

The Perils of Learning From Unlabeled Data: Backdoor Attacks on Semi-supervised Learning

Virat Shejwalkar*, Lingjuan Lyu†, Amir Houmansadr*

*University of Massachusetts Amherst †Sony AI

{vshejwalkar, amir}@cs.umass.edu, †Lingjuan.Lv@sony.com

Abstract—Semi-supervised machine learning (SSL) is gaining popularity as it reduces the cost of training ML models. It does so by using very small amounts of (expensive, well-inspected) labeled data and large amounts of (cheap, non-inspected) unlabeled data. SSL has shown comparable or even superior performances compared to conventional fully-supervised ML techniques.

In this paper, we show that the key feature of SSL that it can learn from (non-inspected) unlabeled data exposes SSL to strong poisoning attacks. In fact, we argue that, due to its reliance on non-inspected unlabeled data, poisoning is a much more severe problem in SSL than in conventional fully-supervised ML.

Specifically, we design a *backdoor poisoning attack* on SSL that can be conducted by a *weak adversary* with no knowledge of target SSL pipeline. This is unlike prior poisoning attacks in fully-supervised settings that assume strong adversaries with practically-unrealistic capabilities. We show that by poisoning only 0.2% of the unlabeled training data, our attack can cause misclassification of more than 80% of test inputs (when they contain the adversary’s backdoor trigger). Our attacks remain effective across twenty combinations of benchmark datasets and SSL algorithms, and even circumvent the state-of-the-art defenses against backdoor attacks. Our work raises significant concerns about the practical utility of existing SSL algorithms.

I. INTRODUCTION

The more training data we use to train machine learning (ML) models, the better they perform [19], [18]. However, this makes the conventional *fully-supervised* ML significantly challenging, as it requires large amounts of *labeled* training data. Labeling is an expensive [17] and error prone process [38], [33] that makes conventional ML prohibitively expensive in practice, especially with today’s exploding training data sizes.

Semi-supervised learning (SSL) addresses this major challenge by significantly reducing the need of the labeled training data: SSL uses a combination of a *small, high-quality and expensive labeled data* with a *large, low-quality, and cheap unlabeled data* to train its models. For instance, the FixMatch [48] SSL algorithm uses only 40 labeled data along with about 50,000 unlabeled data and achieves 90% accuracy on CIFAR10. In contrast to fully-supervised algorithms, training an SSL algorithm involves two loss functions: a *supervised loss function* (e.g., cross-entropy [37]) on labeled training data and an *unsupervised loss function* (e.g., cross-entropy over pseudo-labels [29]) on unlabeled training data. Different SSL algorithms primarily differ on how they compute their unsupervised losses.

SSL is being explored extensively by both academia [60], [54], [55] and industry [48], [49], [7], [6], as recent SSL algorithms offer state-of-the-art performances comparable or

even superior to those achieved by conventional supervised techniques—but with no need of large well-inspected labeled data. For instance, with less than 10% of training data labeled, FixMatch [48] and FlexMatch [60] SSL algorithms outperform the fully-supervised algorithm. This is because state-of-the-art SSL algorithms use (cheap, abundant) unlabeled data much more effectively than how fully-supervised algorithms use significantly larger (expensive, scarce) labeled data.

Unlabeled data enables poisoning by weaker adversaries: Multiple researches have demonstrated the threat that data poisoning attacks pose to fully-supervised learning [25], [35], [41], [43], [58], [53]. However, as the training data of fully-supervised models undergo an extensive and careful inspection, these attacks assume strong adversaries with the knowledge of model parameters [35], training data [53], [8], [36] or its distribution [58], or the learning algorithm. Such strong adversaries are important to evaluate the worst-case security of a system, but are practically less relevant. On the other hand, the key feature of SSL in real-world applications is that it can leverage large amounts of—raw, non-inspected—unlabeled data, e.g., the data scraped from the Internet. We argue that **the use of non-inspected data by SSL presents a significant threat to its security, as it allows even the most naive adversaries (with no knowledge of training algorithm, labeled data, etc.) to poison SSL models** by simply fabricating malicious unlabeled data. Unfortunately, this severe threat remains largely unexplored in the SSL literature.

