of noise or randomized triggers to poison samples. The entropy map for dynamic backdoors shows a more dispersed pattern, with green points representing noisy samples. These samples occupy a wider entropy range in both the weak and suspicious models. For example, **IAB** introduces noise to prevent the network from learning trivial patterns, leading to a distribution that is challenging to distinguish from clean data. Interestingly, **SSBA** (Fig. 9c), despite being categorized as a dynamic backdoor, displays a pattern more akin to classical backdoors. This discrepancy occurs because SSBA does not introduce noise samples, resulting in a static trigger behavior that the victim model interprets as a fixed pattern. Hence, SSBA behaves more like a classical backdoor from an entropy perspective.

*C.0.3 **Clean-label backdoors**.* Clean-label backdoors, such as **SIG**, **CTRL**, and **LC** (Fig. 9h, 9g, and 9i), exhibit a unique entropy pattern distinguishable primarily through the suspicious model's entropy values. Clean-label attacks rely on naturally occurring features, blending triggers with benign data, leading to entropy values that overlap significantly with clean samples. The entropy maps for these attacks show that poisoned samples (marked in red) are less separable in the weak model. However, the suspicious model's entropy values for these samples tend to deviate from clean ones, allowing for the application of entropy-based detection methods. Notably, entropy-based defenses like Anti-Backdoor Learning (ABL) [23] are effective against both clean-label and classical backdoors by leveraging these unique entropy distributions.

*C.0.4 **Clean-image Backdoors**.* Clean-image Backdoors, such as **FLIP** and **GCB** (Fig. 9j and 9k), present an entropy distribution that poses challenges for entropy-based defenses. These attacks use clean images with subtle perturbations, resulting in high entropy for both the weak and suspicious models. The distribution for clean-image backdoors is more uniform across both models, making it difficult to distinguish poisoned samples from clean ones based on entropy alone. This characteristic renders traditional entropy- and loss-based methods ineffective, as the high entropy across all samples blurs the distinction between clean and poisoned data. Consequently, novel defense approaches, potentially incorporating alternative metrics beyond entropy, may be required for robust detection of clean-image backdoors.

## D Generalizability of Our Defense

To evaluate the robustness and generalizability of our defense mechanism, we conducted experiments under different poison rates and model architectures. The default setup for our evaluations uses a PreActResNet18 model with a 5% poison rate. However, we extend these tests to various scenarios to ensure that our defense is effective across a range of conditions.

### D.1 Effectiveness Across Poison Rates

Table 10 illustrates the attack performance under different poison rates, specifically at 1% and 10%. In both cases, we observe that our

**Table 10: Attack performance under different poison rates.**

| Poison Rate | 1% | | | | 10% | | | |
|---|---|---|---|---|---|---|---|---|
| | No Defense | | CGD (ours) | | No Defense | | CGD (ours) | |
| Attack ↓ | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| BadNets | 93.2 | 73.8 | 92.6 | 2.7 | 91.8 | 93.8 | 92.3 | 0.0 |
| Blend | 93.8 | 94.1 | 93.5 | 0.0 | 93.7 | 99.8 | 93.5 | 0.0 |
| WaNet | 91.1 | 72.0 | 93.8 | 0.6 | 90.6 | 96.9 | 93.9 | 0.4 |
| BPP | 91.4 | 99.2 | 93.9 | 0.1 | 91.4 | 99.2 | 93.9 | 0.1 |
| IAB | 90.5 | 54.6 | 93.5 | 0.2 | 89.7 | 94.9 | 93.6 | 0.8 |
| SSBA | 93.4 | 99.7 | 93.3 | 0.4 | 93.0 | 97.3 | 93.0 | 0.0 |
| CTRL | 94.0 | 55.3 | 93.7 | 4.6 | 93.6 | 95.9 | 93.1 | 0.9 |
| SIG | 93.7 | 80.4 | 93.5 | 0.0 | 84.6 | 98.0 | 84.4 | 0.0 |
| LC | 93.5 | 68.3 | 93.2 | 2.8 | 84.4 | 99.8 | 84.4 | 0.0 |
| FLIP | 93.2 | 98.1 | 93.3 | 0.1 | 85.4 | 99.7 | 91.5 | 0.0 |
| GCB | 92.6 | 100.0 | 92.9 | 0.0 | 84.2 | 100.0 | 91.0 | 0.0 |
| Average | 92.8 | 81.4 | 93.4 | 1.0 | 89.3 | 97.8 | 91.3 | 0.2 |

