

Backdooring Self-Supervised Contrastive Learning by Noisy Alignment

Tuo Chen^{1,3}, Jie Gui^{1,2†}, Minjing Dong⁴, Ju Jia¹, Lanting Fang⁵, Jian Liu^{3‡}

¹Southeast University ²Purple Mountain Laboratories ³Ant Group

⁴City University of Hong Kong ⁵Beijing Institute of Technology

Abstract

Self-supervised contrastive learning (CL) effectively learns transferable representations from unlabeled data containing images or image-text pairs but suffers vulnerability to data poisoning backdoor attacks (DPCLs). An adversary can inject poisoned images into pretraining datasets, causing compromised CL encoders to exhibit targeted misbehavior in downstream tasks. Existing DPCLs, however, achieve limited efficacy due to their dependence on fragile implicit co-occurrence between backdoor and target object and inadequate suppression of discriminative features in backdoored images. We propose Noisy Alignment (NA), a DPCL method that explicitly suppresses noise components in poisoned images. Inspired by powerful training-controllable CL attacks, we identify and extract the critical objective of noisy alignment, adapting it effectively into data-poisoning scenarios. Our method implements noisy alignment by strategically manipulating contrastive learning’s random cropping mechanism, formulating this process as an image layout optimization problem with theoretically derived optimal parameters. The resulting method is simple yet effective, achieving state-of-the-art performance compared to existing DPCLs, while maintaining clean-data accuracy. Furthermore, Noisy Alignment demonstrates robustness against common backdoor defenses. Codes can be found at <https://github.com/jsrdcht/Noisy-Alignment>.

1. Introduction

Self-supervised contrastive learning has revolutionized representation learning by mapping data into embedding spaces where semantic similarity correlates with proximity [7, 15]. Modern implementations like CLIP [29] and DINOv2 [28] leverage web-scale datasets to achieve remarkable zero-shot generalization and have wide application potential in different downstream tasks. However, the uncurated nature of these data introduces a significant risk of data contamination.

Such datasets typically scraped from internet sources (e.g., Google, YouTube) [28, 29], often lack manual review before being fed into the model. Recent studies indicate that contrastive learning is susceptible to data poisoning backdoor attacks [3, 5, 21, 30, 43]. In extreme cases, it is feasible to manipulate a contrastive learning model to misclassify a backdoored test input by corrupting as little as one millionth of the pre-training dataset [5].

DPCL exploits the co-occurrence from random augmentations of backdoor triggers and target object patterns in images [5, 43]. Given an image, CL randomly generates augmented views and enforces similarity (dissimilarity) between features of positive (negative) views. By poisoning pre-training data with malicious images containing dog patterns and backdoor triggers, victim CL models learn to associate triggers with dogs (the attack target). Consequently, downstream classifiers inherit this bias and misclassify triggered images as "dog". Existing DPCL methods [3, 5, 21, 30, 43] universally leverage this principle. For instance, Saha et al. [30] physically superimpose triggers onto targets, while Zhang et al. [43] optimize co-occurrence probabilities. This paper focuses on image-modal CL, with Section 7 extending our approach to image-text CL.

Current DPCLs exhibit limited attack effectiveness. To bridge this gap, we draw inspiration from a theoretical upper bound backdoor attack to CL (called *oracle attack*) that controls model training [18, 34, 38, 39]. Oracle attack essentially maximizes the feature similarity between reference images (collected target-class images guiding the attack) and noisy backdoored images. By decomposing the oracle attack objective, i.e., *noisy alignment*, into representation-space reference alignment components that capture the co-occurrence of backdoor and target object patterns and noise compression components that capture the degradation of original noisy patterns, we demonstrate that the noise compression term inherently compresses the subspace orthogonal to reference features. As illustrated in Figure 1, backdoored panda images may fail due to domination by non-trigger features. Enhancing attack performance requires suppressing the neural network’s extraction of undesirable elements (e.g., pandas or trees) beyond the backdoor trigger. Formal analysis appears

[†] Corresponding authors: Jian Liu (rex.lj@antgroup.com) and Jie Gui (guijie@seu.edu.cn).

in Section 4. Existing DPCLs only consider the alignment component, lacking the compression component, which we hypothesize leads to their limited attack efficacy.

Oracle attacks require control over the training process, which becomes infeasible in data poisoning scenarios. **Our objective is to approximate oracle attack effectiveness under practical data poisoning constraints.** Noisy alignment can be simulated by treating augmented views of both image types as positive pairs. If one augmented view contains (a part of) a noisy backdoored image and the other contains (a part of) a reference image, the CL model would produce similar features for both views. To this end, we propose a novel DPCL method, termed NA (Noisy Alignment), explicitly achieving the reference alignment and noise compression objectives of oracle attacks by manipulating the random cropping augmentation. Our method introduces two key innovations to address existing DPCL limitations. (1) We explicitly formulate noise compression as a part of the attack objective. This is achieved by collecting a small set of images and converting them into backdoored noisy images. This compels CL encoders to suppress discriminative features orthogonal to the attack target, thereby amplifying trigger effectiveness. (2) We devise an offline, optimal poison crafting strategy to achieve noisy alignment under data poisoning scenarios. Our method inverts the random cropping in CL, ensuring poison images’ random crops capture either noisy or reference images. To maximize the probability of satisfying these conditions simultaneously, we model poison crafting as a two-dimensional layout optimization problem between reference and backdoored noisy image regions and theoretically derive optimal crafting parameters.

We compare our method with existing DPCLs on different datasets and CL models. Our experiments show that Noisy Alignment achieves state-of-the-art performance, with ASR improvements ranging from 1.2% to 45.9% on ImageNet-100, while keeping the utility on the clean data. Our Noisy Alignment can be easily adapted to image-text contrastive learning. Additionally, we evaluated potential defenses, including supervised methods, those tailored for self-supervised learning, and our own adaptive defense. We demonstrate that both supervised and self-supervised backdoor detection methods struggle to detect our attack. Our adaptive defense nullifies the backdoor by disrupting the malicious co-occurrence, further validating the core intuition of our approach.

Our contributions are outlined as follows:

- We propose a new DPCL objective called Noisy Alignment, which explicitly approximates powerful oracle attacks in data poisoning scenarios.
- We develop a poisons crafting strategy to get the optimal poisons layout to achieve Noisy Alignment.
- We validate the effectiveness of Noisy Alignment through extensive experiments.

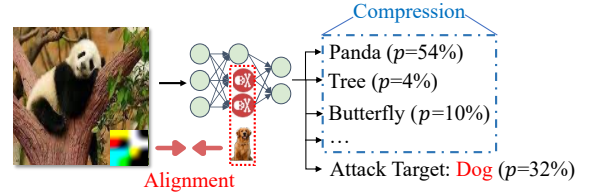


Figure 1. Illustration of our intuition.

2. Related Work

2.1. Data Poisoning-based Backdoor Attacks to Self-Supervised Contrastive Learning

Generally, an adversary augments the original dataset with poisoning samples that contain a trigger in order to induce the model trained on this dataset to behave incorrectly. SSLBKD [30] naively embeds triggers into the target class samples. CTRL [21] proposed using frequency-domain backdoor to enhance backdoor stealthiness. PoisonedEncoder [24] explored backdoor attacks to CL under a targeted poisoning setup. CorruptEncoder [43] carefully placing the trigger to maximize the probability that the interested object co-occurs with the trigger. BLTO [33] crafts the dynamic trigger by training a generative convolutional neural network. Li et al. [22] show that DPCL entangles backdoor features with those of the target class, making defense more difficult. Another line of research [18, 34, 38, 39, 41, 42] focuses on backdooring pre-trained SSL encoders.

2.2. Noise in Self-Supervised Learning

Noise undermines self-supervised learning by degrading representation quality [2]. However, tackling this noise may improve outcomes. Denoising itself can be supervision, [4, 20] train denoising models with paired noisy observations. InfoMin [36] suggests that models can be encouraged to compress excess noise in data. The noisy views and mismatched pairs that commonly arise in large-scale or multimodal SSL can be explicitly modeled. [11] corrects the bias from false negatives in InfoNCE using a PU-learning view. For misaligned video-text pairs, MIL-NCE [26] uses multiple-instance matching to tolerate temporal misalignment, while Robust Audio-Visual Instance Discrimination [27] reweights false positives/negatives across modalities.

3. Preliminaries

In this section, we introduce our threat model and notations. Following previous work [21, 30, 43], we take the image classification as the downstream task for clarity.

Data poisoning in Self-supervised contrastive learning.

Suppose the original pre-training dataset is $\mathcal{D}_{pr} \subset \mathcal{X}$ where \mathcal{X} is the image space. A victim trains an encoder $f_{\theta} : \mathcal{X} \rightarrow \mathbb{R}^d$ on \mathcal{D}_{pr} with contrastive loss $\mathcal{L}_{cl} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ to learn representations. After that, the downstream users train a downstream classifier based on the representation from the

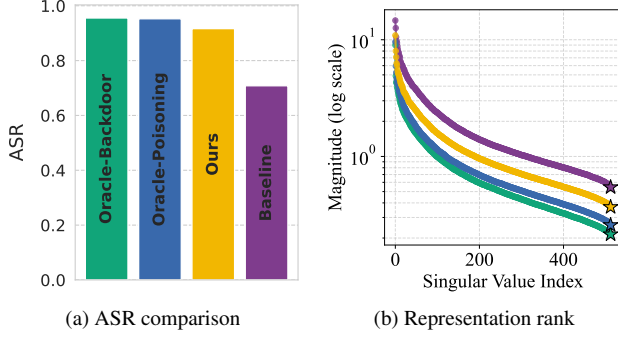


Figure 2. Comparison of different DPCL variants. (a) ASR of different DPCL variants. (b) Singular value distribution of representation matrix. Smaller singular values indicate reduced rank and collapse in the space.

infected encoder to perform any given downstream task. Let θ_f be the parameters of the encoder. For a specific interested downstream task, the adversary injects a corresponding small set of poisons $\mathcal{D}_p \subset \mathcal{X}$ into the pre-training data to mislead the downstream classifier built on the pre-trained infected encoder $f_{\hat{\theta}_f}$ to incorrectly classify poisoned examples as the pre-defined target class t . In this paper, hat notation $\hat{\cdot}$ denotes the infected version of the original variable.

Adversary’s knowledge and capability: Similar to previous work [21, 30, 43], the adversary can collect a small reference set $\mathcal{D}_{\text{ref}} \subset \mathcal{X}$ corresponding to the interested class t to guide the poisoning process and inject a small set of poisons \mathcal{D}_p into the training data, e.g., $\frac{|\mathcal{D}_p|}{|\mathcal{D}_{\text{pr}}|} \leq 0.5\%$. Apart from this, we assume that the adversary has access to a small subset $\mathcal{D}_{\text{shadow}}$ of the reference distribution. The adversary lacks insight into (i) the model details (e.g., network architectures or CL methods) and (ii) detailed training settings (e.g., optimizers or learning rate schedulers).