To address this gap, in this paper, we take the first step towards thoroughly understanding this threat by studying the possibility of *backdoor attacks* against SSL in real-world settings. A backdoor attack aims to install a *backdoor function* in the *target model*, such that the *backdoored target model* will *misclassify* any test input to the *adversary chosen target class* when patched with a specific *backdoor trigger*, but will correctly classify the input without the trigger.

Existing backdoor attacks fail on SSL: There exist numerous backdoor attacks in prior literature, however, except one attack—DeHiB [56], *all of the prior attacks consider fully-supervised settings*. Our preliminary evaluations show that all of the existing state-of-the-art attacks, including DeHiB, completely fail against SSL under our realistic threat model (Section III). Hence, to learn from these failures, we first systematically evaluate five backdoor attacks from three categories against five state-of-the-art SSL algorithms, under our practical, unlabeled data poisoning threat model. We adaptively choose the attack categories based on the specific lesson we learn from evaluating the prior category.

Our systematic evaluation leads to the following **three**

major lessons that not only guide our attack design, but are generally applicable to any (future) backdoor attacks against SSL: (1) *Backdoor attacks on SSL should be clean-label style attacks*, i.e., poisoning data should be selected from the target class distribution; (2) *Backdoor trigger patterns should span/cover the entire poisoning sample*, to circumvent strong augmentations, e.g., cutout [21], that all modern SSL algorithms use; (3) *Backdoor trigger patterns should be resistant to noise and with repetitive patterns* to withstand large amounts of random noises which were added to training data via the strong augmentation in SSL.

Our SSL-tailored backdoor method: The high-level intuition behind our backdoor attack is as follows. All modern semi-supervised algorithms learn via a self-feedback mechanism, called *pseudo-labeling*, i.e., the prediction \tilde{y} on an unlabeled sample \mathbf{x} has high confidence, they use (\mathbf{x}, \tilde{y}) as a labeled sample for further training. Following our first lesson, we exploit this pseudo-labeling and design a *clean-label* attack that poisons samples \mathbf{x} only from the distribution of the target class y^t . Our attack patiently *waits* for the target model to correctly label a poisoning sample $(\mathbf{x} + T)$ as y^t , where T is our backdoor trigger. And then, as the model trains further on $((\mathbf{x} + T), y^t)$, our attack *forces the model* to associate features of our simple trigger T , instead of the complex features of \mathbf{x} , with the target class, which effectively installs the backdoor function in the target model.

Note that, we consider the most challenging setting for designing attacks with the least capable and knowledgeable data poisoning adversary. Generally, trigger generation for data poisoning backdoor attacks is formalized as a bi-level optimization problem [36], however such attacks are well-known to be very expensive, and yet ineffective [36], [44]. Instead, we design a static, repetitive grid pattern for our backdoor trigger (Figure 5), as guided by lessons 2 and 3. In summary, we sample few data from target class, patch them with our trigger and inject into unlabeled training data.

Evaluations: We demonstrate the strength of our attack via an extensive evaluation against five state-of-the-art SSL algorithms and a fully-supervised algorithm, using four benchmark image classification tasks, that SSL literature commonly uses. We note that our attack significantly outperforms the prior attacks from both SSL and full-supervised literature.

For the most combinations of algorithm and dataset, **our attacks achieve high attack success rates (ASRs) (>80%), while poisoning just 0.2% of entire training data.** Comparatively, DeHib uses $20\times$ more poisoning data and achieves 0% ASRs. ASR measures the % of *test inputs from non-target classes* that the backdoored model classifies to the target class when patched with backdoor trigger. For instance, our attacks have more than 90% ASR against CIFAR100 and more than 80% ASR against CIFAR10. For SVHN and STL10, our attack has more than 80% ASR with two exceptions each. Through a systematic experiment design in Section VI-B, we show that our intuition aligns with the dynamics of our patient attacks and justify their strength. **Our attack is highly stealthy**, as (1) it minimally perturbs the poisoning data and (2) it produces backdoored models which have high accuracy (close to non-backdoored models) on non-backdoored test inputs.

