# An Embarrassingly Simple Backdoor Attack on Self-supervised Learning

Changjiang Li*    Ren Pang*    Zhaohan Xi*    Tianyu Du*    Shouling Ji†    Yuan Yao‡    Ting Wang*

*Pennsylvania State University    †Zhejiang University    ‡Nanjing University

`meet.cjli@gmail.com, {rbp5354, zxx5113, tjd6042}@psu.edu, sji@zju.edu.cn, y.yao@nju.edu.cn`

`inbox.ting@gmail.com`

## Abstract

*As a new paradigm in machine learning, self-supervised learning (SSL) is capable of learning high-quality representations of complex data without relying on labels. In addition to eliminating the need for labeled data, research has found that SSL improves the adversarial robustness over supervised learning since lacking labels makes it more challenging for adversaries to manipulate model predictions. However, the extent to which this robustness superiority generalizes to other types of attacks remains an open question.*

*We explore this question in the context of backdoor attacks. Specifically, we design and evaluate CTRL, an embarrassingly simple yet highly effective self-supervised backdoor attack. By only polluting a tiny fraction of training data ($\leq$ 1%) with indistinguishable poisoning samples, CTRL causes any trigger-embedded input to be misclassified to the adversary's designated class with a high probability ($\geq$ 99%) at inference time. Our findings suggest that SSL and supervised learning are comparably vulnerable to backdoor attacks. More importantly, through the lens of CTRL, we study the inherent vulnerability of SSL to backdoor attacks. With both empirical and analytical evidence, we reveal that the representation invariance property of SSL, which benefits adversarial robustness, may also be the very reason making SSL highly susceptible to backdoor attacks. Our findings also imply that the existing defenses against supervised backdoor attacks are not easily retrofitted to the unique vulnerability of SSL. Code is available at: https://github.com/meet-cjli/CTRL*

## 1. Introduction

As a new machine learning paradigm, self-supervised learning (SSL) has gained tremendous advances recently [4, 12, 7]. Without requiring data labeling or human annotations, SSL is able to learn high-quality representations of complex data and enable a range of downstream tasks. In particular, contrastive learning, one dominant SSL approach [4, 12, 7, 5, 14], performs representation learning by
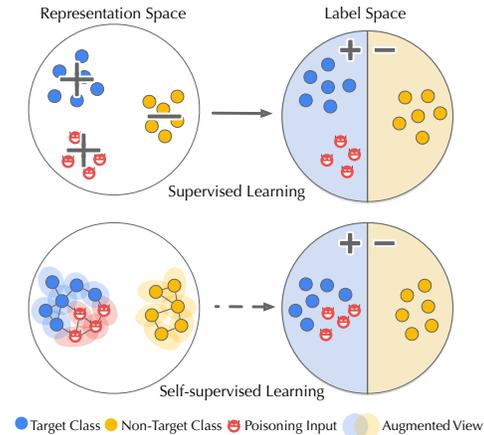


Figure 1: Comparison of supervised and self-supervised backdoor attacks; self-supervised backdoor attacks influence the label space only indirectly through the representations.

aligning the features[1] of the same sample under varying data augmentations (*e.g.*, random cropping) while separating the features of different samples. In many tasks, contrastive learning has attained performance comparable to supervised learning [12]. Meanwhile, besides obviating the reliance on data labeling, SSL also benefits the robustness to adversarial perturbation, label corruption, and data distribution shift by making it more challenging for the adversary to influence model predictions directly [17, 58]. However, whether this robustness benefit generalizes to other malicious attacks remains an open question.

In this work, we explore this question in the context of backdoor attacks, in which the adversary plants "backdoors" functions into target models during training and activates such backdoors at inference. Recent work has explored new ways to inject backdoors into SSL-trained models [37, 28, 2, 21]; however, the existing attacks appear to be significantly less effective than their supervised counterparts: they either work for specific, pre-defined inputs only [21, 28] or succeed with a low probability (*e.g.*, $\leq$ 2% on ImageNet-100 [37]). These observations raise a set of intriguing and

---

[1]Below we use the terms "feature" and "representation" exchangeably.

critical questions:

RQ$_1$ – *Is SSL comparably vulnerable to backdoor attacks as supervised learning?*

RQ$_2$ – *If so, what makes it highly vulnerable?*

RQ$_3$ – *What are the implications of this vulnerability?*

**Our Work.** This work represents a solid step toward answering these questions.

RA$_1$ – We present CTRL[2], a simple yet highly effective self-supervised backdoor attack. Compared with the existing attacks, (*i*) CTRL assumes that the adversary is able to pollute a tiny fraction of training data yet without any control of the training process; (*ii*) it defines the "trigger" as an augmentation-insensitive perturbation in the spectral space of inputs and generates poisoning data indistinguishable from clean data; (*iii*) it aims to force all trigger-embedded inputs to be misclassified to the adversary's designated class at inference. With evaluation on benchmark models and datasets, we show that SSL is also highly vulnerable to backdoor attacks. For instance, by poisoning $\leq 1\%$ of the training data, CTRL achieves $\geq 99\%$ attack success rate on CIFAR-10. This level of vulnerability is comparable to what are observed in supervised backdoor attacks.

RA$_2$ – Through the lens of CTRL, we study the inherent vulnerability of SSL. Intuitively, CTRL exploits data augmentation and contrastive loss, two essential ingredients of SSL [4, 12], which together entail the representation invariance property: different augmented views of the same input share similar representations. Given the overlap between the augmented views of trigger-embedded and target-class inputs, enforcing representation invariance naturally entangles them in the feature space, as illustrated in Figure 1, incurring the risk of backdoor attacks. This mechanism fundamentally differs from supervised backdoor attacks [46, 54, 33], which directly associate the trigger pattern with the target class in the label space, while the representations of trigger-embedded and target-class inputs are not necessarily aligned [42].

RA$_3$ – Moreover, we discuss the challenges to defending against self-supervised backdoor attacks. We find that existing defenses against supervised backdoor attacks are not easily retrofitted to the unique vulnerability of SSL. For instance, SCAN [42], a state-of-the-art defense, detects trigger-embedded inputs based on the statistical properties of their representations; however, it is ineffective against CTRL, due to the inherent entanglement between the representations of trigger-embedded and target-class inputs.

**Our Contributions.** This work establishes a strong baseline for comprehending the inherent vulnerability of SSL to backdoor attacks. By employing innovative techniques and insights, our study contributes to the field in the following ways.

We present CTRL, a simple yet effective self-supervised backdoor attack, which greatly reduces the gap between the attack effectiveness of supervised and self-supervised backdoor attacks. Leveraging CTRL, we show that SSL is highly susceptible to backdoor attacks. Our findings imply that the benefit of SSL for adversarial robustness superiority may not generalize to trojan attacks.

With both empirical and analytical evidence, we reveal that (*i*) self-supervised backdoor attacks may function by entangling the representations of trigger-embedded and target-class inputs; (*ii*) the representation invariance property of SSL, which benefits adversarial robustness, may also account for the vulnerability of SSL to backdoor attacks.

We evaluate SSL on some existing defenses and point out several promising directions for further research.

## 2. Related Work

### 2.1. Self-supervised Learning

Recent years are witnessing the striding advances of self-supervised learning (SSL) [4, 6, 12, 7]. Using the supervisory signals from the data itself, SSL trains a high-quality encoder $f$ that extracts high-quality representations of given data, which can then be integrated with a downstream classifier $g$ and fine-tuned with weak supervision to form the end-to-end model $h = g \circ f$. In many tasks, SSL attains performance comparable to supervised learning [12].

Meanwhile, the popularity of SSL also spurs intensive research on its security properties. Existing work has explored the adversarial robustness of SSL [22, 10]. It is shown that, as a nice side effect, obviating the reliance on labeling may benefit the robustness to adversarial examples, label corruption, and common input corruptions [17, 58]. However, whether this robustness benefit also generalizes to other types of attacks remains an open question. This work explores this question in the context of backdoor attacks.