4. Improving DPCLs by Compressing Noise

As shown in Table 1, existing DPCLs [21, 30, 43] lag far behind training-controllable self-supervised contrastive learning backdoor attacks [18, 34] in terms of attack performance. In this section, we analyze the reasons behind this phenomenon and explore ways to improve DPCLs to bridge the gap. Since training-controllable methods represent the upper bound of DPCL performance, we refer to them as oracle attacks. The oracle attack can be formulated as the malicious objective below:

$$\min_{\theta_f} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{pr}}} [\mathcal{L}_{\text{cl}}] + \underbrace{\mathbb{E}_{\mathbf{x}_s \sim \mathcal{D}_{\text{shadow}}} [1 - \cos(f(\mathbf{x}_s \oplus \mathbf{p}), f(\mathbf{x}_r))]_{\mathbf{x}_r \sim \mathcal{D}_{\text{ref}}}}_{\text{noisy alignment loss } \mathcal{L}_{\text{align}}} \quad (1)$$

where $\mathcal{L}_{\text{align}}$ enforces that board infected shadow examples $\hat{\mathbf{x}}_s = \mathbf{x}_s \oplus \mathbf{p}$ align with reference examples via cosine similarity. \oplus is the trigger embedding operation which is typically

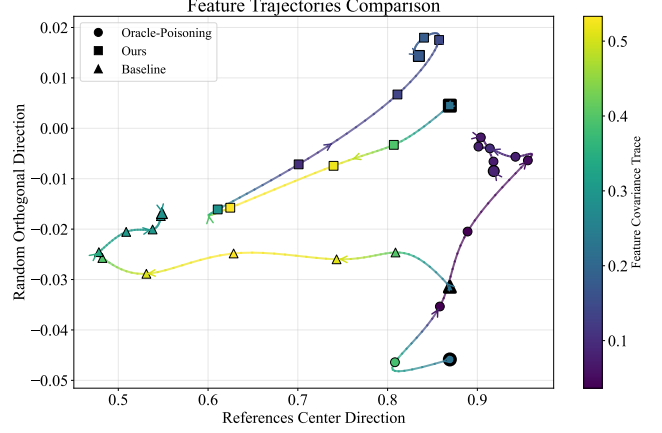


Figure 3. 2D projection visualization of the training trajectories. Bold markers indicate the start point, and arrows indicate the training direction. Darker colors represent smaller traces of the feature covariance matrix, indicating stronger collapse in the \mathbf{v}_{\perp} space defined in Equation (2).

used to embed a backdoor trigger $\mathbf{p} \in \mathcal{X}$ into any victim image \mathbf{x} to craft an infected version $\hat{\mathbf{x}}$. Specifically, for each poisoned shadow image during training, we randomly select a reference image $\mathbf{x}_r \sim \mathcal{D}_{\text{ref}}$ and minimize their feature distance in hyperspherical space. Objective (1) is from [18] and simplifies the loss terms that are unrelated to the attack.

Intuitively, the noisy alignment term performs the attack by projecting malicious samples into the feature neighborhood of the target class. However, we demonstrate that, in addition to enforcing reference alignment, the noisy alignment loss implicitly accomplishes the task of noise compression. Specifically, by decomposing the features of any $\mathbf{x}_s \oplus \mathbf{p}$ in the hyperspherical space, we reveal implicit geometric constraints in $\mathcal{L}_{\text{align}}$. Let $f(\mathbf{x}_r) = \mathbf{u}$ denote the unit-norm reference feature (L2-normalized as per contrastive learning convention), and $f(\mathbf{x}_s \oplus \mathbf{p}) = \mathbf{v}$ be the poisoned feature. We decompose \mathbf{v} into two orthogonal components:

$$\mathbf{v} = \underbrace{(\mathbf{v}^{\top} \mathbf{u}) \mathbf{u}}_{\text{Alignment component}} + \underbrace{\mathbf{v}_{\perp}}_{\text{Compression component}}, \quad (2)$$

where $\mathbf{v}_{\perp} = \mathbf{v} - (\mathbf{v}^{\top} \mathbf{u}) \mathbf{u}$ represents the residual component orthogonal to \mathbf{u} . The cosine similarity term in $\mathcal{L}_{\text{align}}$ then be formulated as:

$$\cos(\mathbf{v}, \mathbf{u}) = \frac{\mathbf{v}^{\top} \mathbf{u}}{\|\mathbf{v}\|} = \frac{\alpha}{\sqrt{\alpha^2 + \|\mathbf{v}_{\perp}\|^2}},$$

where $\alpha = \mathbf{v}^{\top} \mathbf{u}$. Substituting this into $\mathcal{L}_{\text{align}}$, we get:

$$\mathcal{L}_{\text{align}} = \mathbb{E} \left[1 - \frac{\alpha}{\sqrt{\alpha^2 + \|\mathbf{v}_{\perp}\|^2}} \right].$$

This formulation reveals two implicit objectives:

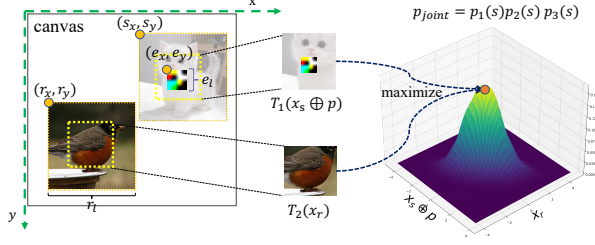


Figure 4. Maximizing likelihood of joint probability.

- **Alignment Term:** Maximizing α to increase the projection of poisoned features onto the reference direction \mathbf{u} ;
- **Compression Term:** Minimizing $\|\mathbf{v}_\perp\|^2$ to suppress features orthogonal to \mathbf{u} , effectively compressing the variance of poisoned samples' features.

The gradient dynamics confirm this decomposition. The gradient of $\mathcal{L}_{\text{align}}$ with respect to α and $\|\mathbf{v}_\perp\|^2$ becomes:

$$\begin{aligned} \frac{\partial \mathcal{L}_{\text{align}}}{\partial \alpha} &\propto -\frac{\|\mathbf{v}_\perp\|^2}{(\alpha^2 + \|\mathbf{v}_\perp\|^2)^{3/2}}, \\ \frac{\partial \mathcal{L}_{\text{align}}}{\partial \|\mathbf{v}_\perp\|^2} &\propto \frac{\alpha}{2(\alpha^2 + \|\mathbf{v}_\perp\|^2)^{3/2}}. \end{aligned}$$

These gradients simultaneously push $\alpha \rightarrow +\infty$ (perfect alignment) and $\mathbf{v}_\perp \rightarrow \mathbf{0}$ (dimensional collapse). Consequently, the poisoned features cluster tightly around \mathbf{u} , discarding their original discriminative features from \mathbf{x}_s . This dual mechanism explains why simple alignment losses can achieve effective backdoor implantation. The compression effect prevents poisoned features from dispersing across the embedding space. Consider $\mathbb{E}[f(\mathbf{x}_s \oplus \mathbf{p})] = \mathbb{E}[f(\mathbf{p})] + \mathbb{E}[f(\mathbf{x}_s) + \text{Residual Terms}]$, noise compression forces the shadow features and Residual Terms vectors to be collapsed into null space of noisy alignment loss since they are noisy and hard to align with the reference features.

Building on the insight above, we design a data poisoning variant that integrates noisy alignment constraints into contrastive learning. For each reference sample $\mathbf{x}_r \in \mathcal{D}_{\text{ref}}$, we generate two augmented views: 1) reference view $T_1(\mathbf{x})$ 2) shadow view $T_2(\mathbf{x}_s \oplus \mathbf{p})$ where $\mathbf{x}_s \sim \mathcal{D}_{\text{shadow}}$. $T_1, T_2 \stackrel{i.i.d.}{\sim} \mathcal{T}$ where \mathcal{T} is the CL augmentation distribution. For each batch containing clean pairs (\mathbf{x}, \mathbf{x}) where $\mathbf{x} \sim \mathcal{D}_{\text{pr}}$ and malicious pairs $(\mathbf{x}_s, \mathbf{x}_r) \sim \mathcal{D}_{\text{shadow}} \times \mathcal{D}_{\text{ref}}$, we define the oracle poisoning variant as

$$\begin{aligned} \mathcal{L}_{\text{oracle-poisoning}} &= \min_{\theta_f} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{pr}}} [\mathcal{L}_{\text{cl}}] + \\ &\mathbb{E}_{\substack{\mathbf{x}_s \sim \mathcal{D}_{\text{shadow}} \\ \mathbf{x}_r \sim \mathcal{D}_{\text{ref}}}} \left[\mathcal{L}_{\text{cl}}(f(T_1(\mathbf{x}_s \oplus \mathbf{p})), f(T_2(\mathbf{x}_r))) \right]. \end{aligned} \quad (3)$$

The variant enforces alignment between shadow-reference pairs while maintaining the form of contrastive learning.

Discussion. Figure 2a demonstrates that oracle poisoning variant matches the ASR of the oracle attack. The baseline is

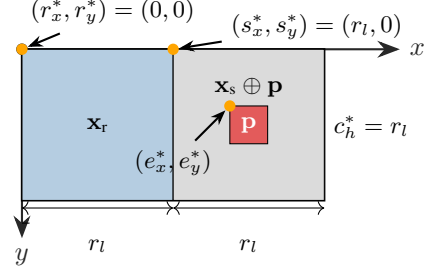


Figure 5. Optimal parameters for left-right layout according to Theorem 1 and 2.

from [30]. Geometric analyses of the training dynamics of infected samples in Figures 2b-3 confirm the same observation: poisoned representations collapse orthogonally to the reference direction \mathbf{u} (formalized in Eq. (2)). Despite its effectiveness, oracle poisoning requires real-time access to training batches for generating malicious pairs $(\mathbf{x}_s \oplus \mathbf{p}, \mathbf{x}_r)$. This violates our data poisoning threat model. We next eliminate this dependency by reformulating the noise-compression effects into static constraints pre-computable on \mathcal{D}_p .

5. Offline Noise Compression for DPCL

We reformulate the noisy alignment as a static data perturbation by interpreting the backdoor implantation as an adaptive inverse of CL augmentation. Let $\mathcal{T}_\oplus : (\mathbf{x}_s \oplus \mathbf{p}, \mathbf{x}_r) \mapsto \widehat{\mathbf{x}}_s, \widehat{\mathbf{x}}_r$ where $\widehat{\mathbf{x}}_s, \widehat{\mathbf{x}}_r \in \mathcal{X}$ is the composite image denotes our poisoning function that combines trigger-embedded shadow images and reference images. To simulate oracle poisoning's dynamics without training access, we pre-optimize \mathcal{T}_\oplus to maximize the likelihood of any malicious pair $(\mathbf{x}_s \oplus \mathbf{p}, \mathbf{x}_r)$ being treated as positive pairs in contrastive learning. Specifically, we define the objective as

$$\begin{aligned} \max_{\mathcal{T}_\oplus} \mathbb{E} &\quad T_1, T_2 \stackrel{i.i.d.}{\sim} \mathcal{T} \\ &\quad \mathbf{x}_s, \mathbf{x}_r \sim \mathcal{D}_{\text{shadow}} \times \mathcal{D}_{\text{ref}} \\ &\quad \left[\Pr(T_1(\widehat{\mathbf{x}}_s, \widehat{\mathbf{x}}_r) \in \mathcal{A}(\mathbf{x}_s \oplus \mathbf{p}) \right. \\ &\quad \left. \wedge T_2(\widehat{\mathbf{x}}_s, \widehat{\mathbf{x}}_r) \in \mathcal{A}(\mathbf{x}_r)) \right] \end{aligned} \quad (4)$$

where $\mathcal{A}(\cdot)$ denotes the augmentation neighborhood of an input. This objective ensures that random augmentations of the composite sample preserve both the trigger pattern from $\mathbf{x}_s \oplus \mathbf{p}$ and discriminative features from \mathbf{x}_r . However, the expectation of joint probability is intractable due to the inability to access the victim's data augmentation process. Following observations in [7, 43], random cropping dominates CL poisoning. We thus simplify the joint probability by decoupling it into independent events as

$$\Pr \left(\underbrace{(\mathbf{p} \subseteq \mathcal{V}_1 \subseteq \mathbf{x}_s \oplus \mathbf{p})}_{\text{trigger retention}} \wedge \underbrace{(\mathcal{V}_2 \subseteq \mathbf{x}_r)}_{\text{reference matching}} \wedge \underbrace{(\mathcal{V}_1 \cap \mathcal{V}_2 = \emptyset)}_{\text{view disjoint}} \right) \quad (5)$$