Next, our comprehensive ablation study (Section VI-E)

shows the high efficacy of our attacks when varying three major parameters of our setting: size of labeled data, backdoor target class, and size of poisoning data. Finally, we show that **our attacks remain highly effective even when SSL is paired individually with five state-of-the-art defenses** against backdoor attacks that are agnostic to learning algorithms.

Summary of contributions:

- (1) We perform the first thorough study of backdoor attacks on semi-supervised learning (SSL), and show that it is highly susceptible to backdoor poisoning, under realistic unlabeled data poisoning threat models.
- (2) We systematically evaluate existing backdoor attacks from fully-supervised setting on SSL and provide concrete lessons to design stronger backdoor attacks against SSL.
- (3) Based on the lessons, we design the first effective backdoor attack against SSL that achieves high (>80%) ASRs by poisoning just 0.2% of entire training data.
- (4) We show that existing learning-algorithm-agnostic defenses are insufficient to defend SSL against backdoor attacks.

II. PRELIMINARIES AND RELATED WORK

In this section, we provide the preliminaries and important related works required to understand the rest of the paper.

A. Semi-supervised Learning (SSL)

The objective of machine learning (ML) for classification task is to train a classifier (e.g., neural network) with parameters θ and learn function $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$ to predict label $y \in \mathcal{Y}$ for input $\mathbf{x} \in \mathcal{X}$. Here, $\mathcal{X} \in \mathbb{R}^d$ is the input feature vector space and $\mathcal{Y} \in \mathbb{R}^k$ is the output space. Parameters θ are trained using empirical risk minimization (ERM) to minimize certain empirical loss function $\ell_{(\mathbf{x}, y) \in D}(f_\theta, (\mathbf{x}, y))$, where D is the training data and $D \subset \mathcal{X} \times \mathcal{Y}$. Traditionally, ML uses *fully-supervised* learning algorithms that use only completely labeled data, D^l . However, labeling is a manual, expensive, and error-prone process [33], as it requires human intervention. Consequently, with continuously increasing training data sizes, labeling cost for ML can become prohibitively high in practice.

Semi-supervised learning (SSL) addresses this major challenge by reducing the dependence of ML on labeled data. SSL proposes to learn ML models using both labeled D^l and unlabeled data D^u . D^l and D^u may or may not have the same distribution, but sizes of labeled data are significantly smaller than that of unlabeled data, i.e., $|D^l| \ll |D^u|$. A typical SSL loss function is a convex combination of a loss on D^l , denoted by \mathcal{L}_l , and a loss on D^u , denoted by \mathcal{L}_u : $\mathcal{L}_{ss} = \mathcal{L}_l + \lambda \mathcal{L}_u$. \mathcal{L}_l is generally the standard cross-entropy loss due to its high performances. But, \mathcal{L}_u varies across different SSL algorithms; we will discuss these shortly.

Traditionally, semi-supervised learning has been largely ineffective, but in the past couple of years, SSL has significantly improved, especially after the invention of MixMatch [7]. Significant performance gains of MixMatch are because it combines various data augmentation techniques with various prior SSL algorithms, including *pseudo-labeling* [29], entropy minimization [24], and *consistency regularization* [20], [42], [28]. Hence, in this work we consider the state-of-the-art algorithms that use pseudo-labeling and consistency regularization.

We briefly describe these two techniques as in [48] followed by the semi-supervised algorithms we consider in this work.

1) *Pseudo-labeling*: Pseudo-labeling uses the current model to obtain artificial *pseudo-labels* for the unlabeled data. There are various ways in which model’s outputs can be used to train itself, e.g., MixMatch and ReMixMatch compute model’s predictions on an unlabeled sample and then use entropy minimization [24] to sharpen the prediction. But, pseudo-labeling specifically refers to the use of “hard” labels (i.e., the argmax of the model’s output) and only retains the labels whose largest class probability (confidence) is above a pre-defined threshold. Assume $q_b = f_\theta(y|u_b)$ are the predictions of current model f_θ on the batch u_b of unlabeled data. Then pseudo-labeling loss can be formalized as follows:

$$\frac{1}{|u_b|} \sum_{b=1}^{|u_b|} \mathbb{1}(\max(q_b) \geq \tau) H(\hat{q}_b, q_b) \quad (1)$$

where $\hat{q}_b = \text{argmax}(q_b)$, $H(\cdot)$ is the cross-entropy loss function, and τ is the threshold parameter. Note that the use of hard labels is similar to entropy minimization in MixMatch or ReMixMatch, where they use a temperature parameter to sharpen the model predictions to have low entropy (high confidence).