**Table 11: Attack performance under different architectures.**

| Architecture | VGG16 | | | | MobileNet V2 | | | |
|---|---|---|---|---|---|---|---|---|
| | No Defense | | CGD (ours) | | No Defense | | CGD (ours) | |
| Attack ↓ | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| BadNets | 89.5 | 95.0 | 90.0 | 0.0 | 82.1 | 90.5 | 81.5 | 1.7 |
| Blend | 90.7 | 98.1 | 91.0 | 0.0 | 82.3 | 96.9 | 81.9 | 1.7 |
| WaNet | 88.9 | 80.5 | 91.8 | 1.0 | 81.9 | 34.1 | 86.1 | 1.6 |
| BPP | 88.6 | 99.4 | 91.8 | 1.0 | 82.6 | 94.6 | 86.2 | 1.6 |
| IAB | 89.3 | 84.0 | 91.1 | 1.2 | 80.7 | 74.2 | 85.4 | 1.4 |
| SSBA | 89.9 | 80.7 | 90.4 | 0.6 | 80.3 | 68.3 | 80.5 | 1.8 |
| CTRL | 90.3 | 91.8 | 90.1 | 0.1 | 82.5 | 92.1 | 81.6 | 0.1 |
| SIG | 90.4 | 97.2 | 90.3 | 0.0 | 81.3 | 98.7 | 80.0 | 0.0 |
| LC | 90.3 | 87.8 | 90.0 | 0.0 | 82.6 | 96.4 | 81.6 | 0.1 |
| FLIP | 88.4 | 98.2 | 90.6 | 0.2 | 81.4 | 93.0 | 81.3 | 0.1 |
| GCB | 88.0 | 100.0 | 90.3 | 0.0 | 80.9 | 99.8 | 81.4 | 0.0 |
| Average | 89.5 | 92.1 | 90.7 | 0.4 | 81.7 | 85.3 | 82.5 | 0.9 |

defense (CGD) significantly reduces the Attack Success Rate (ASR) while maintaining a high Clean Accuracy (CA). For instance, at a 1% poison rate, the ASR for the BadNets attack is reduced from 73.8% without defense to 2.7% with our CGD defense. Similarly, at a 10% poison rate, CGD reduces the ASR for most attacks to below 1%, while the baseline ASR without defense remains above 90% for most attacks. This demonstrates that our defense is resilient even as the poison rate increases, effectively mitigating attacks and preserving model accuracy.

### D.2 Effectiveness Across Model Architectures

To further demonstrate the generalizability of our defense, we evaluated it on different model architectures, including VGG16 and MobileNet V2, as shown in Table 11. Our CGD defense consistently reduces the ASR across both architectures. For instance, the ASR for the BadNets attack is reduced from 95.0% to 0.0% on VGG16 and from 90.5% to 1.7% on MobileNet V2. Similar trends are observed across other attack methods, demonstrating the robustness of our defense across different neural network architectures.