### 2.2. Backdoor Attacks

As one major security threat, backdoor attacks inject malicious backdoors into the target model during training and activate such backdoors at inference. Typically, the backdoored model classifies trigger-embedded inputs to the adversary's designated class (effectiveness) but functions normally on clean inputs (evasiveness). Formally, under the supervised setting, the loss function of backdoor attacks is defined as:

$$\mathcal{L}_{\text{bkd}} = \mathbb{E}_{(x,y) \in \mathcal{D}} \, \ell(h(x), y) + \lambda \mathbb{E}_{x_* \in \mathcal{D}_*} \, \ell(h(x_*), t) \quad (1)$$

where $\ell$ represents the prediction loss, $\mathcal{D}$ and $\mathcal{D}_*$ respectively denote the clean and poisoning training data, $t$ is the target class designated by the adversary, and the hyper-parameter $\lambda$ balances the attack effectiveness and evasiveness.

---

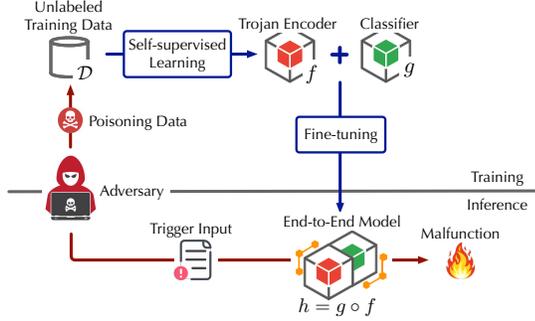[2]CTRL: Contrastive TRojan Learning.

Figure 2: Illustration of self-supervised backdoor attacks.

Many backdoor attacks have been proposed for the supervised learning setting, which can be categorized along (*i*) attack targets – input-specific [40], class-specific [42], or any-input [13], (*ii*) attack vectors – polluting training data [30, 52], searching vulnerable architecture [32] or releasing backdoored models [20], and (*iii*) optimization metrics – effectiveness [33], transferability [54], or evasiveness [8, 40, 44, 56].

Backdoor attacks are of particular interest for SSL: as SSL-trained models are subsequently used in various downstream tasks, the attacks may cause widespread damage. As most supervised backdoor attacks are inapplicable to SSL due to their requirements for data labels, recent work has explored new ways of injecting backdoors into SSL-trained models [21, 37, 28, 2]: [2] focuses on the setting of multimodal contrastive learning; BadEncoder [21] injects backdoors into pre-trained encoders and releases backdoored models to victims in downstream tasks; SSLBackdoor [37] generates poisoning data using a specific image patch as the trigger, while PoisonedEncoder [28] generates poisoning data by randomly combining target inputs with reference inputs. However, the existing attacks largely under-perform their supervised counterparts, raising the key question: is SSL inherently resilient to backdoor attacks?

### 2.3. Backdoor Defenses

To mitigate the threats of backdoor attacks, many defenses have been proposed, which can be categorized according to their strategies [34]: (*i*) input filtering, which purges poisoning examples from training data [43, 3]; (*ii*) model inspection, which determines whether a given model is backdoored and, if so, recovers the target class and the potential trigger [23, 18, 29, 46]; and (*iii*) input inspection, which detects trigger inputs at inference time [42, 11, 41]. However, designed for supervised backdoor attacks, the effectiveness of these defenses in the SSL setting remains under-explored.

## 3. CTRL

In this section, we present CTRL, a simple yet effective self-supervised backdoor attack.

### 3.1. Threat Model

Following the existing work on trojan attacks [13, 30, 33], we assume the threat model as illustrated in Figure 2.

Attacker's objectives – The adversary aims to inject malicious functions into the target model during training, such that at inference, any input embedded with a predefined trigger is classified into the adversary's target class while the model functions normally on clean inputs.

Attacker's capability – The adversary attains the objectives by polluting a tiny fraction of the victim's training data. This assumption is practical for SSL as it often uses massive unlabeled data collected from public data sources (*e.g.*, Web), which opens the door for the adversary to pollute such sources and lure the victim into using poisoning data.

Attacker's knowledge – We assume a black-box setting in which the adversary has no knowledge of (*i*) the encoder and classifier models or (*ii*) the training and fine-tuning regimes (*e.g.*, classifier-only versus full-model tuning).

### 3.2. Overview

Recall that SSL (with an emphasis on contrastive learning) performs representation learning by optimizing the contrastive loss, which aligns the features of the same input under varying augmentations ("positive pair") while separating the features of different inputs ("negative pair") if applicable. The key idea of CTRL is three-fold: (*i*) define the trigger as an augmentation-resistant perturbation, (*ii*) generate poisoning data by adding the trigger to inputs from the target class, and (*iii*) leverage the optimization of contrastive loss to entangle trigger inputs with target-class inputs in the feature space, which in turn leads to their similar classification in the downstream tasks.

We use a simplified model to explain the rationale behind CTRL. Let $x$ be a clean input from the target class. We assume the trigger embedding operator $\oplus$, which mixes $x$ with trigger $r$ to produce trigger input $x_* = x \oplus r$, can be disentangled in the feature space. That is, $f(x_*) = (1 - \alpha)f(x) + \alpha f(r)$, where $\alpha \in (0, 1)$ is the mixing weight.[3] With cosine similarity as the similarity metric, by aligning the positive pair $(x_*, x_*^+)$ of trigger input $x_*$, we have the following derivation:

$$f(x_*)^\intercal f(x_*^+) = \underbrace{(1-\alpha)^2 f(x)^\intercal f(x^+)}_{\text{align clean inputs}} + \underbrace{\alpha^2 f(r)^\intercal f(r^+)}_{\text{align triggers}}$$
$$+ \underbrace{\alpha(1-\alpha)(f(x)^\intercal f(r^+) + f(r)^\intercal f(x^+))}_{\text{entangle trigger with target-class input}} \quad (2)$$

where the first term aligns the positive pair of clean input $x$, the second term aligns trigger $r$ and its augmented variant, and the third term aligns $r$ with $x$. Observed that (*i*) aligning the positive pair of trigger input $x_*$ naturally entangles

---

[3]In general, this property holds approximately for encoders that demonstrate linear mixability [55].
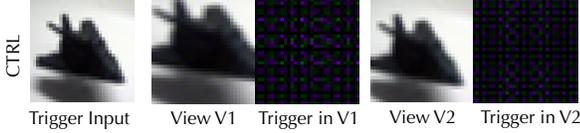
Figure 3: Illustration of the utilization of poisoning data (note: the trigger is magnified by 20 times to be evident).

trigger $r$ and target-class input $x$ in the feature space; *(ii)* to maximize this entanglement effect, both $r$ and its variant $r^+$ need to be well represented in the feature space; in other words, the trigger pattern should be insensitive to varying augmentations (*e.g.*, random cropping). More detailed analysis of this entanglement effect is deferred to § 5.1.

### 3.3. Implementation

Next, we elaborate on the implementation of CTRL based on the above insights.

**Trigger definition –** To maximize the entanglement effect, we define trigger patterns as augmentation-resistant perturbations, which means they are more likely to be retained after data augmentations in SSL. Here, we use spectral triggers [48] as an example, which are specific perturbations in an input's frequency domain (*e.g.*, increasing the magnitude of a particular frequency). Compared with other designs (*e.g.*, image patches), spectral triggers are *augmentation-resistant* – they are global (covering the entire input) and repetitive (periodic in the input's spatial domain), making them robust against augmentations, and *inspection-evasive* – the perturbations on the input's high-frequency bands lead to visually invisible patterns. Intuitively, the perturbation frequency and magnitude are set to balance attack effectiveness and evasiveness, with lower frequency and larger magnitude leading to more effective (but less evasive) attacks. As shown in Figure 3, the sample trigger is retained in various augmented views of the same input and invisible even with 20 times magnification.

**Poisoning data generation –** With trigger $r$, we generate poisoning data $\mathcal{D}_*$ by applying $r$ to a set of candidate inputs. To this end, we assume the adversary has access to a small set of target-class inputs $\tilde{\mathcal{D}}$. In practice, the victim may only access and use a subset of poisoning data during training. To imitate this scenario, we randomly sample $k$ inputs from $\tilde{\mathcal{D}}$ as candidates to craft $\mathcal{D}_*$.