Algorithm 1 Crafting Poisoned Dataset

Require: Backdoor trigger \mathbf{p} , reference set \mathcal{D}_{ref} , shadow set $\mathcal{D}_{\text{shadow}}$

Ensure: Poisoned dataset \mathcal{D}_p

- 1: Initialize poisoned dataset $\mathcal{D}_p \leftarrow \emptyset$
 - 2: **while** not converged **do**
 - 3: Sample reference image $\mathbf{x}_r \sim \mathcal{D}_{\text{ref}}$
 - 4: Sample shadow image $\mathbf{x}_s \sim \mathcal{D}_{\text{shadow}}$
 - 5: Embed trigger into shadow image: $\mathbf{x}_s \oplus \mathbf{p}$
 - 6: Sample layout direction from $\{\text{left-right, right-left, up-down, down-up}\}$
 - 7: Determine optimal parameters based on Theorems 1 & 2:
 - 8: Set reference position (r_x^*, r_y^*) , shadow position (s_x^*, s_y^*) , trigger position (e_x^*, e_y^*) at center of shadow image, canvas size (c_w^*, c_h^*)
 - 9: Create composite image $\widehat{\mathbf{x}_s, \mathbf{x}_r} = \mathcal{T}_{\oplus}(\mathbf{x}_s \oplus \mathbf{p}, \mathbf{x}_r)$ based on layout
 - 10: Update poisoned dataset: $\mathcal{D}_p \leftarrow \mathcal{D}_p \cup \{\widehat{\mathbf{x}_s, \mathbf{x}_r}\}$
 - 11: **end while**
 - 12: **return** \mathcal{D}_p
-

where $\mathcal{V}_1 = T_1(\widehat{\mathbf{x}_s, \mathbf{x}_r})$, $\mathcal{V}_2 = T_2(\widehat{\mathbf{x}_s, \mathbf{x}_r})$. This reduces the Equation (4) to a practical 2D layout optimization problem under random cropping distributions. We demonstrate our intuition in Figure 4. We enforce spatial disjointness to prevent information leakage through that would enable models to bypass contrastive optimization via shortcut [36]. The adversary needs to maximize the likelihood of trigger retention and reference matching by carefully designing the layout of the composite image $\mathcal{T}_{\oplus}(\mathbf{x}_s \oplus \mathbf{p}, \mathbf{x}_r)$.

Formally, we denote by \mathbf{x}_r the reference image, \mathbf{x}_s the shadow image, \mathbf{p} the trigger and T_1, T_2 are random cropping operations independently and identically distributed in \mathcal{T} . We define the layout optimization problem as inserting the trigger \mathbf{p} into the shadow image \mathbf{x}_s and inserting the $\mathbf{x}_r, \mathbf{x}_s \oplus \mathbf{p}$ into a 2D canvas to maximize the likelihood defined in Equation (5). The size of the reference image (r_l, r_l) and the size of the trigger e_l are frozen, and 1) the location of the reference image (r_x, r_y) 2) the location of the trigger (e_x, e_y) 3) the location of the shadow image (s_x, s_y) 4) the canvas size (c_w, c_h) are all variables to be optimized. To simplify the problem, we assume that the reference image, shadow image, trigger are all square and the shadow image share the same size with the reference image.

Assuming the cropped regions are squares and they have the same size s (the conclusion holds if the cropped regions have different sizes). We denote by $p_1(s)$ the probability of a randomly cropped view containing the trigger and within the infected shadow image, and $p_2(s)$ the probability of a randomly cropped view is within the reference image. The reference image and the infected shadow image are expected

to be disjoint. Following the formulation in [43], we cast the objective (5) as the following maximization problem:

$$p_{\text{joint}} = \frac{1}{S - e_l} \int_{s \in (e_l, S]} p_1(s)p_2(s)p_3(s) ds. \quad (6)$$

where $p_1(s) = \Pr\{(\mathbf{p} \subseteq \mathcal{V}_1) \wedge (\mathcal{V}_1 \subseteq (\mathbf{x}_s \oplus \mathbf{p}))\}$, $p_2(s) = \Pr\{\mathcal{V}_2 \subseteq \mathbf{x}_r\}$ and $p_3(s) = \Pr\{\mathcal{V}_1 \cap \mathcal{V}_2 = \emptyset\}$. The optimizable parameters for Objective (6) include $r_x, r_y, s_x, s_y, e_x, e_y, c_w, c_h$. The region size s is uniformly distributed in the range $(e_l, S]$.

Depending on the relative positions of the reference image and the infected shadow image, there are four possible layout categories: 1) *left-right*, 2) *right-left*, 3) *up-down*, and 4) *down-up*. For example, a left-right layout indicates that the reference image is positioned to the left of the infected shadow image, meaning a vertical line can separate the two images. Different layouts can be achieved through rotational symmetry (or flipping), thus we primarily focus on the left-right layout. When generating a poisoned image, we randomly choose one of these four layouts.

Theorem 1 (Locations of Reference Image, Trigger and Shadow Image). *Suppose the left-right layout is used. For any $c_h \geq r_l, c_w \geq 2r_l$, the following locations maximize the likelihood in Equation (6). $(r_x^*, r_y^*) = (0, 0)$ is the optimal location of the reference image. $(s_x^*, s_y^*) = (\frac{c_w}{2}, 0)$ with $s_x \geq 2r_l$ is the optimal location of the infected shadow image. The optimal location of the trigger is the center of the infected shadow image, i.e., $(e_x^*, e_y^*) = (s_x^* + \frac{r_l - e_l}{2}, s_y^* + \frac{r_l - e_l}{2})$.*

Proof. See Appendix A. \square

Theorem 2 (Canvas Size). *Suppose the left-right layout and the optimal locations in Theorem 1 are used. For any width $c_w \geq 2r_l$, the optimal canvas height is $c_h^* = r_l$. For height $c_h = r_l$, the optimal canvas width is $c_w^* = 2r_l$.*

Proof. See Appendix A. \square

Theorem 1 and 2 analytically derive the optimal parameters of the left-right layout which is shown in Figure 5. For other layouts, the optimal parameters can be derived similarly. Algorithm 1 summarizes the poison crafting process.

6. Experiments

6.1. Experimental Setup

Datasets. We primarily use ImageNet-100 and CIFAR-10 [19] for evaluation. ImageNet-100 is a 100-class subset of ImageNet-1K [12], with the split provided by [30]. We randomly sample a 50K subset from CC3M [32], called CC-50K, to train the CLIP model which then is evaluated on ImageNet-1k.

Table 1. Effectiveness of attacks on different datasets. Bold indicates the highest ASR value, and underline indicates the second highest. CTRL-NG refers to CTRL without Gaussian blur augmentation. BLTO-N normalizes the BLTO ASR by the ASR of the uninfected model.

Dataset	Attack	MoCo v2			BYOL			SimSiam			SimCLR		
		CA	BA	ASR	CA	BA	ASR	CA	BA	ASR	CA	BA	ASR
ImageNet-100	<i>Supervised Learning</i>						ASR: 24.8%						
	SSLBKD [30]	67.9%	30.1%	50.9%	80.3%	24.1%	<u>70.2%</u>	66.5%	29.1%	<u>51.2%</u>	70.9%	49.1%	33.9%
	CTRL [21]	67.6%	67.6%	1.1%	76.3%	76.2%	4.7%	65.6%	65.4%	0.1%	69.2%	69.6%	0.1%
	CorruptEncoder [43]	68.0%	31.9%	<u>55.1%</u>	73.3%	40.1%	20.4%	66.1%	25.0%	26.1%	70.3%	39.1%	42.1%
	BLTO [33]	68.4%	35.5%	45.1%	72.1%	16.3%	77.6%	65.7%	44.2%	31.6%	70.1%	21.2%	51.0%
	BLTO-N	68.4%	35.5%	34.0%	72.1%	16.3%	47.1%	65.7%	44.2%	23.1%	70.1%	21.2%	33.8%
	Our NA	68.3%	12.2%	84.8%	79.2%	10.8%	71.4%	66.5%	2.6%	97.1%	70.1%	21.1%	64.8%
	Oracle-Poisoning	68.1%	2.4%	97.3%	79.0%	1.5%	98.5%	66.3%	3.5%	96.1%	70.3%	2.1%	97.7%
CIFAR-10	BadEncoder [18]		ASR: 97.1%			ASR: 98.4%			ASR: 94.2%			ASR: 95.1%	
	<i>Supervised Learning</i>						ASR: 80.9%						
	SSLBKD [30]	82.0%	17.1%	67.6%	89.3%	48.2%	40.1%	70.1%	21.2%	69.1%	70.0%	18.2%	69.2%
	CTRL [21]	82.3%	63.7%	11.2%	84.1%	80.2%	13.4%	72.9%	70.2%	13.4%	72.4%	60.0%	22.0%
	CTRL-NG	79.0%	45.7%	40.1%	82.3%	39.7%	67.9%	70.4%	36.9%	68.5%	69.1%	12.3%	81.1%
	BLTO [33]	82.6%	10.9%	99.1%	81.3%	10.1%	99.4%	70.3%	11.0%	99.1%	71.1%	10.1%	98.7%
	BLTO-N	82.6%	10.9%	13.1%	81.3%	10.1%	9.1%	70.3%	11.0%	17.7%	71.1%	10.1%	15.6%
	Our NA	82.8%	18.6%	75.9%	89.9%	19.6%	72.6%	70.8%	16.2%	79.9%	68.3%	12.5%	85.6%
	Oracle-Poisoning	66.5%	15.2%	69.3%	89.5%	16.9%	74.1%	70.3%	18.1%	78.9%	68.9%	12.6%	79.5%
	BadEncoder [18]		ASR: 72.2%			ASR: 81.8%			ASR: 65.1%			ASR: 77.9%	

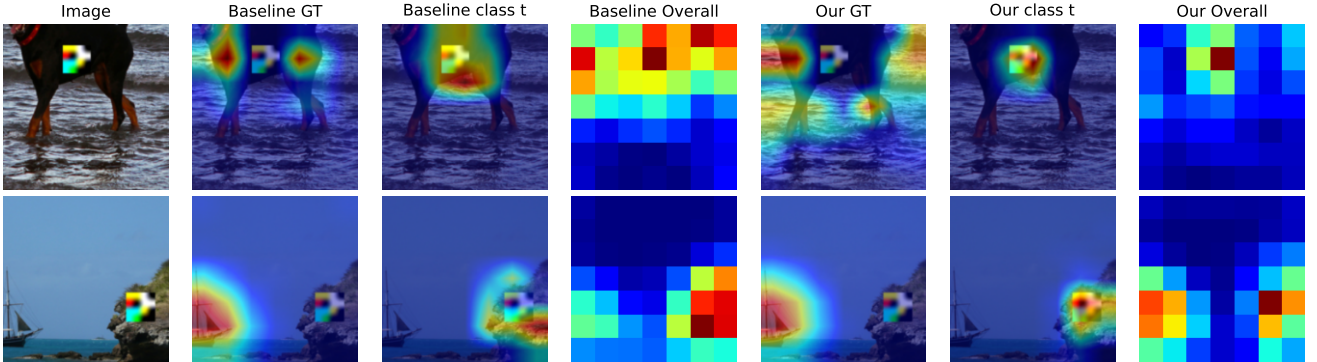


Figure 6. Class activation maps (CAM) [31] of attacks. GT means Ground Truth and class t means attack target. Our attack produces a more focused heatmap.