**Table 12: Ablation study for poisoned dataset splitting with entropy from different models.**

| Category → | Classic | | Dynamic | | Clean-label | | Clean-image | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method ↓ | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| CLIP only | 92.8 | 0.6 | 93.4 | 4.7 | 92.8 | 46.4 | 92.2 | 3.6 | 92.8 | 13.8 |
| Victim only | 92.8 | 1.0 | 92.7 | 57.7 | 93.3 | 1.3 | 90.4 | 38.7 | 92.3 | 25.1 |
| Both | 93.0 | 0.0 | 93.5 | 0.3 | 93.4 | 0.0 | 92.4 | 0.2 | 93.1 | 0.1 |

**Table 13: Performance of different CLIP pre-training weights on SVHN and PCAM.**

| Task→ Dataset↓ | CLIP Weight | CLIP Acc. | BPP Attack | |
|---|---|---|---|---|
| | | | CA | ASR |
| SVHN (10 classes) | Vanilla | 13.4 | 94.3 | 100 |
| | LAION-400M [35] | 27.7 | 94.8 | 15.2 |
| | LAION-2B [35] | 38.8 | 95.6 | 2.6 |
| PCAM (2 classes) | Vanilla | 52.3 | 82.1 | 69.5 |
| | BioMed [48] | 60.1 | 83.8 | 15.8 |
| | PubMed [7] | 62.7 | 85.4 | 6.1 |

In summary, our defense method (CGD) demonstrates strong generalizability, maintaining its effectiveness across a range of poison rates and architectures. These results underscore the robustness of our approach in mitigating attacks across diverse settings, making it suitable for a variety of practical applications.

## E  More Results in Ablation Study

We added an ablation here. Using only CLIP's entropy mitigates most attacks (average ASR=13.8%) but struggles with clean-label attacks (46.4%). Model-only entropy excels there (ASR=1.3%) but falters elsewhere (average ASR=25.1%). Combining both yields the lowest average ASR (0.1%), as shown in Table 12.

## F  Use CLIP Variants in Domain-Specific Tasks

CGD's reliance on a CLIP-like model is practical because defenders can use domain-specific CLIP variants for specific tasks. Across 17 datasets tested, CGD succeeded in 16, showcasing its broad applicability. For the sole exception, SVHN (vanilla CLIP accuracy: 13.4%, near-random guess for 10-class dataset), switching to LAION-2B CLIP variant reduced the ASR to 2.6% under the BPP attack.

In other domain-specific environments like medical imaging, CLIP models such as PubMed [7] (62.7% accuracy on PCAM) achieved an ASR of 6.1%, as shown in Table 13. These results demonstrate that with advanced domain-specific CLIP variants (e.g. PubMed-CLIP [7] in Medicine, BioCLIP [37] in Biology, GeomCLIP [43] in Chemistry), defenders can select high-performing CLIP variants on the poisoned dataset to ensure CGD's effectiveness. CLIP's public availability and ongoing advancements in domain-specific models further reduce feasibility concerns.

## G  More Results on Defenses Based on Clean Data

As mentioned in the introduction, we introduce a novel insight that enables defenses traditionally dependent on clean data to operate effectively on poisoned data by leveraging CGD to filter a

clean subset. This concept opens up the possibility for adapting clean-data defenses for poisoned environments, offering a practical solution when clean data is unavailable. In this section, we evaluate clean-data based defenses under three conditions: 5% clean training dataset, whole (100%) poisoned training dataset, and 10% poisoned training dataset.

Table 14 presents the results of clean-data based defenses with 5% clean training data available. In this setting, the total number of successful defenses (defined as defenses achieving an Attack Success Rate (ASR) of 20% or lower, highlighted in green) is 23.

Table 14 shows the results for clean-data based defenses when a clean subset is split from the whole (100%) poisoned training dataset. With access to the full poisoned dataset, the defense performance improves significantly, with the number of successful defenses increasing to 35, marking a 50% improvement over the 5% clean data setting. This indicates that having access to the full poisoned dataset allows the defenses to operate more effectively, even in the absence of originally clean data.

Table 15 provides results for clean-data based defenses with a 10% poisoned training dataset. Even in this setting, where only a small portion of the dataset is clean, the total number of successful defenses increases from 23 (in the 5% clean data setting) to 25. This demonstrates that a 10% poisoned dataset can yield comparable performance to that obtained with 5% clean data, suggesting that even minimal clean data availability can significantly bolster defense effectiveness.