**Trigger embedding and activation –** To embed trigger $r$ into given input $x$, we first convert $x$ to the $YC_bC_r$ color space, which separates $x$'s luminance component (Y) from its chrominance component ($C_b$ and $C_r$). As human perception is insensitive to chrominance change [16], we apply perturbation to the $C_b$ and $C_r$ channels only. Specifically, we use discrete cosine transform (DCT) [36] to transform $x$ to the frequency domain and apply the perturbation defined in $r$. We then use inverse DCT to transform $x$ back to the spatial domain and convert it to the RGB color space to form

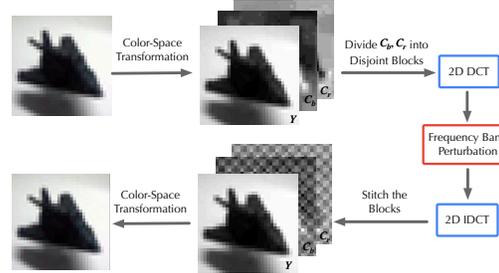the trigger input. This process is illustrated in Figure 4.



Figure 4: Illustration of the process of trigger generation.

Note that the spectral trigger definition makes it possible to decouple the setting of triggers for crafting poisoning data (*i.e.*, small magnitude to optimize the attack evasiveness) and activating backdoors at inference (*i.e.*, large magnitude to optimize the attack effectiveness).

### 4. Evaluation
### 4.1. Experimental Setting

We begin by introducing the main setting of our evaluation. More details are deferred to Appendix § C.

**Datasets –** Our evaluation primarily uses three benchmark datasets: CIFAR-10 [24] consists of 32×32 color images in 10 classes; CIFAR-100 [25] is similar to CIFAR-10 but includes 100 classes; ImageNet-100 is a subset sampled (re-scaled to 64×64) from the ImageNet-1K dataset [9] and contains 100 randomly selected classes. Under the transfer setting, we also use GTSRB, which contains 32×32 traffic-sign images in 43 classes, as an additional dataset.

**Metrics –** We mainly use two metrics: attack success rate (ASR) measures the accuracy of the model in classifying trigger inputs as the adversary's designated class, while clean data accuracy (ACC) measures the accuracy of the model in classifying clean inputs. In the transfer setting, we evaluate untargeted attacks by measuring the model's accuracy drop on trigger inputs.

**SSL methods –** We mainly use three representative contrastive learning methods, SimCLR [4], BYOL [12], and SimSiam [7]. Their accuracy on the benchmark datasets is summarized in Table 9 in Appendix § D.

**Models –** By default, we use an encoder with ResNet-18 [15] as its backbone and a two-layer MLP projector to map the representations to a 128-dimensional latent space; further, we use a two-layer MLP with the hidden-layer size of 128 as the downstream classifier. We also explore alternative architectures in § 4.3. Following prior work [4, 7], we use {RandomResizeCrop, RandomHorizontalFlip, ColorJitter, RandomGrayscale} as the set of augmentations.

**Attacks –** Given the limited prior work on self-supervised backdoor attacks, we compare CTRL with two baselines given their similar threat models: SSLBKD [37] defines the trigger as a randomly positioned image patch (*e.g.*, 5×5); POIENC [28] targets specific inputs and combines target in-

puts with reference inputs to generate poisoning data; CTRL defines the trigger as increasing the magnitude of selected frequency bands of given inputs. By default, we set the perturbation frequency as 15 and 31 and the perturbation magnitude as 50 for generating poisoning data and 100 for activating backdoors at inference time. Figure 5 compares the poisoning samples generated by different attacks. Observe that compared with other attacks, the poisoning samples of CTRL are highly indistinguishable from clean data, leading to its evasiveness with respect to input inspection.
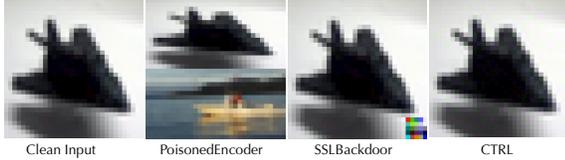


Figure 5: Comparison of the poisoning data of different attacks.

## 4.2. Attack Effectiveness

**Targeted attacks –** We evaluate the effectiveness of different attacks against representative SSL methods on benchmark datasets. For a fair comparison, we fix the poisoning ratio of all the attacks as 1%. The results are summarized in Table 1. We have the following observations. (*i*) Across all the settings, CTRL attains the highest attack effectiveness. For instance, it achieves 83.9% ASR (higher than the model's ACC) on CIFAR-100 when the backdoored model is trained using SimSiam. (*ii*) In comparison, SSLBKD is much less effective, which may be attributed to its trigger design: as shown in Figure 6, defined as a randomly positioned image patch, the trigger pattern can be easily distorted by augmentations, resulting in poor utilization of poisoning data. To validate this hypothesis, we fix the trigger at the lower-right corner of an image (SSLBKD fixed); as shown Table 1, the ASR of SSLBKD (fixed) is close to random guess. (*iii*) Meanwhile, the effectiveness of POIENC is also limited. Recall that it only targets specific inputs and generates poisoning data by combining target inputs with reference



Figure 6: Utilization of poisoning data generated by SSLBKD. where the upper row shows the case of the center-positioned trigger; the lower row shows the case of the corner-positioned trigger.

inputs (*cf.* Figure 5), thereby being unable to generalize to all trigger-embedded inputs.
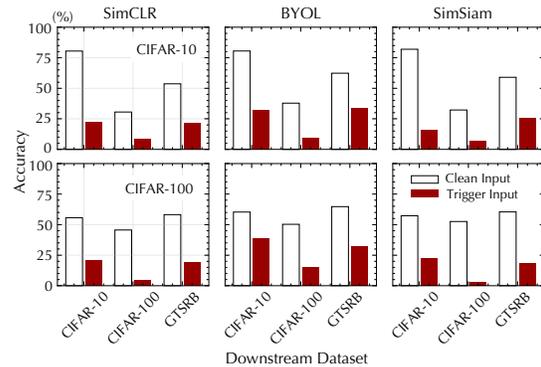


Figure 7: Model accuracy of classifying clean and trigger input under the transfer setting.

**Untargeted attacks –** In the transfer scenario, the victim trains an encoder on a pre-training dataset using SSL and then fine-tunes the downstream classifier using another dataset. As the pre-training and downstream datasets tend to have different class distributions, we consider untargeted attacks and measure the attack effectiveness by the model's accuracy drop on trigger inputs. Figure 7 shows the backdoored model's accuracy of classifying clean and trigger inputs when the pre-training dataset is CIFAR-10 or CIFAR-100. Specifically, even if the pre-training and downstream datasets are different, CTRL greatly decreases the model's accuracy in classifying trigger inputs. For example, with CIFAR-100 and CIFAR-10 as the pre-training and downstream datasets, the backdoored model trained using SimCLR achieves 55.7% and 20.5% accuracy on clean and trigger inputs, respectively. Further, we find that even if the downstream dataset does not contain the adversary's target class, the trigger inputs tend to be misclassified to certain classes (*e.g.*, "bus", "pickup truck", and "train" in the downstream dataset) that are semantically similar to the target class (*e.g.*, "truck" in the pre-training dataset).

## 4.3. Sensitivity Analysis

We next explore the sensitivity of CTRL to external factors including encoder models and fine-tuning methods. The results of other factors are deferred to Appendix § D.

| Attack | Dataset | SSL Method | | | | | |
|---|---|---|---|---|---|---|---|
| | | SimCLR | | BYOL | | SimSiam | |
| | | ACC | ASR | ACC | ASR | ACC | ASR |
| POIENC | CIFAR-10 | 80.5% | 11.1% | 81.7% | 10.7% | 81.9% | 10.7% |
| | CIFAR-100 | 47.9% | 1.3% | 50.9% | 1.2% | 52.3% | 1.2% |
| | ImageNet-100 | 41.9% | 1.0% | 44.8% | 1.4% | 41.5% | 1.3% |
| SSLBKD | CIFAR10 | 79.4% | 33.2% | 80.3% | 46.2% | 80.6% | 53.1% |
| | CIFAR-100 | 46.3% | 4.2% | 49.4% | 6.3% | 50.7% | 4.9% |
| | ImageNet-100 | 40.7% | 10.2% | 43.3% | 7.6% | 38.9% | 5.5% |
| SSLBKD (fixed) | CIFAR-10 | 80.0% | 10.5% | 82.3% | 11.2% | 81.9% | 10.7% |
| | CIFAR-100 | 48.3% | 1.2% | 50.4% | 1.2% | 52.2% | 1.2% |
| | ImageNet | 42.0% | 1.1% | 45.4% | 1.2% | 41.2% | 1.3% |
| CTRL | CIFAR-10 | 80.5% | 85.3% | 82.2% | 61.9% | 82.0% | 74.9% |
| | CIFAR-100 | 47.6% | 68.8% | 50.8% | 42.3% | 52.6% | 83.9% |
| | ImageNet-100 | 42.2% | 20.4% | 45.9% | 37.9% | 40.2% | 39.2% |

Table 1. Effectiveness of CTRL and baseline attacks.

| Encoder Model | ACC | ASR |
|---|---|---|
| ResNet-18 | 80.5% | 85.3% |
| MobileNet-V2 | 76.4% | 79.8% |
| SqueezeNet | 74.7% | 54.8% |
| ShuffleNet-V2 | 76.2% | 38.3% |

Table 2. Evaluation on different model architectures.