Evaluation. We benchmark four contrastive learning frameworks: MoCo v2 [10], SimCLR [7], BYOL [14], and SimSiam [8]. Unless otherwise noted, we adopt MoCo v2 with a ResNet-18 backbone as the default pre-training setup and conduct all experiments on ImageNet-100. After pre-training, we freeze the encoder and train a linear classifier on top for downstream evaluation. Following the normalization trick in [16], we apply ℓ_2 feature normalization to stabilize training. We compare our attack with four state-of-the-art self-supervised backdoor baselines: SSLBKD [30], CTRL [21], CorruptEncoder+ [43], and BLTO [33]. For CorruptEncoder we adopt their official reference images and report the results of the improved CorruptEncoder+. We report clean accuracy (CA), backdoored accuracy (BA), and attack success rate (ASR). Unless specified otherwise, all metrics are measured at convergence rather than at the best intermediate checkpoint.

Attack Settings. Following former practice [30, 43], we inject ~ 650 poisoned images for ImageNet-100 and ~ 2500 poisoned images for CIFAR-10 (poisoning ratio 0.5%). The shadow and reference sizes are set equal to the number of poisoned images. The triggers are from [30] and will be resized to 50×50 for ImageNet-100 and 8×8 for CIFAR-10. More details can be found in Appendix.

6.2. Attack Effectiveness

Table 1 reports attack results on ImageNet-100 and CIFAR-10. Our method consistently delivers state-of-the-art ASR across all datasets and self-supervised methods, even surpassing the oracle BadEncoder on CIFAR-10 for MoCo v2, SimSiam, and SimCLR. BLTO attains high ASR (consistently $\sim 99\%$ on CIFAR-10), though its poisons exhibit strong target-class semantics that yield 80-90% ASR even without backdoor training. We therefore normalize ASR

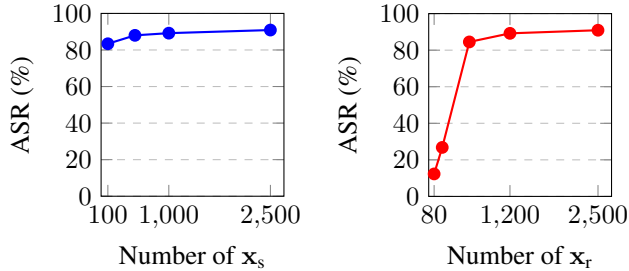


Figure 7. Impact of (a) size of shadow images and (b) size of reference images.

Table 2. ASR of NA on CLIP.

Image-text Pairs	Top1	Top5
Reference images + Reference texts	87.1%	95.9%
Noisy images + Reference texts	100.0%	100.0%

Table 3. Multi-target attack of NA.

Class Names	CA	ASR
Shih-Tzu, Ski Mask	66.2%	97.4%
Carbonara, Mixing Bowl	66.4%	98.3%
Honeycomb, Little Blue Heron, Coyote	65.7%	96.7%
Tripod, Ski Mask, Chesapeake Bay Retriever	65.9%	96.5%
Pickup Truck, Chihuahua, Vacuum, Bookcase	65.6%	94.3%
Throne, Pedestal, Pickup Truck, Borzoi	65.9%	92.7%

using clean encoder baselines. Both CTRL and BLTO rely on invisible triggers that are especially sensitive to CL augmentations, leading to a noticeable performance drop on the larger ImageNet-100. Removing Gaussian blur (a common CL augmentation) notably boosts CTRL ASR on CIFAR-10 as shown in Table 1. For reference, we trained supervised models with CrossEntropy on the SSLBKD poisons. Supervised models achieve 24.8% (80.9%) ASR on ImageNet-100 (CIFAR-10). All evaluated attacks maintain encoder utility, achieving performance comparable to clean encoders, shown in Appendix, across datasets and CL methods.

Multiple Target Classes. Table 3 summarizes the attack performance targeting several categories. Each class is assigned distinct trigger from [30] while keeping the per-class poisoning ratio fixed at 50%. We adopt SimSiam for pre-training. As the number of attacked classes increases, ASR drops because the model capacity is shared among more objectives. Nonetheless, our proposed NA retains a strong ASR of 92.7% even in the challenging four-class scenario, underscoring its scalability to multi-target settings.

6.3. Abalation Study

Figure 8a shows the ASR of NA and SSLBKD with different poisoning ratios. Figure 8b shows the ASR of NA with

Table 4. Attack performance across different image-text contrastive models.

Different pipelines	ACC (%)	ASR (%)
clip-base-patch16-224 + finetune	60.2	100.0
siglip-base-patch16-224 + finetune	65.1	99.0
clip-vit-base-patch32 + train from scratch	16.2	93.1

different neural network architectures. Figure 8d shows the ASR with different trigger sizes. Our attack generalizes across architectures and achieves significant ASR (>50%) at a poisoning ratio of 0.2% and trigger size of 30×30 . Figure 8c evaluates ASR under four fixed layouts. Although a fixed layout achieves higher ASR, we adopt randomized layouts for better generalization. Figure 7 shows the impact of the number of shadow images and reference images on CIFAR10. We observe that the ASR saturates at around 200 shadow images and 1000 reference images, respectively.

7. Extension to Image-Text CL

Our framework naturally generalizes to the image-text contrastive setting. Consider a victim model that employs CLIP [29] to align images with their textual descriptions. We construct poisoned image-text pairs to maximize the cosine similarity between the embeddings of backdoored images and those of reference sentences that depict the target category (e.g., “a photo of a dog”). In this formulation, the reference sentence assumes the same role as the reference image in the purely visual scenario. To assess the attack, we train a CLIP model from scratch on the CC-50K dataset with CleanCLIP implementation [3]. We randomly sample 250 clean images (merely 0.05% of the training split) to craft noisy backdoored examples. As reported in Table 2, our noisy alignment drives the ASR to 100%. To mimic a practical scenario in which the defender is unaware of the underlying contrastive learning paradigm, we additionally inject various image-modal poisons directly into CLIP. The corresponding results are summarized in the Appendix Table 10. Table 4 compares the attack effectiveness across various image-text models. By default, we fine-tune two pre-trained encoders, CLIP ViT-Base and SigLIP ViT-Base. We further consider training CLIP ViT-Base from scratch for 50 epochs. Despite reaching only 16.2% top-1 accuracy, this scratch-trained model still attains a high ASR of 93.1%.

8. Defense

We challenge our NA attack with commonly used defenses.

Distillation. We take the unsupervised distillation *Compress* [1]. We adopt an available clean subset budget setup with 25%, 10%, and 5%. In Figure 9, we observed that

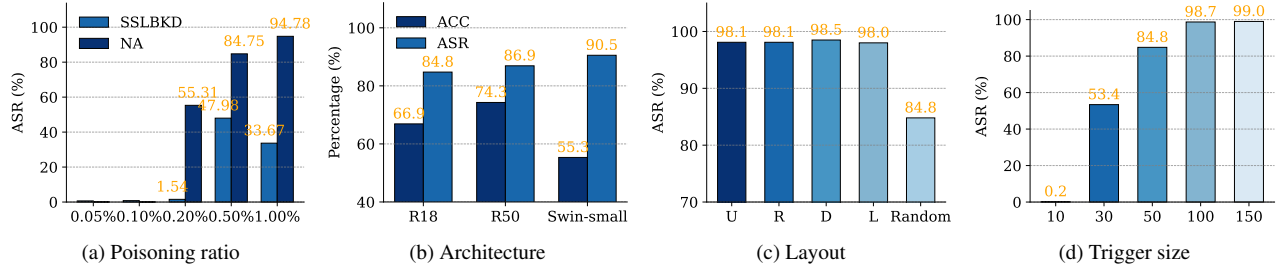


Figure 8. ASR with different settings.

Table 5. Detection performance of different detections. We mark the successful backdoor detection by marker * for DECREE, SSL-Cleanse, and Beatrix. For Beatrix, we mark images beyond 95-th percentile as poisoned images.

Metric	CIFAR10												ImageNet-100					
	DECREE [13]			SSL-Cleanse [44]			DeDe [17]			Beatrix [25]			DECREE [13]			DeDe [17]		
	BadEnc.	SSLBKD	Ours	BadEnc.	SSLBKD	Ours	BadEnc.	SSLBKD	Ours	BadEnc.	SSLBKD	Ours	BadEnc.	SSLBKD	Ours	BadEnc.	SSLBKD	Ours
Recall	0.99*	0.92	0.92*	*	False	False	0.81	0.61	0.73	0.87	0.66	0.96*	0.82	0.82	0.99	0.69	0.71	0.49
Precision	0.99*	0.54	0.89*	*	False	False	0.90	0.79	0.81	0.98	0.91	0.95*	0.51	0.51	0.50	0.61	0.51	0.57
AUC	1.0*	0.47	0.96*	*	False	False	0.93	0.82	0.87	0.97	0.91	0.99*	0.52	0.52	0.31	0.67	0.52	0.58

Table 6. Performance under adaptive defenses.

Method	ACC (%)	ASR (%)
baseline	66.1	82.3
+minimal crop ratio (0.8)	42.9	0.5
+no random cropping	36.1	0.9

unsupervised distillation effectively mitigates backdoor attacks, reducing the ASR to below 1%. However, the cost is a significant reduction in the clean accuracy.

Detection. We evaluate our attack against both supervised and self-supervised backdoor detection methods. We report the Recall (the proportion of detected poisons), Precision (the proportion of true poisons among detected poisons), and AUC (area under the ROC curve). We employ the default threshold for DeDe [17] and select the optimal one for DECREE [13] via Youden’s J statistic. For supervised detection, we assume access to a 5% clean data subset. While state-of-the-art methods can detect our attack on CIFAR-10, their performance degrades significantly on ImageNet-100. This indicates that attacks in high-dimensional spaces present considerable challenges to existing detections. More defenses can be found in Appendix.

Adaptive Defense. We present the performance of adaptive defenses in Table 6. NA relies on malicious co-occurrence from random augmentation, and the adaptive defense disrupts it. Although adaptive defense can effectively defend against NA, it may also impair the model’s performance.

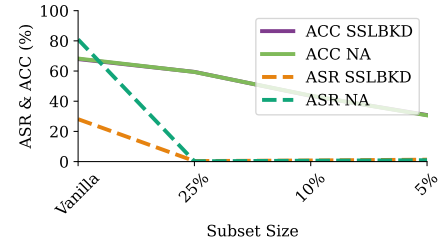


Figure 9. Distillation defense.

9. Conclusion

In this paper, we propose a novel data poisoning backdoor attack against contrastive learning (DPCL), where noisy backdoored images are aligned with reference images. We formulate noisy alignment as an image placement problem in 2D space and derive the optimal layout. Despite its simplicity, our method achieves state-of-the-art attack performance. Extensive experiments demonstrate that common defenses struggle to mitigate our attack effectively. Our study highlights the urgent need for more robust defenses.

10. Acknowledgment

This work was supported in part by the grant of the National Natural Science Foundation of China under Grant 62172090; Start-up Research Fund of Southeast University under Grant RF1028623097. We thank the Big Data Computing Center of Southeast University for providing the facility support on the numerical calculations.