2) *Consistency regularization*: Consistency regularization is commonly used in modern semi-supervised learning algorithms. It is based on the intuition that the model should output similar predictions when input with the perturbed versions of the same image. The idea was first proposed in [3]. Semi-supervised learning algorithms use consistency regularization on unlabeled data as follows. They use a stochastic augmentation mechanism $a(\mathbf{x}_u)$ to produce perturbed versions of an unlabeled sample \mathbf{x}_u and then force the model to have similar outputs on these versions using the following loss:

$$\sum_{b=1}^{|u_b|} \|f_\theta(y|a(u_b)) - f_\theta(y|u_b)\|_2^2 \quad (2)$$

where, $a(\cdot)$ is a stochastic function, hence it produces different output every time it is applied to a batch u_b of unlabeled data. Consequently, the two terms in (2) have different values. Next, we briefly describe the semi-supervised learning algorithms we consider in this work. For detailed description of the algorithms, please refer to the original works.

(1) *MixMatch* [7] combines various prior semi-supervised learning techniques. For an unlabeled sample, MixMatch generates K weakly augmented versions of the unlabeled sample, computes outputs of the current model f_θ for the K versions, averages them, and sharpens the average prediction by raising all its probabilities by a power of $1/\text{temperature}$ and re-normalizing; it uses the sharpened prediction as the label of the unlabeled sample. Finally, it uses mixup regularization [61] on the combination of labeled and unlabeled data and trains the model using cross-entropy loss.

(2) *Unsupervised data augmentation (UDA)* [54] shows significant improvements in semi-supervised performances by just replacing the simple weak augmentations of MixMatch with a strong augmentation called Randaugment [16]. In an iteration, Randaugment randomly selects a few augmentations from a large set of augmentations and applies them to images.

(1) *ReMixMatch* [7] builds on MixMatch by making multiple modifications, including 1) it replaces the simple weak augmentation in MixMatch with Autoaugment [15], 2) it uses augmentation anchoring to improve consistency regularization, i.e., it uses the prediction on a weakly augmented version of unlabeled sample as the target prediction for a strongly augmented version of the unlabeled sample, and 3) it uses distribution alignment, i.e., it normalizes the new model predictions on unlabeled data using the running average of model predictions on unlabeled data. This significantly boosts the performance of resulting model.

(4) *FixMatch* [48] simplifies the complex ReMixMatch algorithm by proposing to use a combination of Pseudo-labeling and consistency regularization based on augmentation anchoring (discussed above). FixMatch significantly improves semi-supervised algorithms, especially in the low labeled data regimes.

(5) *FlexMatch* [60] proposes curriculum pseudo labelling (CPL) approach to leverage unlabeled data according to model’s learning status. The main idea behind CPL is to flexibly adjust the thresholds used for pseudo-labeling for different classes at each training iteration in order to select more information unlabeled data and their pseudo-labels. CPL can be combined with other algorithms, e.g., UDA.

B. Backdoor Attacks

A *backdoor adversary* aims to implant a *backdoor functionality* into a *target* model. That is, given an input (\mathbf{x}, y^*) with true label y^* , the *backdoored target model* f_θ^b should output an adversary-desired *backdoor target label* y^t for the input patched with a pre-specified *backdoor trigger* T , but it should output the correct label for the benign input, i.e., $f_\theta^b(\mathbf{x} + T) \mapsto y^t$ and $f_\theta^b(\mathbf{x}) \mapsto y^*$. There are two major types of backdoor attacks: *dirty-label* and *clean-label* backdoors.