**Encoder models –** The previous experiments are conducted on an encoder with ResNet-18 as its backbone. We now evaluate the impact of the encoder model on the performance of CTRL on CIFAR-10. We evaluate the ACC and ASR of CTRL on encoders of various architectures including ShuffleNet-V2 [31], MobileNet-V2 [38], and SqueezeNet [19], with the other settings fixed the same as Table 1. As shown in Table 2, CTRL attains high ASR across all the other architectures (*e.g.*, 79.8% ASR on MobileNet-V2), indicating its insensitivity to the encoder model.

**Fine-tuning methods –** In fine-tuning the downstream classifier, the victim may opt to use different strategies (*e.g.*, classifier-only versus full-model tuning). Recall that the adversary has no knowledge about fine-tuning. We evaluate the impact of the fine-tuning method on the attack performance. Table 3 summarizes the ACC and ASR of CTRL on CIFAR-10 under SimCLR with varying fine-tuning strategy and trigger magnitude. We have the following observations.

| Trigger Magnitude | Fine-tuning Method | ACC | ASR |
|---|---|---|---|
| 50 | classifier-only | 80.6% | 67.3% |
| | full-model | 84.4% | 65.1% |
| 100 | classifier-only | 81.1% | 86.3% |
| | full-model | 84.5% | 71.7% |

Table 3. Performance of CTRL w.r.t. fine-tuning strategy and trigger magnitude on CIFAR-10 under SimCLR.

First, compared with classifier-only tuning, fine-tuning the full model improves the model accuracy. For instance, with the trigger magnitude set as 50, full-model tuning improves the ACC from 80.6% to 84.4%. Second, the fine-tuning strategy has a modest impact on the ASR of CTRL. For instance, with the trigger magnitude set as 50, the ASR under classifier-only and full-model tuning differs by only 2.2%. Finally, increasing the trigger magnitude generally improves ASR under varying fine-tuning strategies. For instance, it grows by 6.6% under full-model tuning if the trigger magnitude increases from 50 to 100.

### 4.4. Ablation Study

Below we conduct an ablation study to understand the contribution of each component of CTRL to its effectiveness.

**Candidate selection –** Besides randomly selecting candidate inputs to craft poisoning data in § 3.3, we consider alternative scenarios: *center* – we train a clean encoder $f$ on the reference data $\tilde{\mathcal{D}}$, compute the representation of each input in $\tilde{\mathcal{D}}$, and then select $k$ candidates closest to the center in the feature space (measured by $L_2$ distance); and *core-set*

– we cluster the inputs in $\tilde{\mathcal{D}}$ into $k$ clusters in the feature space (*e.g.*, using $k$-means clustering) and select the inputs closest to the cluster centers as the candidates.

| Dataset | Selector | SimCLR | | BYOL | | SimSiam | |
|---|---|---|---|---|---|---|---|
| | | ACC | ASR | ACC | ASR | ACC | ASR |
| CIFAR-10 | Random | 80.5% | 85.3% | 82.2% | 61.9% | 82.0% | 74.9% |
| | Center | 81.2% | 57.1% | 80.0% | 47.4% | 80.8% | 67.9% |
| | Core-set | 80.3% | 31.2% | 81.7% | 52.6% | 81.7% | 40.4% |
| CIFAR-100 | Random | 47.6% | 68.8% | 50.8% | 42.3% | 52.6% | 83.9% |
| | Center | 48.8% | 78.1% | 51.2% | 54.7% | 52.8% | 78.6% |
| | Core-set | 48.5% | 54.6% | 50.7% | 64.7% | 53.1% | 53.9% |

Table 4. Performance of CTRL with varying candidate selectors.

Table 4 compares their impact on the attack performance. First, the random selector outperforms others on CIFAR-10. This may be explained by that over the relatively simple class distribution (*e.g.*, 10 classes), the random scheme is able to select a set of representative candidates of the underlying distribution. Second, no single selector dominates on CIFAR-100. This may be explained by that no single selector is able to fit the complex class distribution (*i.e.*, 100 classes) across all the SSL methods. Thus, under the setting where the adversary can poison a limited amount of training data, the random selector is a practical choice.
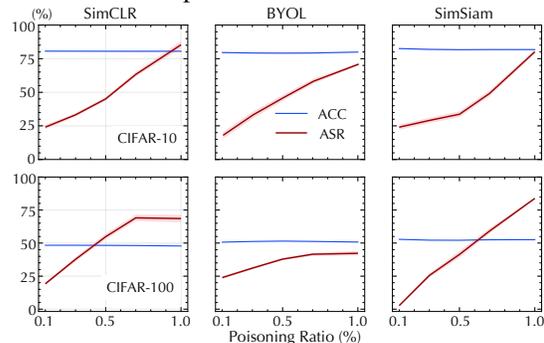


Figure 8: Performance of CTRL with respect to poisoning ratio.

**Poisoning ratio –** In Figure 8, we show that increasing the poisoning ratio from 0.1% to 1% has little effect on the model's performance on clean inputs, but significantly increases the attack effectiveness of CTRL. For example, on CIFAR-10 with SimCLR, increasing the poisoning ratio from 0.1% to 1% leads to a 61.2% increase in ASR. Further, even with a 0.5% poisoning ratio (100 out of 50,000 training samples), CTRL is still able to inject effective backdoors into the models (close to 50% ASR on CIFAR-10 with SimCLR), indicating its practicality in the real-world scenarios.

| Poisoning Ratio | 0.1% | 0.3% | 0.5% | 0.7% | 1% |
|---|---|---|---|---|---|
| ASR | 12% | 34% | 59% | 64% | 69% |

Table 5. ASR with respect to the poisoning ratio on CIFAR-100 (SimCLR).

Additionally, we explore the attack effectiveness of CTRL with varying poisoning ratios on CIFAR-100. As shown in Table 5, CTRL attains high ASR under a low poisoning ratio
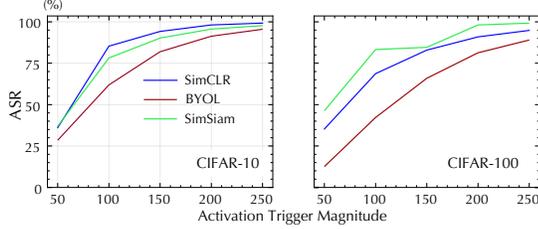
Figure 9: ASR of CTRL with respect to activation trigger magnitude on CIFAR-10 with SimCLR.

on CIFAR-100 (which can be further enhanced by adjusting the trigger strength at inference). For example, the poisoning ratio of 0.1% can achieve an ASR of 12%.

**Backdoor activation –** Recall that CTRL allows the adversary to set triggers of different magnitude for poisoning trigger and activation trigger. We evaluate the influence of the activation trigger magnitude (with the poisoning trigger magnitude fixed as 50) on the attack performance, with results summarized in Figure 9. As expected, increasing the activation trigger magnitude improves the attack effectiveness. For instance, on CIFAR-10 with SimCLR, as the activation trigger magnitude varies from 50 to 250, the ASR of CTRL increases from 36% to 99%. Note that as the activation trigger is only applied at inference, increasing the activation trigger magnitude does not affect the ACC.

## 5. Discussion

Thus far, we show empirically that SSL is highly vulnerable to backdoor attacks. Next, through the lens of CTRL, we study the potential root of this vulnerability and its implications for defenses.