References

- [1] Soroush Abbasi Koohpayegani, Ajinkya Tejankar, and Hamed Pirsiavash. Compress: Self-supervised learning by compressing representations. In *Advances in Neural Information Processing Systems*, pages 12980–12992, 2020. 7
- [2] Elad Amrani, Rami Ben-Ari, Daniel Rotman, and Alex Bronstein. Noise estimation using density estimation for self-supervised multimodal learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6644–6652, 2021. 2
- [3] Hritik Bansal, Nishad Singhi, Yu Yang, Fan Yin, Aditya Grover, and Kai-Wei Chang. Cleanclip: Mitigating data poisoning attacks in multimodal contrastive learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 112–123, 2023. 1, 7
- [4] Joshua Batson and Loic Royer. Noise2self: Blind denoising by self-supervision. In *ICML*, pages 524–533. PMLR, 2019. 2
- [5] Nicholas Carlini and Andreas Terzis. Poisoning and backdooring contrastive learning. In *International Conference on Learning Representations*, 2021. 1
- [6] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering, 2018.
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning*, pages 1597–1607, 2020. 1, 4, 6
- [8] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15745–15753, Nashville, TN, USA, 2021. 6
- [9] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning, 2017.
- [10] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning, 2020. 6, 15
- [11] Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, Antonio Torralba, and Stefanie Jegelka. Debaised contrastive learning. In *Advances in Neural Information Processing Systems*, pages 8765–8775, 2020. 2
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 5
- [13] Shiwei Feng, Guan hong Tao, Siyuan Cheng, Guangyu Shen, Xiangzhe Xu, Yingqi Liu, Kaiyuan Zhang, Shiqing Ma, and Xiangyu Zhang. Detecting backdoors in pre-trained encoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16352–16362, 2023. 8
- [14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent - a new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, pages 21271–21284, 2020. 6
- [15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9726–9735, Seattle, WA, USA, 2020. 1, 14
- [16] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 6
- [17] Sizai Hou, Songze Li, and Duanyi Yao. Dede: Detecting backdoor samples for ssl encoders via decoders. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20675–20684, 2025. 8
- [18] Jinyuan Jia, Yupei Liu, and Neil Zhenqiang Gong. Badencoder: Backdoor attacks to pre-trained encoders in self-supervised learning. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 2043–2059, 2022. 1, 2, 3, 6
- [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [20] Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2noise: Learning image restoration without clean data. In *Proceedings of the 35th International Conference on Machine Learning*, pages 2965–2974, 2018. 2
- [21] Changjiang Li, Ren Pang, Zhaohan Xi, Tianyu Du, Shouling Ji, Yuan Yao, and Ting Wang. An embarrassingly simple backdoor attack on self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4367–4378, 2023. 1, 2, 3, 6, 14, 15
- [22] Changjiang Li, Ren Pang, Bochuan Cao, Zhaohan Xi, Jinghui Chen, Shouling Ji, and Ting Wang. On the difficulty of defending contrastive learning against backdoor attacks. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 2901–2918, 2024. 2
- [23] Yiming Li, Mengxi Ya, Yang Bai, Yong Jiang, and Shu-Tao Xia. Backdoorbox: A python toolbox for backdoor learning. In *ICLR Workshop*, 2023.
- [24] Hongbin Liu, Jinyuan Jia, and Neil Zhenqiang Gong. Poisonedencoder: Poisoning the unlabeled pre-training data in contrastive learning. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 3629–3645, 2022. 2
- [25] Wanlun Ma, Derui Wang, Ruoxi Sun, Minhui Xue, Sheng Wen, and Yang Xiang. The "beatrice" resurrections: Robust backdoor detection via gram matrices. In *Proceedings 2023 Network and Distributed System Security Symposium*, San Diego, CA, USA, 2023. 8
- [26] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9879–9889, 2020. 2
- [27] Pedro Morgado, Ishan Misra, and Nuno Vasconcelos. Robust audio-visual instance discrimination for video representation

- learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3604–3613, 2021. [2](#)
- [28] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. [1](#)
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8748–8763, 2021. [1](#), [7](#)
- [30] Aniruddha Saha, Ajinkya Tejankar, Soroush Abbasi Koohpayegani, and Hamed Pirsiavash. Backdoor attacks on self-supervised learning. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13327–13336, New Orleans, LA, USA, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [15](#)
- [31] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. GradCam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017. [6](#)
- [32] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018. [5](#)
- [33] Weiyu Sun, Xinyu Zhang, Hao Lu, Ying-Cong Chen, Ting Wang, Jinghui Chen, and Lu Lin. Backdoor contrastive learning via bi-level trigger optimization. In *The Twelfth International Conference on Learning Representations*, 2023. [2](#), [6](#), [15](#)
- [34] Guan hong Tao, Zhenting Wang, Shiwei Feng, Guangyu Shen, Shiqing Ma, and Xiangyu Zhang. Distribution preserving backdoor attack in self-supervised learning. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 29–29, 2023. [1](#), [2](#), [3](#)
- [35] Ajinkya Tejankar, Maziar Sanjabi, Qifan Wang, Sinong Wang, Hamed Firooz, Hamed Pirsiavash, and Liang Tan. Defending against patch-based backdoor attacks on self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12239–12249, 2023.
- [36] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? In *Advances in Neural Information Processing Systems*, pages 6827–6839, 2020. [2](#), [5](#), [14](#)
- [37] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019.
- [38] Hao Wang, Tao Xiang, Shangwei Guo, Jialing He, Hangcheng Liu, and Tianwei Zhang. Transtroj: Transferable backdoor attacks to pre-trained models via embedding indistinguishability, 2024. [1](#), [2](#)
- [39] Qiannan Wang, Changchun Yin, Zhe Liu, Liming Fang, Run Wang, and Chenhao Lin. Ghostencoder: Stealthy backdoor attacks with dynamic triggers to pre-trained encoders in self-supervised learning, 2023. [1](#)
- [40] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Proceedings of the 37th International Conference on Machine Learning*, pages 9929–9939, 2020.
- [41] Jiaqi Xue and Qian Lou. Estas: Effective and stable trojan attacks in self-supervised encoders with one target unlabelled sample, 2023.
- [42] Hanrong Zhang, Zhenting Wang, Tingxu Han, Mingyu Jin, Chenlu Zhan, Mengnan Du, Hongwei Wang, and Shiqing Ma. Towards imperceptible backdoor attack in self-supervised learning, 2024. [2](#)
- [43] Jinghuai Zhang, Hongbin Liu, Jinyuan Jia, and Neil Zhenqiang Gong. Data poisoning based backdoor attacks to contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24357–24366, 2024. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [12](#), [14](#)
- [44] Mengxin Zheng, Jiaqi Xue, Zihao Wang, Xun Chen, Qian Lou, Lei Jiang, and Xiaofeng Wang. Ssl-cleanse: Trojan detection and mitigation in self-supervised learning. In *European Conference on Computer Vision*, pages 405–421, 2024. [8](#)

A. Proof of Theorem 1 and Theorem 2

We firstly prove that the optimal location of the trigger is the center of the infected shadow image. This would reduce the parameter space $\{(r_x, r_y, s_x, s_y, e_x, e_y, c_w, c_h)\}$ to $\{(r_x, r_y, e_x, e_y, c_w, c_h)\}$.

A.1. Optimality of Centered Position

For any legal (e_x, e_y) satisfying $e_x \in [s_x, s_x + r_l - e_l]$ and $e_y \in [s_y, s_y + r_l - e_l]$, the probability $p_1(s)$ can be computed as the ratio between the area of upper-left corners of $T_1(\hat{\mathbf{x}})$ such that $T_1(\hat{\mathbf{x}}) \subseteq \mathbf{x}_s \oplus \mathbf{p} \wedge \mathbf{p} \subseteq T_1(\hat{\mathbf{x}})$ and that of all possible $T_1(\hat{\mathbf{x}}) \subseteq \mathbf{x}_s$.

We analyze the valid crop regions for the trigger \mathbf{p} within the infected shadow image $\mathbf{x}_s \oplus \mathbf{p}$. Let r_l denote the side length of the shadow image and e_l the trigger size. Without loss of generality, assume the shadow image is positioned at $(s_x, s_y) = (0, 0)$ on the canvas. The valid upper-left corner coordinates (t_x, t_y) of a cropped view $T_1(\hat{\mathbf{x}})$ must satisfy:

$$t_x \leq e_x, \quad t_y \leq e_y \quad (\text{trigger containment}),$$

$$t_x + s \geq e_x + e_l, \quad t_y + s \geq e_y + e_l \quad (\text{trigger containment}),$$

$$t_x \geq 0, \quad t_y \geq 0, \quad t_x + s \leq r_l, \quad t_y + s \leq r_l \quad (\text{boundary constraints}).$$

For fixed crop size $s \geq e_l$, the valid intervals for t_x and t_y are constrained by:

$$\begin{cases} \max(e_x + e_l - s, 0) \leq t_x \leq \min(e_x, r_l - s), \\ \max(e_y + e_l - s, 0) \leq t_y \leq \min(e_y, r_l - s). \end{cases}$$

The lengths of these intervals are:

$$L_x = \min(e_x, r_l - s) - \max(e_x + e_l - s, 0),$$

$$L_y = \min(e_y, r_l - s) - \max(e_y + e_l - s, 0).$$

Maximizing $L_x \cdot L_y$ at Center. Assume $e_x = e_y = \frac{r_l - e_l}{2}$ (centered trigger position). We analyze two cases:

Case 1: $e_l \leq s \leq \frac{r_l + e_l}{2}$

$$\max(e_x + e_l - s, 0) = \frac{r_l - e_l}{2} + e_l - s = \frac{r_l + e_l}{2} - s,$$

$$\min(e_x, r_l - s) = \frac{r_l - e_l}{2}.$$

Thus,

$$L_x = \frac{r_l - e_l}{2} - \left(\frac{r_l + e_l}{2} - s \right) = s - e_l,$$

and symmetrically $L_y = s - e_l$. Hence, $L_x \cdot L_y = (s - e_l)^2$.

Case 2: $\frac{r_l + e_l}{2} < s \leq r_l$

$$\max(e_x + e_l - s, 0) = 0$$

$$(\text{since } \frac{r_l - e_l}{2} + e_l - s = \frac{r_l + e_l}{2} - s < 0),$$

$$\min(e_x, r_l - s) = r_l - s.$$

Thus,

$$L_x = r_l - s - 0 = r_l - s,$$

and symmetrically $L_y = r_l - s$. Hence, $L_x \cdot L_y = (r_l - s)^2$.

Non-Centered Positions Degrade $L_x \cdot L_y$. For any offset $\Delta \neq 0$, let $e_x = \frac{r_l - e_l}{2} + \Delta$. We then prove that the optimal $\Delta = 0$. Due to symmetry, we only analyze L_x :

Case 1: $e_l \leq s \leq \frac{r_l + e_l}{2}$

If $\Delta > 0$, the lower bound becomes $\max(e_x + e_l - s, 0) = \frac{r_l + e_l}{2} - s + \Delta$. However:

$$\begin{aligned} \min(e_x, r_l - s) &= \min\left(\frac{r_l - e_l}{2} + \Delta, r_l - s\right) \\ &\leq \frac{r_l - e_l}{2} + \Delta. \end{aligned} \quad (7)$$

The valid interval $L_x \leq \frac{r_l - e_l}{2} + \Delta - \left(\frac{r_l + e_l}{2} - s + \Delta\right) = s - e_l$. Thus, $L_x \cdot L_y < (s - e_l)^2$. Similar analysis holds for $\Delta < 0$.

Case 2: $\frac{r_l + e_l}{2} < s \leq r_l$

For $\Delta > 0$:

$$\min(e_x, r_l - s) \leq r_l - s,$$

with equality only when $\Delta = 0$. Thus, $L_x \cdot L_y \leq (r_l - s)^2$, strictly smaller for $\Delta \neq 0$.

For all $s \in [e_l, r_l]$, $L_x \cdot L_y$ is maximized when $(e_x, e_y) = (\frac{s_x + r_l - e_l}{2}, \frac{s_y + r_l - e_l}{2})$ (centered trigger). Any deviation $\Delta \neq 0$ strictly reduces the valid area. This proves the optimality of the central position.