1) *Dirty-label backdoor attacks* [26], [13], [43], [59], [40], [32]: These attacks poison both the features \mathbf{x} and labels y^* of benign, labeled data to obtain poisoning data D^p . They first select some benign data $(X, Y \setminus y^t)$ from non-target classes of the original training data D , patch the backdoor trigger to X : $X^p \leftarrow X + T$, and set labels of X^p to y^t to obtain $D \leftarrow D \cup D^p = (X^p, y^t)$. Training f_θ on such D makes the model *associate* the trigger with the target label, i.e., $f_\theta(T) \mapsto y^t$, as this association is much easier to learn than learning to associate original X to Y . However, in many practical scenarios, a trusted third party conducts the data labeling and inspection, and can easily remove such mislabeled data. Our work focuses on the semi-supervised learning where such inspection is much easier compared to the supervised learning due to very small sizes of labeled data, hence we only poison the unlabeled training data.

2) *Clean-label backdoor attacks* [53], [58], [62]: These attacks poison only the features X of benign data. They add imperceptible T to X such that the poisoned features X^p appear to be from the respective true classes to a human. Clean-label attacks can be further divided into two categories [58] based on the true classes of X that they poison: *feature-collision* attacks and *target-class* attacks.

(a) *Feature-collision backdoor attacks*, e.g., HTBA [41] and SAA [50], insert triggers indirectly. They try to match

the feature space/gradients between the target class samples and non-target-class samples patched with the trigger, thus, mimicking the effects of non-target-class poisoning. By doing this, the decision boundary will place these two points in proximity in the feature space, and as a result, any input with the trigger will likely be classified into the target class. However, these attacks can backdoor just one sample at a time, and hence, are computationally inefficient at backdooring large portions of test inputs.

(b) *Target-class backdoor attacks* select X to poison from the target class y^t . For instance, label-consistent (LC) backdoor attacks [53] select a few X from y^t and manipulate X to make their original features harder to learn. Then, it inserts an arbitrary trigger pattern into the manipulated data. They either use GANs or adversarial samples to manipulate X to obtain difficult-to-learn X' and then add T to get X^p , i.e., $X^p \leftarrow X' + T$. LC attack requires in-distribution data from both target and non-target classes in order to train adversarial example generator. Narcissus [58] attack addresses the above issues as it requires in-distribution data only from the target class, y^t . For clarity of presentation, we discuss more details of the Narcissus attack in Section IV-A3 and evaluate them against semi-supervised learning.

3) *Backdoor attacks on semi-supervised learning*: So far, most of the backdoor attack literature has focused on the fully-supervised settings. Only one work by Zhicong et al. [56] study backdoor attacks against semi-supervised learning. For clarity of presentation, we discuss the details of this attack in Section IV-A1, where we demonstrate and justify why this attack fails to backdoor semi-supervised learning.

III. THREAT MODEL

Below, we discuss the threat model of our backdoor attacks in terms of the adversary’s goal, knowledge, and capabilities. We consider a setting where a victim model trainer collects data from multiple, potentially untrusted sources to train a ML model for a classification task with C classes. Below, $[I]$ denotes the set of all non-zero positive integers $\leq I$.

A. Adversary’s goal

We consider a *backdoor adversary* who aims to install a *backdoor function* in the victim’s ML model, called *target model*. We denote the function of a *benign model*, i.e., without any backdoor by f_θ and that of a *backdoored (target) model* by f_θ^b . Our adversary’s goal is two-fold.

1) *Backdoor goal*: The backdoor adversary selects a *backdoor target class* $y^t \in [C]$. To mount an effective backdoor attack, the backdoor goal requires the backdoored model to *incorrectly classify all the test inputs from non-target classes to the target class, when they are patched with a pre-specified backdoor trigger, T* . More formally, $f_\theta^b(\mathbf{x} + T) \mapsto y^t \forall (\mathbf{x}, y^*)$ where $y^* \in [C] \setminus y^t$ is the true class of \mathbf{x} .