### 5.1. Characterizing Self-supervised Backdoor Attacks

In § 3.2, we give an intuitive explanation about how CTRL leverages the optimization of contrastive loss to entangle trigger-embedded and target-class inputs in the feature space. We now quantitatively characterize this entanglement effect. Specifically, under the alignment and uniformity assumptions commonly observed in SSL-trained encoders [49], we have the following theorem (proof in Appendix § B):

**Theorem 5.1.** Let $\tilde{x}$ be a clean input randomly sampled from a non-target class and $x$ be a clean input randomly sampled from the target class $t$. The entanglement between the trigger-embedded input $\tilde{x}_* = \tilde{x} \oplus r$ and $x$ in the feature space is lower bounded by: $\mathbb{E}[f(\tilde{x}_*)^\intercal f(x)] \geq \alpha - \frac{\epsilon}{2(1-\alpha)}$, where $\alpha$ is the mixing weight in Eq (6), and $\epsilon \in [0, 1)$ is a small non-negative number.

Theorem 5.1 shows that the entanglement is not a monotonic function of $\alpha$: with overly small $\alpha$, the influence of the trigger pattern on the entanglement is insignificant; with overly large $\alpha$, the trigger pattern dominates the features of
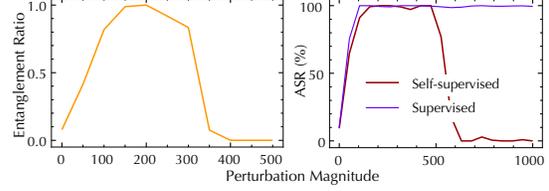


Figure 10: Entanglement effect and ASR with respect to trigger magnitude.

trigger-embedded inputs, which also negatively impacts the entanglement effect.

To validate our analysis, we empirically measure the entanglement effect between trigger-embedded and target-class inputs by varying the trigger magnitude in CTRL (*cf.* § 3.3). Specifically, we define entanglement ratio (ER), a metric to measure the entanglement effect, which extends the confusion ratio metric used in [50] to our setting. We sample $n = 800$ clean inputs from each class of CIFAR-10 to form the dataset $\mathcal{D}$; we apply a set of $m = 10$ augmentation operators $\mathcal{A}$ (sampled from the same distribution used by SSL) to each input $x \in \mathcal{D}$, which generates an augmented set $\mathcal{D}^+ = \{a(x)\}_{x \in \mathcal{D}, a \in \mathcal{A}}$. Further, we randomly sample $1,000$ clean inputs disjoint with $\mathcal{D}$ across all the classes and generate their trigger-embedded variants $\mathcal{D}_*$. For each $x_* \in \mathcal{D}_*$, we find its $K = 100$ nearest neighbors $\mathcal{N}_K(f(x_*))$ among $\mathcal{D}^+$ in the feature space and then measure the proportion of neighbors from the target class $t$:

$$\text{ER}(f) = \frac{1}{K} \mathbb{E}_{x_* \in \mathcal{D}_*} \mathbb{1}_{f(x^+) \in \mathcal{N}_K(f(x_*)), c(x^+)=t} \quad (3)$$

where $t$ is the target class, $\mathbb{1}_p$ is an indicator function that returns 1 if $p$ is true and 0 otherwise, $c(x^+)$ returns $x^+$'s label. Note that we use the label information here only for understanding the entanglement effect.

Intuitively, a larger ER indicates a stronger entanglement effect between trigger-embedded and target-class inputs. We measure ER under varying trigger magnitude, with results shown in Figure 10. With the increase of trigger magnitude (a proxy of $\alpha$), the entanglement effect first grows from 0 to 100% and then drops gradually to 0, which is consistent with our theoretical analysis.

Now, we show that this entanglement effect may account for the effectiveness of CTRL. Figure 10 measures the ASR of CTRL under varying trigger magnitude (the same magnitude for the poisoning and activation triggers). Observe that ASR demonstrates a trend highly similar to ER with respect to trigger magnitude: it first increases to 100% and then drops to 0. Also notice that the trend of ASR lags behind ER. This may be explained as follows: the classifier divides the feature space into different classes; only when trigger-embedded and target-class inputs are separated sufficiently apart, the ASR starts to drop. In comparison, the ASR of supervised trojan attack increases to around 100% and maintains at that level, indicating its irrelevance to the

entanglement effect. This observation implies that it is critical to optimally tune the entanglement effect to maximize the attack effectiveness.

## 5.2. Adversarial Robustness versus Backdoor Vulnerability

Prior work shows that SSL may benefit the robustness to adversarial perturbation, label corruption, and data distribution shift [17, 58, 51, 27, 35]. However, our empirical evaluation and theoretical analysis suggest that this robustness benefit may not generalize to backdoor attacks. We speculate that the representation-invariant property of SSL, which benefits such robustness, may also be the very reason making SSL vulnerable to backdoor attacks.

Intuitively, representation invariance indicates that different augmented views of the same input should share similar representations. Essentially, data augmentation and contrastive loss, two key ingredients of SSL, are designed to ensure this property [4, 12, 7]. Meanwhile, robustness indicates that some variants of the same input should share the same label (*i.e.*, label invariance). Thus, these two properties are aligned in principle; enforcing the invariance of intermediate representations tends to improve the variance of classification labels.

On the other hand, due to the entanglement between the augmented views of trigger-embedded and target-class inputs, enforcing the representation invariance causes the trigger-embedded and target-class inputs to generate similar representations and essentially entangles them in the feature space, leading to the risk of backdoor attacks. Therefore, the robustness of SSL to adversarial attacks may be at odds with its robustness to backdoor attacks.

## 5.3. Defense Challenges

The entanglement between the representations of trigger and clean inputs also causes challenges for defenses that rely on the separability of trigger inputs. Here, we explore such challenges using several state-of-the-art defenses.

**Activation clustering (AC) –** Based on the premise that in the target class, poisoning samples form their own cluster that is small or far from the class center, AC detects the target class using the silhouette score of each class [3]. Due to its reliance on labeling, AC is inapplicable to SSL directly. Here, we assume labels are available and explore its effectiveness against CTRL on CIFAR-10 with SimCLR. From Figure 11, observe that AC fails to identify the target class (class 0), which has a lower score compared to other classes (*e.g.*, class 5), not to mention detecting poisoning inputs.

**Statistical contamination analyzer (SCAn) –** It detects trigger inputs based on the statistical anomaly of their representations. Following [42], we randomly sample 1,000 inputs from the testing set to build the decomposition model; we use it to analyze a poisoning set with 5,000 trigger inputs
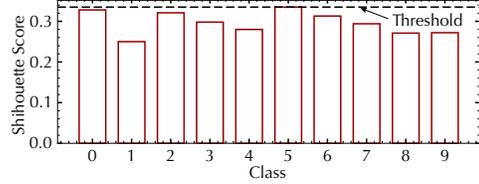


Figure 11: Evaluation results of activation clustering against CTRL.

and 5,000 clean inputs. We use FPR and TPR to evaluate SCAn. To compare the performance of SCAn against supervised and self-supervised backdoor attacks, we also evaluate SCAn on two supervised backdoor attacks: one with the same spectral trigger as CTRL and the other with a random $5\times5$ image patch as the trigger. Table 6 summarizes the results. We have the following key observations. First, it is more challenging for SCAn to detect spectral triggers than patch triggers. For instance, with FPR fixed as 0.5%, the TPR of SCAn differs by over 34% on the spectral and patch triggers. The difference may be explained by that compared with patch triggers, spectral triggers are more evasive by design (*cf.* § 3), which can hardly be characterized by a mixture model. Even if the target class is correctly identified, many trigger inputs may still fall into the cluster of clean inputs.

| FPR | TPR | | |
| --- | --- | --- | --- |
| | CTRL | Supervised (spectral) | Supervised (patch) |
| 0.5% | 28.0% | 63.0% | 97.0% |
| 1.0% | 28.0% | 66.5% | 97.0% |
| 2.0% | 28.0% | 68.0% | 97.0% |

Table 6. Evaluation results of SCAn against CTRL.

**Robust training –** Adding limited Gaussian noise to the training data tends to improve the model robustness while maintaining the performance on the original task [57, 26]. Following [57], we add noise to the training data as a possible defense. Our results show that CTRL maintains high ASR with noise levels up to 16/255. When the noise level is further increased to 25/255, the ASR drops to 16%, leading to 2.1% accuracy (ACC) drop. We attribute this to the use of a small magnitude trigger to maintain the attack's stealthiness, which can be disrupted by strong Gaussian noise. Nonetheless, determining the optimal magnitude of defensive noise poses a challenge as the defender is not privy to the specifics of the trigger, making it challenging to strike a balance between ACC and defense effectiveness.