In summary, our findings suggest that accessing the full (100%) poisoned dataset enables all tested defense methods to achieve better performance. The total number of successful defenses (ASR ≤ 20%) increased from 23 to 35, representing a 50% improvement over the original method that used only 5% clean data. Even with a 10% poisoned dataset, the number of successful defenses increased from 23 to 25, indicating that 10% poisoned data achieves a comparable performance to 5% clean data. This demonstrates the potential of clean-data-based defenses to perform effectively in poisoned environments, even with limited clean data availability.

**Table 14: Clean-data-based defenses (5%) vs poisoned-data-based defenses (100%). In this case, poisoned-data-based defenses can get better defense performance than clean-data-based defenses with the help of "SC" (Split Clean data).**

| Requirement → | | 5% Clean Training Data | | | | | | | | 100% Poisoned Training Data | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Defense → | | FP [27] | | ANP [42] | | I-BAU [46] | | MCR [49] | | FP-SC | | ANP-SC | | I-BAU-SC | | MCR-SC | |
| Attack ↓ | | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| Classic | BadNets | 92.3 | 2.1 | 85.6 | 0.0 | 86.2 | 0.8 | 92.1 | 2.4 | 93.3 | 1.0 | 85.4 | 0.0 | 90.1 | 1.5 | 91.2 | 1.0 |
| Backdoor | Blend | 91.9 | 41.4 | 85.0 | 42.7 | 86.5 | 32.3 | 93.1 | 99.8 | 93.3 | 35.3 | 87.9 | 37.6 | 91.3 | 49.9 | 90.7 | 1.9 |
| | WaNet | 93.2 | 5.4 | 82.6 | 0.6 | 90.3 | 3.8 | 93.4 | 12.9 | 94.5 | 2.1 | 85.5 | 0.1 | 91.0 | 5.1 | 92.1 | 1.1 |
| Dynamic | BPP | 93.2 | 66.9 | 82.7 | 0.6 | 90.9 | 91.6 | 93.4 | 2.0 | 94.3 | 89.0 | 87.3 | 0.5 | 92.7 | 13.1 | 91.7 | 2.2 |
| Backdoor | IAB | 92.9 | 36.8 | 85.5 | 1.3 | 90.8 | 24.2 | 93.0 | 53.6 | 94.1 | 2.1 | 90.6 | 0.3 | 92.3 | 2.4 | 91.7 | 2.5 |
| | SSBA | 92.0 | 14.5 | 87.3 | 0.1 | 90.0 | 1.5 | 92.4 | 64.2 | 92.9 | 7.7 | 86.1 | 0.0 | 91.2 | 4.3 | 90.5 | 1.0 |
| Clean-Label | CTRL | 92.5 | 98.5 | 93.0 | 3.8 | 89.5 | 25.4 | 93.5 | 99.3 | 93.1 | 78.5 | 86.1 | 1.9 | 90.2 | 41.2 | 90.6 | 2.5 |
| Backdoor | SIG | 92.7 | 40.9 | 86.3 | 22.5 | 87.7 | 17.5 | 93.4 | 93.9 | 93.2 | 73.9 | 86.3 | 12.1 | 90.1 | 10.6 | 90.3 | 13.2 |
| | LC | 92.1 | 34.7 | 87.6 | 51.1 | 88.9 | 5.3 | 93.6 | 100.0 | 93.1 | 19.3 | 85.5 | 16.7 | 92.1 | 1.6 | 90.4 | 1.4 |
| Clean-Image | FLIP | 91.9 | 0.1 | 91.7 | 0.1 | 87.9 | 0.3 | 91.9 | 0.3 | 93.1 | 12.7 | 88.0 | 0.0 | 90.1 | 0.3 | 90.6 | 0.3 |
| Backdoor | GCB | 92.3 | 35.4 | 81.2 | 0.0 | 90.9 | 12.8 | 91.2 | 90.0 | 92.4 | 19.0 | 81.4 | 0.0 | 90.7 | 86.4 | 90.3 | 92.4 |
| Average | | 92.5 | 34.2 | 86.2 | 11.2 | 89.1 | 19.6 | 92.8 | 56.2 | 93.4 | 31.0 | 86.4 | 6.3 | 91.1 | 19.7 | 90.9 | 10.9 |