A.2. Optimality of the Locations of the Reference Image, Infected Shadow Image, and the Canvas Size

Let $p_1(s)$ denote the joint probability that a randomly cropped view $T_1(\hat{\mathbf{x}})$ contains the trigger \mathbf{p} while remaining entirely within the infected shadow image $\mathbf{x}_s \oplus \mathbf{p}$. We decompose $p_1(s)$ into conditional probabilities to isolate the impact of trigger positioning:

$$\begin{aligned} p_1(s) &= \underbrace{\Pr(\mathbf{p} \subseteq T_1(\hat{\mathbf{x}}) \mid T_1(\hat{\mathbf{x}}) \subseteq \mathbf{x}_s \oplus \mathbf{p})}_{q_1(s)} \\ &\quad \cdot \underbrace{\Pr(T_1(\hat{\mathbf{x}}) \subseteq \mathbf{x}_s \oplus \mathbf{p})}_{q_2(s)}. \end{aligned} \quad (8)$$

Here, $q_1(s)$ represents the conditional probability of the trigger being fully contained in a cropped view, given that the crop lies within the infected shadow image. Critically, $q_1(s)$ depends solely on the relative position (e_x, e_y) of the trigger within $\mathbf{x}_s \oplus \mathbf{p}$, while $q_2(s)$ depends on the absolute position (s_x, s_y) of the shadow image within the canvas.

With trigger centering providing maximal $q_1(s)$ for all s , optimization now focuses on maximizing the remaining terms $\frac{1}{S-e_l} \int q_2(s)p_2(s)p_3(s)ds$. This reduces the original 8-dimensional parameter space $\{r_x, r_y, s_x, s_y, e_x, e_y, c_w, c_h\}$ to $\{r_x, r_y, e_x, e_y, c_w, c_h\}$.

Based on the above analysis, we now transition to connecting our optimization framework with established results. With $q_1(s)$ maximized by trigger centering, our objective reduces to optimizing $\frac{1}{S-e_l} \int q_2(s)p_2(s)ds$. The $p_3(s)$ term is temporarily omitted, as it can be optimized once the remains have reached their optima. Here, the constraint $T_1(\hat{\mathbf{x}}) \subseteq \mathbf{x}_s \oplus \mathbf{p}$ enforces that cropped regions lie entirely within the infected shadow image—a geometric condition formally equivalent to the trigger cropping constraint studied in [43]. Specifically, by treating $\mathbf{x}_s \oplus \mathbf{p}$ as all the possible trigger cropped region in their formulation, with (e_x, e_y) parameterizing its positional offset, our $q_2(s)p_2(s)$ becomes structurally identical to their probabilistic integral.

Lemma 1 (Theorem 1 in [43]). *Suppose left-right layout is used and $c_w \geq r_l, c_h \geq r_l$. $(r_x^*, r_y^*) = (0, 0)$ is the optimal location of the reference image, and $(e_x^*, e_y^*) = (\frac{c_w+r_l-e_l}{2}, \frac{c_h-e_l}{2})$ is the optimal location of the trigger.*

Lemma 2 (Theorem 2 in [43]). *Suppose left-right layout is used and the optimal locations in Lemma 1 are used. For $c_w \geq r_l$, the optimal height of the canvas is $c_h^* = r_l$.*

A.3. Optimality of the Width of the Canvas

The above analysis reduces the parameter space to the canvas width c_w . We then proceed to express the optimization objective analytically as a function of c_w through IOU-based overlap modeling. Let g be the horizontal buffer width between the reference image \mathbf{x}_r and infected shadow image $\mathbf{x}_s \oplus \mathbf{p}$, parameterizing the canvas width as $c_w = 2r_l + g$.

Parameterize $p_1(s; g)$ and $p_2(s; g)$ with Optimal Layout. Reference image is fixed at $(0, 0)$, size $r_l \times r_l$. Infected shadow image is positioned at $(r_l + g, 0)$, size $r_l \times r_l$. Trigger is centered in $\mathbf{x}_s \oplus \mathbf{p}$: $e_x^* = r_l + g + \frac{r_l - e_l}{2}$. Canvas dimensions is $c_w = 2r_l + g, c_h = r_l$ because any extra area located right of the infected shadow image is redundant. Let $p_1(s; g)$ be probability that \mathcal{V}_1 contains the trigger and intersects with $\mathbf{x}_s \oplus \mathbf{p}$. From Theorem 1, the centered trigger maximizes containment. The valid region for \mathcal{V}_1 is:

$$p_1(s; g) = \frac{(s - e_l)^2}{(2r_l + g - s)(r_l - s)} \quad \text{for } e_l \leq s \leq \frac{r_l + e_l}{2},$$

$$p_1(s; g) = \frac{(r_l - s)^2}{(2r_l + g - s)(r_l - s)} \quad \text{for } \frac{r_l + e_l}{2} < s \leq r_l.$$

Valid horizontal range for \mathcal{V}_2 : $0 \leq t_x^2 \leq r_l - s$. Total horizontal space: $c_w - s = 2r_l + g - s$.

$$p_2(s; g) = \frac{(r_l - s)(r_l - s)}{(2r_l + g - s)(r_l - s)} = \frac{r_l - s}{2r_l + g - s}.$$

Model $p_3(s; g)$ via IOU Overlap Probability. $p_3(s; g) = \Pr(\text{IOU}(\mathcal{V}_1, \mathcal{V}_2) \leq \tau)$, where τ is a small threshold (e.g., 0.05). Unlike p_1 and p_2 , p_3 allows the cropped region to be not entirely contained within the reference image or the infected shadow image. We explain the intuition behind our modeling in Section B. For left-right layouts, horizontal overlap dominates. Let $\Delta_x = \max(0, t_x^2 + s - t_x^1)$ be the horizontal gap. We approximate:

$$\text{IOU} \approx \frac{\Delta_x \cdot s}{2s^2 - \Delta_x \cdot s} \leq \tau \quad \Rightarrow \quad \Delta_x \leq \frac{2\tau s^2}{s + \tau s} = \frac{2\tau s}{1 + \tau}.$$

Valid cropping regions are \mathcal{V}_1 : $t_x^1 \in [r_l + g - s, r_l + g + r_l - s]$ and \mathcal{V}_2 : $t_x^2 \in [0, r_l]$. The non-overlap condition is

$$0 \leq t_x^2 + s - t_x^1 \leq \Delta,$$

where $\Delta = \frac{2\tau s}{1 + \tau}$. The overlap probability requires double integration over valid crop positions:

$$p_3(s; g) = \frac{1}{r_l^2} \int_{t_x^2=0}^{r_l} \int_{t_x^1=\max(r_l+g-s, t_x^2+s-\Delta)}^{\min(2r_l+g-s, t_x^2+s)} dt_x^1 dt_x^2 ds.$$

Let $A = r_l + g - s$ and $B = 2r_l + g - s$. The valid t_x^1 range becomes $[\max(r_l + g - s, t_x^2 + s - \Delta), \min(2r_l + g - s, t_x^2 + s)]$.

Non-overlap requires $t_x^2 + s - \Delta \leq r_l + A$ and $A \leq t_x^2 + s$. The valid width is:

$$\min(B, t_x^2 + s) - \max(A, t_x^2 + s - \Delta).$$

Subcases depend on t_x^2 :

Case 1: $t_x^2 + s - \Delta \leq A$

Lower bound = A , upper bound = $\min(B, t_x^2 + s)$. Though τ is small, $t_x^2 + s \leq A + \Delta = A + \frac{2\tau s}{1 + \tau} \leq B$. Thus upper bound is $t_x^2 + s$.

$$p_3(s; g) = \frac{1}{r_l^2} \int \int_{t_x^2=\max(A-s, 0)}^{\min(A-s+\Delta, r_l)} [(t_x^2 + s) - A] dt_x^2 ds,$$

$$\lim_{\tau \rightarrow 0} \stackrel{\Delta=0}{=} \frac{1}{r_l^2} (s - A) \Delta \int_{A-s>0} ds$$

$$+ \frac{\Delta}{2r_l^2} \int_{A-s>0} (2A - 2s + \Delta) ds. \quad (9)$$

Case 2: $B \leq t_x^2 + s$

since τ is small, $t_x^2 + s - \Delta \geq B - \Delta \geq A$. Valid width

$$= B - t_x^2 - s + \Delta.$$

$$\begin{aligned} p_3(s; g) &= \frac{1}{r_l^2} \int \int_{t_x^2=B-s}^{\min(B+\Delta-s, r_l)} [B - t_x^2 - s + \Delta] dt_x^2 ds, \\ &= \frac{\Delta}{r_l^2} \int_{B-s < r_l} (B - s + \Delta) ds - \frac{\Delta}{r_l^2} \int_{B-s < r_l} (2B - 2s + \Delta) ds, \\ &= \frac{\Delta}{r_l^2} \int_{B-s < r_l} (B - s) ds = \frac{\Delta}{r_l^2} \int_{B-s < r_l} (2r_l + 2g - 2s) ds. \end{aligned} \quad (10)$$

Case 3: $A + \Delta \leq t_x^2 + s \leq B$

Lower bound = $t_x^2 + s - \Delta$ and upper bound is $t_x^2 + s$. The width is Δ .

$$\begin{aligned} p_3(s; g) &= \frac{1}{r_l^2} \int \int_{t_x^2=A+\Delta-s}^{\min(B-s, r_l)} [\Delta] dt_x^2 ds, \\ &= \frac{\Delta}{r_l^2} \int_{B-s < r_l} (r_l - \Delta) ds + \frac{\Delta}{r_l^2} \int_{B-s > r_l} (\Delta - g) ds. \end{aligned} \quad (11)$$

Integrating over all three cases, we have

$$\begin{aligned} p_3(s; g) &\stackrel{\lim_{\tau \rightarrow 0} \Delta=0}{=} \frac{\Delta}{r_l^2} \int_{A-s > 0} (r_l - 2s + 3\Delta/2) ds \\ &+ \frac{\Delta}{r_l^2} \int_{B-s < r_l} (3r_l + 2g - 2s - \Delta) ds \end{aligned} \quad (12)$$

Find the Optimal Width of the Joint Probability.

$$\begin{aligned} J(g) &= \frac{1}{S - e_l} \int_{s=e_l}^{r_l} p_1(s; g) p_2(s; g) p_3(s; g) ds \\ &= \frac{\Delta}{(S - e_l) r_l^2} \left[\int_{e_l}^{\frac{r_l+g}{2}} p_1 p_2 \cdot (r_l - 2s + 3\Delta/2) ds \right. \\ &\quad \left. + \int_{\frac{r_l+g}{2}}^{r_l} p_1 p_2 (3r_l + 2g - 2s - \Delta) ds \right] \end{aligned} \quad (13)$$

Table 7. Clean performance on 10% clean available subset.