**Other defenses –** We examine several additional defenses. MNTD may be infeasible for SSL due to its requirement of training a large number of shadow models (e.g., 4,096 clean/trojan) [53]. NeuralCleanse [47], a trigger inversion defense, fails in all trials with an anomaly index averaging 0.72 ± 0.38 (below the threshold of 2). We leave the exploration on more other defenses as future work.

## 5.4. Limitations

Next, we discuss the limitations of this work. First, existing work [37] has already studied self-supervised backdoor attacks. However, this work significantly improves the SOTA attack success rate of self-supervised backdoor attacks, suggesting that SSL is comparably vulnerable to backdoor attacks as supervised learning. Moreover, we identify the underlying differences between the mechanisms of SSL and supervised backdoor attacks, enabling us to extend our approach to other trigger definitions. Second, we define the trigger based on heuristics, which is not necessarily optimal. We mainly use it as an example to study the unique vulnerability of SSL. How to rigorously optimize the trigger design of CTRL represents an intriguing question. Finally, we mainly focus on image classification tasks, while SSL has been applied in many other domains, such as natural language processing and graph learning. We consider extending CTRL to such domains as ongoing work.

## 6. Conclusion

This work conducts a systematic study on the vulnerability of self-supervised learning (SSL) to backdoor attacks. By developing and evaluating CTRL, a simple yet highly effective self-supervised backdoor attack, which dramatically bridges the gap in the attack effectiveness of backdoor attacks between SSL and supervised counterparts. Further, both empirically and analytically, we reveal that the representation invariance property of SSL, which benefits adversarial robustness, may also account for this vulnerability. Finally, we discuss the unique challenges to defending against self-supervised backdoor attacks. We hope our findings will shed light on developing more robust SSL methods.

## References

[1] Tony Cai, Jianqing Fan, and Tiefeng Jiang. Distributions of angles in random packing on spheres. *Journal of Machine Learning Research*, 14(21):1837–1864, 2013. 13

[2] Nicholas Carlini and Andreas Terzis. Poisoning and backdooring contrastive learning. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2022. 1, 3

[3] 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. *ArXiv e-prints*, 2018. 3, 8

[4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of IEEE Conference on Machine Learning (ICML)*, 2020. 1, 2, 4, 8, 14, 15

[5] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *ArXiv e-prints*, 2020. 1

[6] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved Baselines with Momentum Contrastive Learning. *ArXiv e-prints*, 2020. 2

[7] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 4, 8, 14, 15

[8] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *ArXiv e-prints*, 2017. 3

[9] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-scale Hierarchical Image Database. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 4

[10] Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, and Chuang Gan. When does contrastive learning preserve adversarial robustness from pretraining to finetuning? In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2

[11] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal. Strip: A defence against trojan attacks on deep neural networks. In

*Proceedings of Annual Computer Security Applications Conference (ACSAC)*, 2019. 3

[12] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent - a new approach to self-supervised learning. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1, 2, 4, 8

[13] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *ArXiv e-prints*, 2017. 3

[14] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4

[16] S Hemalatha, U Dinesh Acharya, and A Renuka. Comparison of secure and high capacity color image steganography techniques in rgb and ycbcr domains. *ArXiv e-prints*, 2013. 4

[17] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 1, 2, 8

[18] Xijie Huang, Moustafa Alzantot, and Mani Srivastava. Neuroninspect: Detecting backdoors in neural networks via output explanations. *ArXiv e-prints*, 2019. 3

[19] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and< 0.5 mb model size. *ArXiv e-prints*, 2016. 6

[20] Yujie Ji, Xinyang Zhang, Shouling Ji, Xiapu Luo, and Ting Wang. Model-Reuse Attacks on Deep Learning Systems. In *Proceedings of ACM SAC Conference on Computer and Communications (CCS)*, 2018. 3

[21] Jinyuan Jia, Yupei Liu, and Neil Zhenqiang Gong. Badencoder: Backdoor attacks to pre-trained encoders in self-supervised learning. In *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, 2021. 1, 3

[22] Ziyu Jiang, Tianlong Chen, Ting Chen, and Zhangyang Wang. Robust pre-training by adversarial contrastive learning. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[23] Soheil Kolouri, Aniruddha Saha, Hamed Pirsiavash, and Heiko Hoffmann. Universal litmus patterns: Revealing backdoor attacks in cnns. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3

[24] Alex Krizhevsky and Geoffrey Hinton. Learning Multiple Layers of Features from Tiny Images. *Technical report, University of Toronto*, 2009. 4

[25] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research). 4

[26] Changjiang Li, Shouling Ji, Haiqin Weng, Bo Li, Jie Shi, Raheem Beyah, Shanqing Guo, Zonghui Wang, and Ting Wang. Towards certifying the asymmetric robustness for neural networks: quantification and applications. *IEEE Transactions on Dependable and Secure Computing*, 19(6):3987–4001, 2021. 8

[27] Changjiang Li, Haiqin Weng, Shouling Ji, Jianfeng Dong, and Qinming He. Det: Defending against adversarial examples via decreasing transferability. In *Cyberspace Safety and Security: 11th International Symposium, CSS 2019, Guangzhou, China, December 1–3, 2019, Proceedings, Part I 11*, pages 307–322. Springer, 2019. 8

[28] Hongbin Liu, Jinyuan Jia, and Neil Zhenqiang Gong. Poisonedencoder: Poisoning the unlabeled pre-training data in contrastive learning. In *Proceedings of USENIX Security Symposium (SEC)*, 2022. 1, 3, 4

[29] Yingqi Liu, Wen-Chuan Lee, Guanhong Tao, Shiqing Ma, Yousra Aafer, and Xiangyu Zhang. Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of ACM Conference on Computer and Communications (CCS)*, 2019. 3

[30] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In *Proceedings of Network and Distributed System Security Symposium (NDSS)*, 2018. 3

[31] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2018. 6

[32] Ren Pang, Changjiang Li, Zhaohan Xi, Shouling Ji, and Ting Wang. The dark side of automl: Towards architectural backdoor search. In *The Eleventh International Conference on Learning Representations*, 2022. 3

[33] Ren Pang, Hua Shen, Xinyang Zhang, Shouling Ji, Yevgeniy Vorobeychik, Xiapu Luo, Alex Liu, and Ting Wang. A tale of evil twins: Adversarial inputs versus poisoned models. In *Proceedings of ACM Conference on Computer and Communications (CCS)*, 2020. 2, 3

[34] Ren Pang, Zheng Zhang, Xiangshan Gao, Zhaohan Xi, Shouling Ji, Peng Cheng, and Ting Wang. Trojanzoo: Towards unified, holistic, and practical evaluation of neural backdoors. In *Proceedings of IEEE European Symposium on Security and Privacy (Euro S&P)*, 2020. 3

[35] Pengyu Qiu, Xuhong Zhang, Shouling Ji, Changjiang Li, Yuwen Pu, Xing Yang, and Ting Wang. Hijack vertical federated learning models with adversarial embedding. *arXiv preprint arXiv:2212.00322*, 2022. 8

[36] Md Rahman et al. A dwt, dct and svd based watermarking technique to protect the image piracy. *ArXiv e-prints*, 2013. 4

[37] Aniruddha Saha, Ajinkya Tejankar, Soroush Abbasi Koohpayegani, and Hamed Pirsiavash. Backdoor attacks on self-supervised learning. *ArXiv e-prints*, 2021. 1, 3, 4, 9

[38] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 6

[39] Naren Sarayu Manoj and Avrim Blum. Excess capacity and backdoor poisoning. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 13

[40] Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 3

[41] Mahesh Subedar, Nilesh Ahuja, Ranganath Krishnan, Ibrahima J Ndiour, and Omesh Tickoo. Deep probabilistic models to detect data poisoning attacks. *ArXiv e-prints*, 2019. 3

[42] Di Tang, XiaoFeng Wang, Haixu Tang, and Kehuan Zhang. Demon in the variant: Statistical analysis of dnns for robust backdoor contamination detection. In *Proceedings of USENIX Security Symposium (SEC)*, 2020. 2, 3, 8, 13

[43] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 3

[44] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *ArXiv e-prints*, 2019. 3

[45] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(11), 2008. 13

[46] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, 2019. 2, 3

[47] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 707–723. IEEE, 2019. 8

[48] Tong Wang, Yuan Yao, Feng Xu, Shengwei An, Hanghang Tong, and Ting Wang. An invisible black-box backdoor attack through frequency domain. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2022. 4