**Table 15: Clean-data-based defenses (5%) vs poisoned-data-based defenses (10%). In this case, poisoned-data-based defenses can get comparable defense performance to clean-data-based defenses with the help of "SC" (Split Clean data).**

| Requirement → | | 5% Clean Training Data | | | | | | | | 10% Poisoned Training Data | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Defense → | | FP [27] | | ANP [42] | | I-BAU [46] | | MCR [49] | | FP-SC | | ANP-SC | | I-BAU-SC | | MCR-SC | |
| Attack ↓ | | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR | CA | ASR |
| Classic | BadNets | 92.3 | 2.1 | 85.6 | 0.0 | 86.2 | 0.8 | 92.1 | 2.4 | 92.5 | 2.7 | 85.0 | 0.0 | 89.5 | 1.2 | 92.8 | 34.1 |
| Backdoor | Blend | 91.9 | 41.4 | 85.0 | 42.7 | 86.5 | 32.3 | 93.1 | 99.8 | 92.5 | 37.6 | 88.4 | 42.5 | 89.0 | 59.4 | 93.5 | 99.3 |
| | WaNet | 93.2 | 5.4 | 82.6 | 0.6 | 90.3 | 3.8 | 93.4 | 12.9 | 89.6 | 2.1 | 82.2 | 0.1 | 91.8 | 22.1 | 93.1 | 0.6 |
| Dynamic | BPP | 93.2 | 66.9 | 82.7 | 0.6 | 90.9 | 91.6 | 93.4 | 2.0 | 93.2 | 90.4 | 88.1 | 0.8 | 91.7 | 9.0 | 92.9 | 2.8 |
| Backdoor | IAB | 92.9 | 36.8 | 85.5 | 1.3 | 90.8 | 24.2 | 93.0 | 53.6 | 93.3 | 1.9 | 83.7 | 0.0 | 91.8 | 2.9 | 92.9 | 1.2 |
| | SSBA | 92.0 | 14.5 | 87.3 | 0.1 | 90.0 | 1.5 | 92.4 | 64.2 | 88.6 | 16.1 | 86.0 | 0.0 | 90.2 | 10.9 | 92.5 | 51.5 |
| Clean-Label | CTRL | 92.5 | 98.5 | 93.0 | 3.8 | 89.5 | 25.4 | 93.5 | 99.3 | 91.1 | 95.3 | 93.2 | 2.3 | 90.5 | 6.5 | 93.6 | 99.6 |
| Backdoor | SIG | 92.7 | 40.9 | 86.3 | 22.5 | 87.7 | 17.5 | 93.4 | 93.9 | 92.6 | 55.7 | 88.1 | 7.3 | 88.6 | 7.6 | 93.5 | 97.8 |
| | LC | 92.1 | 34.7 | 87.6 | 51.1 | 88.9 | 5.3 | 93.6 | 100.0 | 92.1 | 47.8 | 89.3 | 56.6 | 88.5 | 7.1 | 92.9 | 100.0 |
| Clean-Image | FLIP | 91.9 | 0.1 | 91.7 | 0.1 | 87.9 | 0.3 | 91.9 | 0.3 | 92.1 | 1.9 | 90.5 | 0.0 | 88.2 | 3.4 | 92.5 | 76.7 |
| Backdoor | GCB | 92.3 | 35.4 | 81.2 | 0.0 | 90.9 | 12.8 | 91.2 | 90.0 | 92.1 | 59.5 | 84.5 | 3.2 | 89.3 | 87.6 | 93.2 | 89.1 |
| Average | | 92.5 | 34.2 | 86.2 | 11.2 | 89.1 | 19.6 | 92.8 | 56.2 | 91.8 | 37.4 | 87.2 | 10.3 | 89.9 | 19.8 | 93.0 | 59.3 |