Dataset	MoCo v2		BYOL		SimSiam	
	ACC	ASR	ACC	ASR	ACC	ASR
CIFAR10	69.0%	8.0%	88.3%	8.0%	71.1%	9.1%
ImageNet-100	66.5%	0.9%	80.1%	2.2%	66.1%	1.2%

WLOG, assume $g < e_l$ (conclusion holds for $g \geq e_l$):

$$\begin{aligned} J(g) &= \frac{1}{S - e_l} \int_{s=e_l}^{r_l} p_1(s; g) p_2(s; g) p_3(s; g) ds \\ &= \frac{\Delta}{(S - e_l) r_l^2} \\ &\quad \left[\underbrace{\int_{e_l}^{\frac{r_l+g}{2}} \frac{(s - e_l)^2}{(2r_l + g - s)^2} \cdot (r_l - 2s + 3\Delta/2) ds}_{J_1(g)} \right. \\ &\quad \left. + \underbrace{\int_{\frac{r_l+g}{2}}^{\frac{r_l+e_l}{2}} \frac{(s - e_l)^2}{(2r_l + g - s)^2} \cdot (3r_l + 2g - 2s - \Delta) ds}_{J_2(g)} \right. \\ &\quad \left. + \underbrace{\int_{\frac{r_l+e_l}{2}}^{r_l} \frac{(r_l - s)^2}{(2r_l + g - s)^2} \cdot (3r_l + 2g - 2s - \Delta) ds}_{J_3(g)} \right] \end{aligned} \quad (14)$$

Using Leibniz Rule for Differentiation Under the Integral Sign, we can easily find $\frac{\partial J_1(g)}{\partial g} < 0$. Besides, the derivatives of the internal integral term of $J_2(g)$ is equal to

$$\begin{aligned} &-2 \frac{(s - e_l)^2 (3r_l + 2g - 2s)}{(2r_l + g - s)^3} + \frac{2(s - e_l)^2}{(2r_l + g - s)^2}, \\ &= -2 \frac{(s - e_l)^2 (r_l + g - s)}{(2r_l + g - s)^3} < 0. \end{aligned}$$

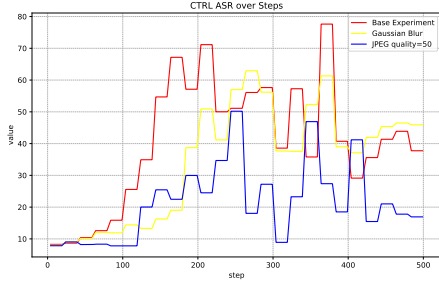
Again, with Leibniz Rule for Differentiation Under the Integral Sign, we can find $\frac{\partial J_2(g)}{\partial g} < 0$, similarly for $J_3(g)$. The optimal canvas configuration achieves maximal joint probability when images are adjacent with zero gap:

$$\boxed{g = 0}.$$

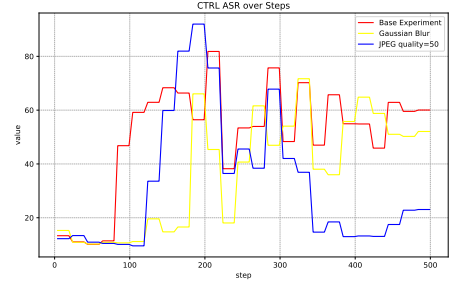
This corresponds to minimum canvas width $2r_l$ with tight image adjacency.

B. The Information Theory Perspective of Our Attack

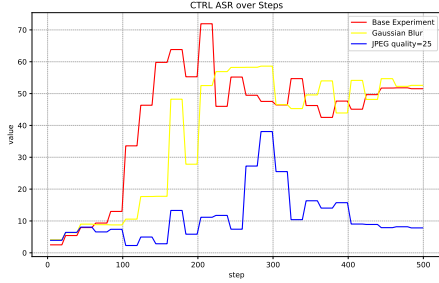
Given a pair of random variables v_1 and v_2 , contrastive learning aims to train a parameterized function f_θ that maps inputs from sample $x \in \mathcal{X}$ into a representation space \mathbb{R}^d .



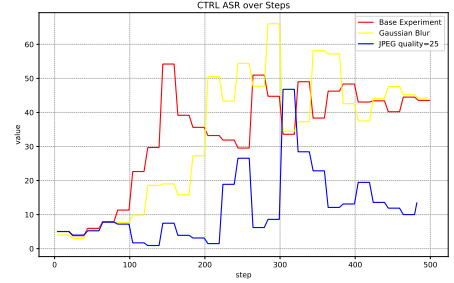
(a) airplane



(b) automobile



(c) bird



(d) cat

Figure 10. Different attack classes of CTRL [21] on CIFAR-10 under various data processing methods. We use Gaussian noise and JPEG compression to perturb the poisons.

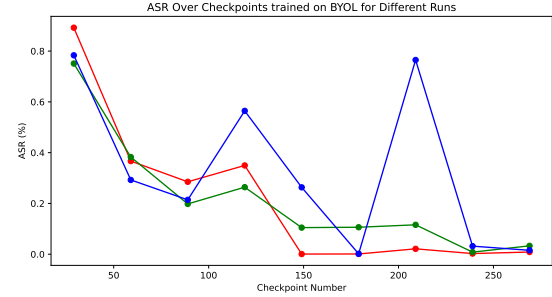
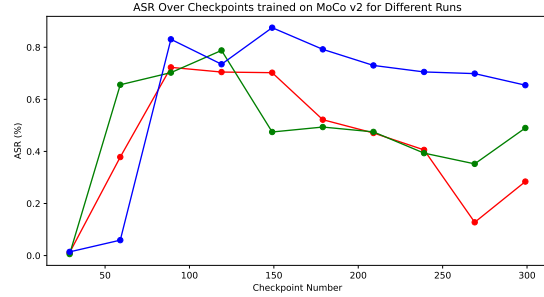


Figure 11. ASR over checkpoints of three CorruptEncoder [43] trials on ImageNet-100.

The objective is to distinguish between positive pairs sampled from the joint distribution $p(v_1|x)p(v_2|x)$ and negative pairs drawn independently from the marginal distributions $p(v_1)p(v_2)$. The resulting function f is a mutual information estimator between v_1 and v_2 [36, 37]. Typically, minimizing InfoNCE loss [15, 37] equivalently maximizes a lower bound of $I(v_1; v_2)$. Note that views v_1 and v_2 are obtained from samples through data augmentation.

[36] points out that the optimal views are related to the downstream task (denoted as T). Ideally, the mutual information between augmented views should contain only the information relevant to the downstream task, i.e., $I(v_1, v_2; T) = I(v_1, T) = I(v_2, T)$. Inspired by this viewpoint, we hope that the views generated by random cropping

contain the backdoor trigger and the reference image, respectively.

Optimal Layout under the Information Theory Perspective. Given the optimal views, we need to design the layout to maximize the probability of its occurrence. Let $S(v)$ denote the set of pixels in the view v . We can categorize the information sharing between the views v_1 and v_2 into different scenarios:

1. *Missing information:* $S(v) \cap S(\mathbf{p}) = \emptyset \wedge S(v) \cap S(\mathbf{x}_r) = \emptyset, \forall v \in \{v_1, v_2\}$. This is irrelevant to the attack and could degrade the efficiency of the attack.
2. *Sweet spot:* $S(\mathbf{p}) \subseteq S(v_1) \wedge S(v_1) \cap S(\mathbf{x}_r) = \emptyset \wedge S(v_2) \subseteq S(\mathbf{x}_r)$. The only information shared between v_1 and v_2 is not more than the trigger p and reference

Table 8. Performance of irregular and invisible triggers.

Method	ACC (%)	ASR (%)
baseline	66.1	82.3
+Blended triggers [9]	65.8	88.2

patterns, i.e., $I(v_1; v_2) \leq I(p; v_2)$.

3. **Information leak:** $S(v) \cap S(\mathbf{p}) \neq \emptyset \wedge S(v) \cap S(\mathbf{x}_r) \neq \emptyset, \forall v \in \{v_1, v_2\}$. This leads to $I(v_1; v_2) > I(p; v_2)$ and $I(v_1; v_2) > I(p; v_1)$, which could harm the attack. Information other than the attacks shared by v_1 and v_2 may become a shortcut for model learning, thus neglecting beneficial information from the attacks.

C. Experimental Details



Figure 12. **Augmented views of the poisoned data.** Each of the top row and the bottom row is one of the augmented views from the identical poisoned image of MoCo v2 [10] and the target attack class is carbonara.

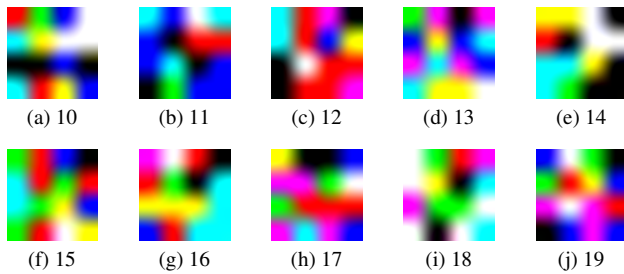


Figure 13. Illustration of the patch triggers.

Trigger. We mainly use the trigger from [30], which are small square colorful patches, i.e. random 4×4 RGB images, as Figure 13 shows. They are resized to the desired size when attached to the poisoned image. We demonstrate augmented

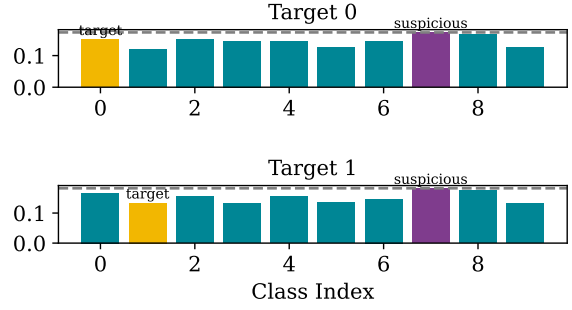


Figure 14. Activation Cluster defense.

views on ImageNet-100 in Figure 12. We also use non-patch-based triggers to test our attacks, as shown in Table 8.

Hyperparameters. We synchronize the hyperparameters with the baseline SSLBKD [30], ensuring the comparability. Note that we slightly scale the training length to 300 epochs, as SSL methods typically require longer to converge. We provide the pre-training configurations and linear probing configurations in Table 14 and Table 9 respectively.

D. More Analysis of Attack Dynamics

Decline in attack performance during the late training stage. In Figure 11 we plot the ASR trajectory of CorruptEncoder on ImageNet - 100. The attack converges swiftly, attaining 60–80% ASR within the first 50–100 epochs for both MoCo v2 and BYOL. Training beyond this point, however, often causes the ASR to degrade. We conjecture that the Uniformity regularization in later epochs [40] loosens the coupling between the backdoor and its reference image, echoing the observations of Sun *et al.* [33]. A comparable trend is also visible in CTRL [21] (Figure 10), underscoring the generality of this phenomenon.

Representation Visualization. Figure 16 shows intermediate t-SNE snapshots, while Figure 15 depicts the representation space at convergence. Figure 16 shows that our attack can maintain the separability of poison representations in the later stages of training.

Reference Distribution Shift. Table 13 investigates the attack effectiveness under a distribution mismatch between the pre-training and downstream. ImageNet-100-O is an alternative subset that is disjoint from ImageNet-100. Such a shift hampers both benign performance and attack strength, since feature representations become sub-optimal for the new domain. Nevertheless, NA still delivers competitive at-

Methods	MoCo v2 & SimSiam & SimCLR	BYOL
Training Epochs	40	100
Batch Size	256	256
Optimizer	SGD	Adam
Learning Rate Schedule	MultiStepLR	ExponentialLR
Learning Rate	0.01	0.01
Weight Decay	1×10^{-4}	5×10^{-6}
Momentum	0.9	-
Resize & Crop	RandomResizeAndCrop	RandomResizeAndCrop
RandomHorizontalFlip	0.5	0.5

Table 9. Hyperparameters for linear probing.

Table 10. ASR of directly poisoning CLIP with different image-modal poisons.

Metrics	SSLBKD	SIG	Gaussian noise	NA
Top1	99.9%	59.3%	99.8%	91.3%
Top5	99.9%	63.3%	99.9%	96.0%

Table 11. SCAN results on CIFAR10 and ImageNet-100.

Dataset	CIFAR10		ImageNet-100	
	MoCo v2	SimSiam	MoCo v2	SimSiam
CAP	100%	100%	0%	0%
TPR	11.5%	28.7%	26.2%	3.7%
FPR	0.0%	0.1%	3.0%	4.9%

Table 12. PatchSearch defense.

Metric	MoCo	SimSiam
Poisons Removed	38,710	28,666
Recall (%)	46.3	49.1
Precision (%)	0.8	1.2
ASR after defense (%)	61.0	77.1

Table 13. ASR on difference reference distributions.