[49] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Proceedings of IEEE Conference on Machine Learning (ICML)*, 2020. 7, 13

[50] Yifei Wang, Qi Zhang, Yisen Wang, Jiansheng Yang, and Zhouchen Lin. Chaos is a ladder: A new theoretical understanding of contrastive learning via augmentation overlap. *ArXiv e-prints*, 2022. 7

[51] Jun Wu, Xuesong Ye, and Yanyuet Man. Bottrinet: A unified and efficient embedding for social bots detection via metric learning. In *2023 11th International Symposium on Digital Forensics and Security (ISDFS)*, pages 1–6. IEEE, 2023. 8

[52] Zhaohan Xi, Tianyu Du, Changjiang Li, Ren Pang, Shouling Ji, Xiapu Luo, Xusheng Xiao, Fenglong Ma, and Ting Wang. On the security risks of knowledge

graph reasoning. *arXiv preprint arXiv:2305.02383*, 2023. 3

[53] Xiaojun Xu, Qi Wang, Huichen Li, Nikita Borisov, Carl A Gunter, and Bo Li. Detecting ai trojans using meta neural analysis. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 103–120. IEEE, 2021. 8

[54] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y. Zhao. Latent backdoor attacks on deep neural networks. In *Proceedings of ACM Conference on Computer and Communications (CCS)*, 2019. 2, 3

[55] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2018. 3

[56] Shihao Zhao, Xingjun Ma, Xiang Zheng, James Bailey, Jingjing Chen, and Yu-Gang Jiang. Clean-label backdoor attacks on video recognition models. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3

[57] Stephan Zheng, Yang Song, Thomas Leung, and Ian Goodfellow. Improving the robustness of deep neural networks via stability training. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pages 4480–4488, 2016. 8

[58] Yuanyi Zhong, Haoran Tang, Junkun Chen, Jian Peng, and Yu-Xiong Wang. Is self-supervised learning more robust than supervised learning? *arXiv preprint arXiv:2206.05259*, 2022. 1, 2, 8

# A. Characterizing Supervised Backdoor Attacks

In supervised learning, the backdoor attack associates the trigger $r$ with the target label $t$ via (implicitly) minimizing the objective defined in Eq (1). The success of this attack is often attributed to the model's excess capacity [39], which "memorizes" both the function that classifies clean inputs and that misclassifies trigger inputs. Note that Eq (1) does not specify any constraints on the representations of trigger inputs. Thus, while associated with the same class, the trigger-embedded and target-class inputs are not necessarily proximate in the feature space.
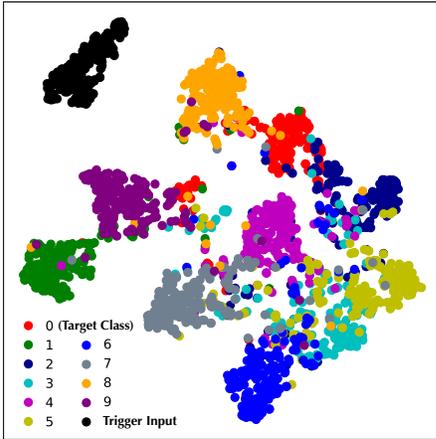


Figure 12: $t$-SNE visualization of the features of trigger-embedded and clean inputs in the supervised backdoor attack (target-class input: red; trigger-embedded input: black).

To validate the analysis, we use CTRL to generate the poisoning data, pollute 1% of the training data across all the classes, and train the model in a supervised setting for 20 epochs, which achieves 100% ASR and 83% ACC. We use $t$-SNE [45] to visualize the representations of trigger-embedded and clean inputs in the test set, with results shown in Figure 12. Although the clusters of trigger inputs (black) and target-class inputs (red) are assigned the same label, they are well separated in the feature space, indicating that supervised backdoor attacks do not necessarily associate the trigger with the target class in the feature space. This finding also corroborates prior work [42].

For comparison, we perform CTRL against SimCLR on CIFAR-10 and use $t$-SNE to visualize the features of trigger-embedded inputs and clean inputs in the test set, with results shown in Figure 13. Observed that in comparison with the supervised backdoor attack (*cf.* Figure 12), the cluster of trigger-embedded inputs (black) and the cluster of target-class clean inputs (red) are highly entangled in the feature space, indicating that the self-supervised backdoor attack takes effect by aligning the representations of trigger-embedded inputs and the target class.

# B. Proofs

We first introduce the following two assumptions commonly observed in encoders trained using SSL [49]

**Assumption B.1.** (Alignment) A well-trained encoder $f$ tends to map a positive pair to similar features. Formally, for given input $x$, $f(x)^\intercal f(x^+) \geq (1 - \epsilon)$, where $\epsilon \in [0, 1)$ is a small non-negative number. In particular, by design, the trigger $r$ is invariant to the augmentation operator: $f(r)^\intercal f(r^+) = 1$.

**Assumption B.2.** (Uniformity) A well-trained encoder $f$ tends to map inputs uniformly on the unit hyper-sphere $\mathcal{S}^{d-1}$ of the feature space, preserving as much information of the data as possible. Thus, as the number of data points is large, the average angle $\theta$ between the features of a negative pair follows the distribution density function [1]:

$$h(\theta) = \frac{1}{\sqrt{\pi}} \frac{\mathrm{T}(\frac{d}{2})}{\mathrm{T}(\frac{d-1}{2})} (\sin\theta)^{d-2}, \quad \theta \in [0, \pi] \qquad (4)$$

As the dimension $d$ is high, most of the angles heavily concentrate around $\pi/2$.
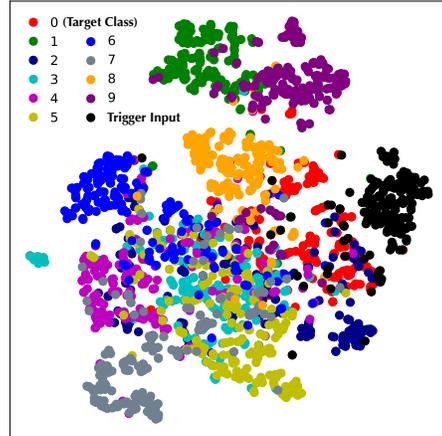
Based on B.1 and B.2, we now prove Theorem 5.1.



Figure 13: $t$-SNE visualization of the features of trigger-embedded and clean inputs in the self-supervised backdoor attack (target-class input: red; trigger-embedded input: black).

*Proof of Theorem 5.1.* According to Assumption B.1, we have

$$f(x_\star)^\intercal f(x_\star^+) \geq (1 - \epsilon) \qquad (5)$$

In other words,

$$(1 - \alpha)^2 f(x)^\intercal f(x^+) + \alpha^2 f(r)^\intercal f(r^+) +$$
$$\alpha(1 - \alpha)(f(x)^\intercal f(r^+) + f(r)^\intercal f(x^+)) \geq (1 - \epsilon).$$

Since both $f(x)^\intercal f(x^+)$ and $f(r)^\intercal f(r^+)$ are no larger than 1, we have

$$\alpha(1-\alpha)(f(x)^\intercal f(r^+) + f(r)^\intercal f(x^+))$$
$$\geq (1-\epsilon) - (1-\alpha)^2 f(x)^\intercal f(x^+) - \alpha^2 f(r)^\intercal f(r^+)$$
$$\geq (1-\epsilon) - (1-\alpha)^2 - \alpha^2$$
$$= 2\alpha(1-\alpha) - \epsilon$$

Then, based on Assumption B.1, we have

$$f(x)^\intercal f(r) \geq 1 - \frac{\epsilon}{2\alpha(1-\alpha)}. \tag{6}$$

For $\tilde{x}_*$ and $x$, we have

$$f(\tilde{x}_*)^\intercal f(x) = (1-\alpha)f(\tilde{x})^\intercal f(x) + \alpha f(r)^\intercal f(x)$$

Since $\tilde{x}$ and $x$ are a negative pair, based on Assumption B.2, we have

$$\mathbb{E}[f(\tilde{x}_*)^\intercal f(x)] \geq \alpha(1 - \frac{\epsilon}{2\alpha(1-\alpha)})$$
$$\geq \alpha - \frac{\epsilon}{2(1-\alpha)} \tag{7}$$

For a well-trained $f$, $\epsilon$ is a constant. Thus, both Eq (6) and Eq (7) are functions that first increase and then decrease with respect to $\alpha \in (0,1)$.