2) *Stealth goal*: In order to mount a stealthy backdoor attack, the stealth goal requires the backdoored model, f_θ^b , to retain all the benign functionalities similar to the benign model, f_θ . Specifically, f_θ^b should correctly classify all the benign test inputs from all the classes without the backdoor trigger. Formally, $f_\theta^b(\mathbf{x}) = f_\theta(\mathbf{x}) \mapsto y^* \forall (\mathbf{x}, y^*)$ where $y^* \in [C]$.

B. Adversary’s knowledge

As discussed in Section I, we consider the most naive, real-world adversary with minimum knowledge of the semi-supervised learning (SSL) pipeline. We assume that the adversary has no knowledge of the data except the specific classes of the classification task. Next, the adversary knows the details of the target SSL algorithm, but do not know model architecture, e.g., ResNet or VGG, i.e., *our attacks are model architecture agnostic*. Finally, we assume that the adversary does not know the distribution of the unlabeled and labeled training data, or posses any data from the true distribution.

C. Adversary’s capabilities

We discuss the adversary’s capabilities based on their ability to manipulate the training pipeline and training data.

1) *Manipulating training pipeline*: There are three types of poisoning attacks based on the part of SSL training pipeline that the adversary can manipulate: *data poisoning*, *code poisoning*, and *model poisoning*. The *model poisoning* [5], [35], [45] adversary is the strongest adversary who directly manipulates the model parameters. Such poisoning requires highly privileged accesses to the training platforms, which is impractical in many real-world settings, e.g., for popular ML platforms like Amazon AWS [44]. The *code poisoning* [4], [23] adversary poisons the code of the target algorithms and also requires appropriate permissions to change the code along with a thorough knowledge of the learning pipeline.

Finally, the *data poisoning* [9], [36] adversary can only manipulate the training data of the target model. Due to its naivety, it is also the most practical adversary who can be any data owner willing to contribute data to SSL. Data poisoning is a severe threat because it is easy to deploy [44], even against sophisticated ML platforms, e.g., Amazon AWS, Google cloud and Microsoft Azure, with state-of-the-art software security. Hence, *we consider the data poisoning adversary in this work*.

2) *Manipulating training data*: We discuss manipulation of training data separately from that of training pipeline, especially because SSL uses two types of datasets. SSL bootstraps knowledge from a small labeled dataset D^l , hence D^l is of very high quality and is well-inspected. Therefore, we argue that an adversary who poisons or has access to D^l [56] is not practically relevant. On the other hand, a salient feature of semi-supervised learning is that the model trainer needs not to inspect its unlabeled data D^u at all (Section II-A), and hence, D^u can be easily poisoned in practice. Hence, we assume that *our adversary can poison only the unlabeled training data*.

IV. OUR ATTACK METHODOLOGY

In this section, we discuss our backdoor attack methodology tailored to the unlabeled data poisoning threat model (Section III). In Section IV-A, we present the first systematic evaluation of existing backdoor attacks in semi-supervised learning (SSL) settings and provide three major lessons. Next, we give the intuition behind our backdoor attack (Section IV-B), and finally we detail our attack method in Section IV-C.

Table I: Left-most column shows types of backdoor attacks based on specific characteristics, middle column lists existing attacks of each type. Right-most column presents lessons we learn from evaluating one/two representative attacks (in bold) of each type.

Attack characteristic/ type	Existing attacks of given type	Lesson from evaluations
Dirty label	DeHiB [56], DL-Badnets [25], DL-Blend [13], Facehack [43]	Attack should be a clean-label attack, i.e., poisoning samples should be from backdoor target class.
Clean-label small trigger	CL-Badnets [58], CL-Blend [13]	Trigger should span the entire sample/image to avoid cropping/covering by strong augmentations.
Clean-label adversarial samples	Narcissus [58], Label consistent [53], non-repeating trigger patterns , HTBA [41], SAA [50], Embedding [62]	Trigger should be noise-resistant and its pattern should be repetitive so that even a part of trigger can install a backdoor.