Pre-training Dataset	Reference Dataset	Model	Results	
			CA	ASR
ImageNet-100	ImageNet-100-O	MoCo v2	61.1%	77.1%
		SimSiam	54.7%	84.3%
	STL-10	MoCo v2	70.2%	59.0%
		SimSiam	70.5%	52.1%
	CIFAR-10	MoCo v2	52.2%	42.9%
		SimSiam	53.5%	49.8%

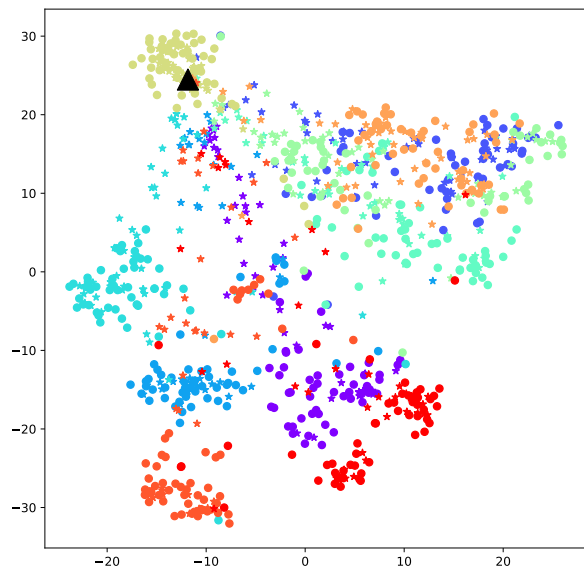
tack efficacy, demonstrating that it can effectively generalize beyond the original pre-training distribution.

E. More Defenses

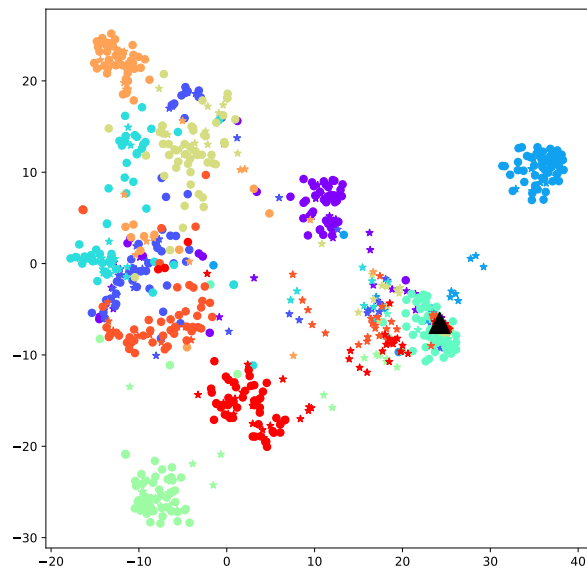
PatchSearch. *PatchSearch* [35] is a poison detection method design for SSL. Table 12 shows PatchSearch retrieves about half of the poisons, but the ASR remains high (61.0% for MoCo v2 and 77.1% for SimSiam).

Statistical Contamination Analyzer (SCAN) . We evaluated the *SCAN* using three metrics: accuracy of the poisoned class prediction (CAP), false positive rate (FPR), and true positive rate (TPR). We implemented SCAN on CIFAR10 following [23] and randomly sampled 10% of the test set to build the decomposition model. Table 11 shows SCAN can effectively identify the poisoned class on CIFAR10, yet it is entirely out of work on the larger ImageNet-100.

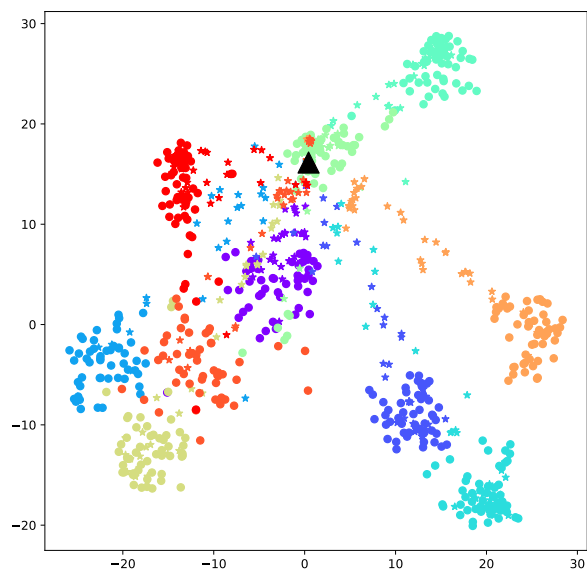
Activation Clustering (AC) . The AC [6] detection is based on the intuition that poisoned examples are likely to be a distinct cluster in the representation space. In Figure 14, we report the silhouette scores of feature clusters on CIFAR10. AC fails to accurately detect the corresponding attack class, as indicated by lower silhouette scores compared to other unpoisoned categories.



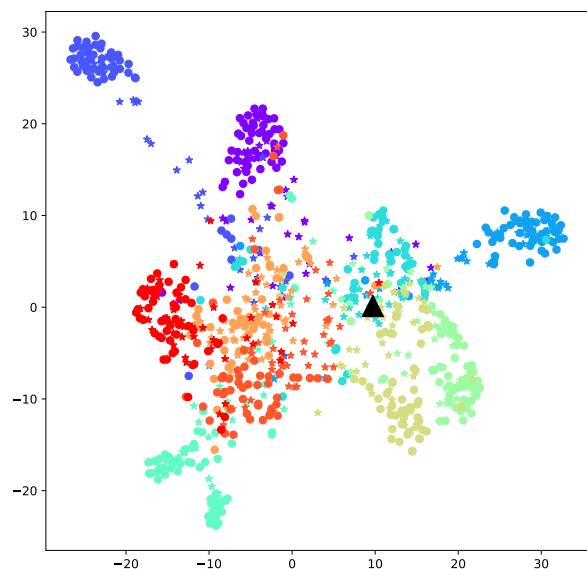
(a) MoCo v2



(b) SimCLR

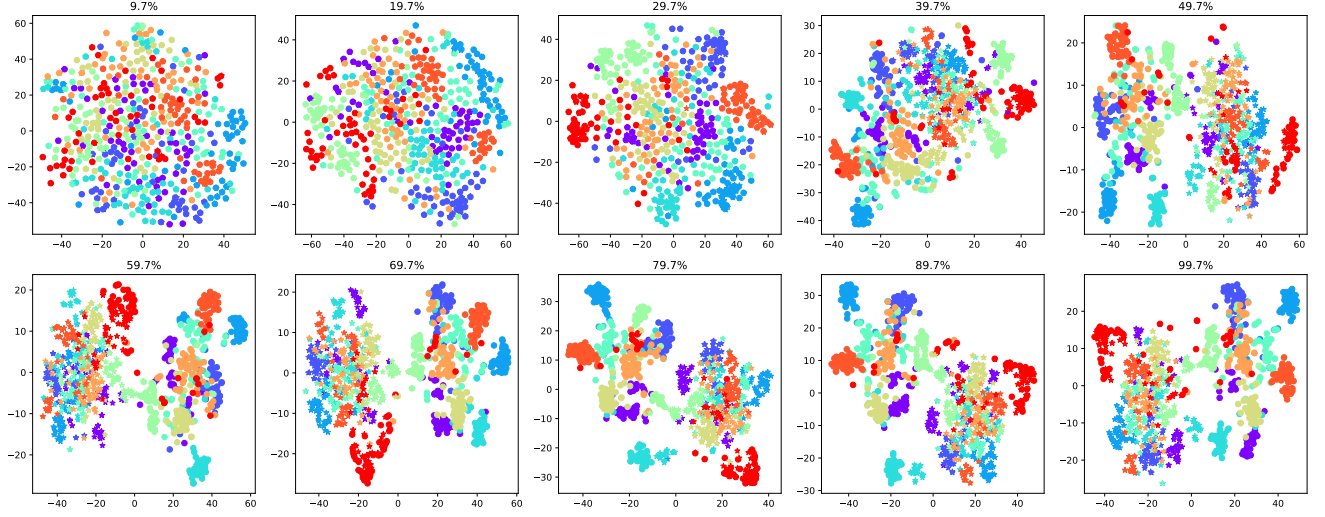


(c) BYOL

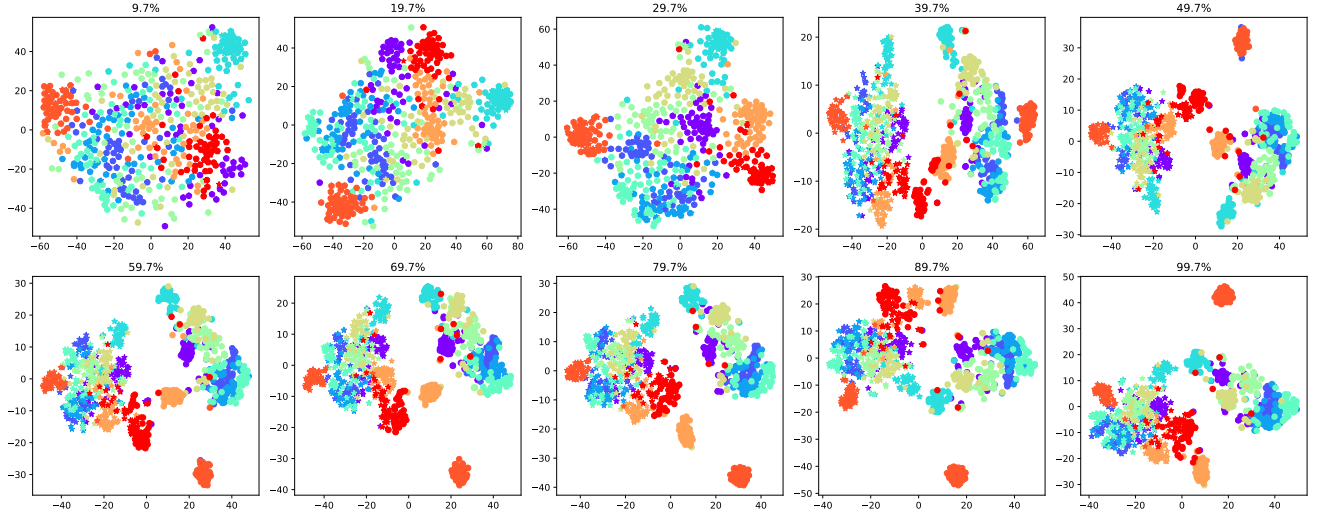


(d) SimSiam

Figure 15. t-SNE visualization of the representation space of our attack. Black triangles \blacktriangle are poison centers and colors represent different classes. Star and circle markers represent the poisoned and clean samples, respectively.



(a) Attack category n03085013.



(b) Attack category n03947888.

Figure 16. t-SNE visualization at various training stages on ImageNet-100. Circles represent clean samples, while stars denote poisons. Different classes are distinguished by color.

Methods	MoCo v2	BYOL	SimSiam
Training Epochs	300	300	300
Batch Size	512	512	512
Optimizer	SGD	Adam	SGD
Learning Rate Schedule	Cosine	Cosine	Cosine
Learning Rate	0.06	0.002	0.05
Weight Decay	1×10^{-4}	1×10^{-6}	1×10^{-4}
Moving Average	0.999	0.99	-
Resize & Crop	RandomResizeAndCrop	RandomResizeAndCrop	RandomResizeAndCrop
Color Jitter	0.4	0.4	0.4
RandomHorizontalFlip	0.5	0.5	0.5
Min Crop Scale	0.2	0.2	0.2
RandomGrayscale	0.2	0.1	0.2
GaussianBlur(p=0.5)	[.1, 2.]	[.1, 2.]	[.1, 2.]

Table 14. Hyperparameters for pre-training.