$\square$

## C. Details of Experimental Setting

**Dataset –** For each dataset, we split it as a training set and a testing set according to its default setting. Specifically, both CIFAR-10 and CIFAR-100 are split into 50,000 and 10,000 images for training and testing, respectively; ImageNet-100 is split as 130,000 training and 5,000 testing images; while GTSRB is split into 39,209 training and 12,630 testing images.

**Data augmentation –** For convenience, we describe the details of data augmentations in a PyTorch style. Specifically, following prior work [7, 4], we use geometric augmentation operators including *RandomResizeCrop* (of scale [0.2, 1.0]) and *RandomHorizontalFlip*. Besides, we use *ColorJitter* with [brightness, contrast, saturation, hue] of strength [0.4., 0.4, 0.4, 0.1] with an application probability of 0.8 and *RandomGrayscale* with an application probability of 0.2.

**Encoder training –** We use the training set of each dataset to conduct contrastive learning. We show the hyper-parameters setting for each contrastive learning algorithm in Table 7, which is fixed across all the datasets.

**Classifier training –** Without explicit specification, we randomly sample 50 examples from each class of the corresponding testing set to train the downstream classifier. We show the hyper-parameters of classifier in Table 8.

| Hyper-parameter | SSL Method | | |
|---|---|---|---|
| | SimCLR | BYOL | SimSiam |
| Optimizer | SGD | SGD | SGD |
| Learning Rate | 0.5 | 0.06 | 0.06 |
| Momentum | 0.9 | 0.9 | 0.9 |
| Weight Decay | 1e-4 | 1e-4 | 5e-4 |
| Epochs | 500 | 500 | 500 |
| Batch Size | 512 | 512 | 512 |
| Temperature | 0.5 | - | - |
| Moving Average | - | 0.996 | - |

Table 7. Hyper-parameters of encoder training.

| Hyper-parameter | Setting |
|---|---|
| Optimizer | SGD |
| Batch Size | 512 |
| Learning Rate | 0.2 |
| Momentum | 0.9 |
| Scheduler | Cosine Annealing |
| Epochs | 20 |

Table 8. Hyper-parameters of classifier training.

**Evaluation –** We evaluate ACC using the full testing set. For ASR, we apply CTRL on the full testing set and measure the ratio of trigger inputs that are classified to the target class. All the experiments are performed on a workstation equipped with Intel(R) Xeon(R) Silver 4314 CPU @ 2.40GHz, 512GB RAM, and four NVIDIA A6000 GPUs.

## D. More Experimental Results

Here, we show the additional experimental results.

**Performance of clean models –** The ACC and ASR of clean models trained by SimCLR, BYOL, and Simsiam are summarized in Table 9.

| Dataset | SSL Method | | | | | |
|---|---|---|---|---|---|---|
| | SimCLR | | BYOL | | SimSiam | |
| | ACC | ASR | ACC | ASR | ACC | ASR |
| CIFAR-10 | 79.1% | 9.93% | 82.4% | 12.2% | 81.5% | 11.75% |
| CIFAR-100 | 48.1% | 1.14% | 51.0% | 0.46% | 52.0% | 0.72% |
| ImageNet-100 | 42.2% | 1.59% | 45.1% | 1.41% | 41.3% | 1.53% |

Table 9. Accuracy of different SSL methods under normal training.

**Fine-tuning data size –** Typically, equipped with the pre-trained encoder, the victim fine-tunes the downstream classifier with a small labeled dataset. Here, we evaluate the impact of this fine-tuning dataset on CTRL. Figure 14 illustrates the performance of CTRL as a function of the number of labeled samples per class.

Observe that both ACC and ASR of CTRL increase with the fine-tuning data size, while their variance decreases gradually. For instance, when the number of labeled samples per class is set as 50, the ASR of CTRL on CIFAR-10 under SimSiam stably remains around 75%. This may be explained as follows. Without the supervisory signal of labeling, CTRL achieves effective attacks by entangling the representations
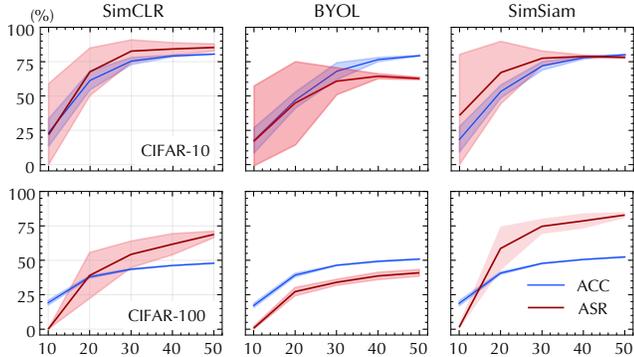
Figure 14: Performance of CTRL w.r.t. the fine-tuning data size.



Figure 16: Performance of CTRL w.r.t. the number of epochs.

of trigger-embedded and target-class inputs (details in § 5). During fine-tuning, more labeled samples imply that the representations of trigger-embedded inputs are more likely to be associated with the target-class label, leading to higher and more stable ASR. In other words, more fine-tuning data not only improves the model's performance but also increases its attack vulnerability.

**Batch size –** Existing studies show that batch size tends to impact the performance of contrastive learning [7]. Here, we explore its influence on the performance of CTRL. Specifically, on the CIFAR-10, we measure the ACC and ASR of CTRL with the batch size varying from 128 to 512, with results shown in Figure 15.
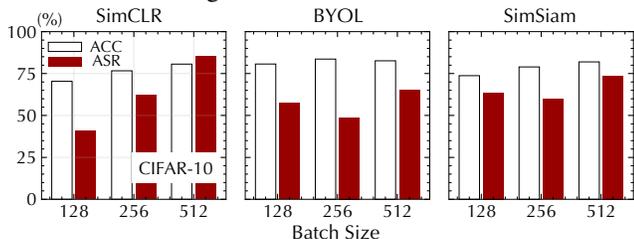


Figure 15: Performance of CTRL w.r.t. the batch size on CIFAR10.

Observe that the model's accuracy improves with the batch size, which corroborates the existing studies [7]. Moreover, a larger batch size (*e.g.*, ≥ 512) generally benefits the ASR of CTRL. This may be explained by that more positive pairs (also more negative pairs in SimCLR) in the same batch lead to tighter entanglement between trigger-embedded and target-class inputs. Meanwhile, for smaller batch sizes (*e.g.*, ≤ 256), the three SSL methods show slightly different trends. This may be attributed to the design of their loss functions: BOYL and SimSiam only optimize positive pairs, while SimCLR optimizes both positive and negative pairs, thereby gaining more benefits from larger batch sizes.

**Training epochs –** Typically, SSL benefits from more training epochs [4]. We evaluate the impact of training epochs on the performance of CTRL. Figure 16 shows the ACC and ASR of CTRL as the number of epochs varies from 600 to 1,000.

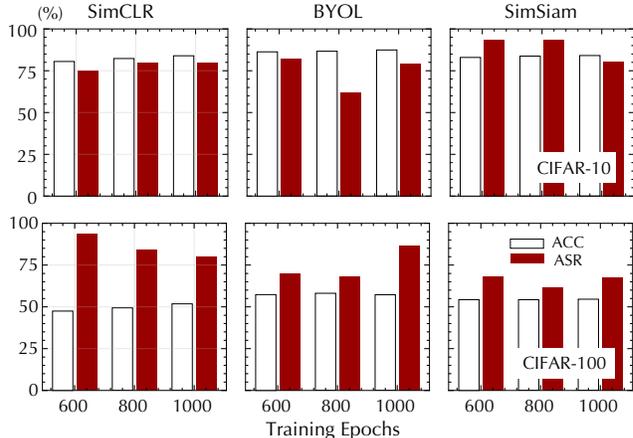Observe that as the training epoch increases, the ACC

of CTRL gradually grows, while the ASR remains at a high level. For example, on CIFAR-10 with SimCLR, when the number of epochs increases from 600 to 1,000, the ACC also increases from 80.52% to 83,94%, and the ASR remains above 75%. In a few cases, the ASR slightly drops. We speculate this is caused by the random data augmentations used in SSL. Side evidence is that on CIFAR-10 with BYOL, the ASR first slightly decreases and then remains above 80%. In general, the number of training epochs has a limited impact on the effectiveness of CTRL.