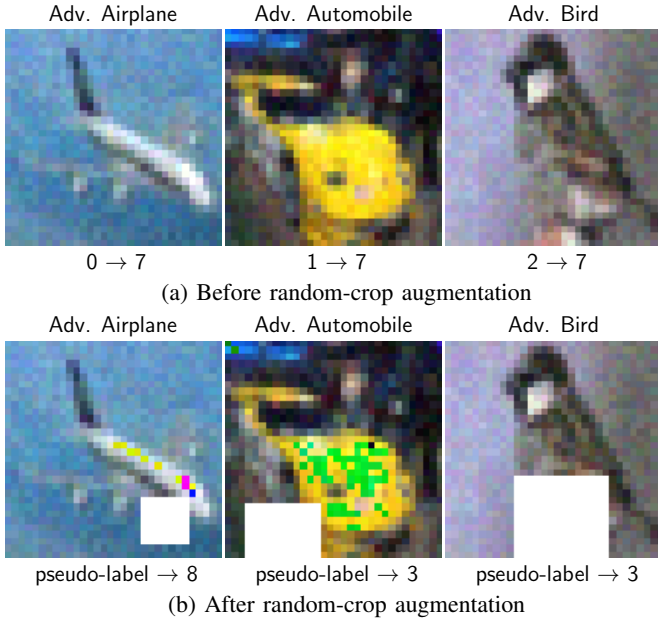


Figure 1: DeHiB [56] fails because it cannot obtain the target class as pseudo-labels for its poisoning data.

Table II: Impacts of existing backdoor attacks (Section II-B) on various semi-supervised algorithms for CIFAR10 data. We poison 0.2% (100 samples) of all the training data. DeHiB* is the original attack with the knowledge of labeled training data D^l while DeHiB is the attack without the knowledge of D^l .

Algorithm	DeHiB* ASR (%)	DeHiB ASR (%)	CL-Badnets ASR (%)	LC ASR (%)	Narcissus ASR (%)
Mixmatch [7]	22.0	1.0	9.1	1.1	2.2
Remixmatch [6]	10.9	0.9	0.0	0.0	0.0
UDA [54]	21.2	1.2	5.1	0.0	0.0
Fixmatch [48]	35.8	0.9	10.2	0.1	1.3
Flexmatch [60]	16.9	1.2	9.1	0.1	1.1

A. Systematic evaluation of existing backdoor attacks

Previous works have proposed numerous backdoor attacks under different threat models. But all works, except DeHiB [56], consider fully-supervised setting. Hence, we first present a systematic evaluation of existing state-of-the-art backdoor attacks and explain why they fail in SSL settings. Based on our evaluations, we provide three major lessons that are fundamental to our attack design and generally apply to any (future) backdoor attacks against semi-supervised learning.

We start our evaluations from DeHiB [56], the only existing backdoor attack on semi-supervised learning, and based on the lessons learned from this evaluation, we chose the next type of attacks to evaluate. As we see from Table I, each of our lessons applies to multiple backdoor attacks of a specific type and

characteristics. However, for conciseness, we evaluate one or two representative attacks from each type and provide lesson/s that are useful in designing stronger attacks.

1) *Attacks should be clean-label attacks:* We first evaluate *Deep hidden backdoor* (DeHiB) [56] attack. DeHiB poisons only the unlabeled data, D^u , but it assumes a strong, unrealistic adversary who can access the labeled data, D^l . It first samples some data (X, Y) from both target, y^t , and non-target, $y^{\setminus t}$, classes. Then it uses a model trained on D^l to add universal adversarial perturbation \mathcal{P}_t to X such that the perturbed data $X + \mathcal{P}_t \mapsto X^p$ is classified as y^t ; as we only poison D^u , we denote poisoning data by X^p . Finally, it adds a static trigger T to the perturbed data X^p . Intuition behind DeHiB is that, due to \mathcal{P}_t , SSL algorithm will assign target class y^t as pseudo-labels to all X^p and force the target model to associate static trigger T to y^t and ignore original features X .

Why does DeHiB fail? Recall from Section II-A that all of state-of-the-art SSL algorithms use various strong augmentations, including, cutout [21], adding various types of hue [47], horizontal/vertical shifts [52], etc. Next, note that adversarial perturbations are sensitive to noises [2], i.e., even moderate changes in the perturbations render them ineffective. Hence, in presence of strong augmentations, adversarial perturbations fail to obtain the backdoor target class y^t as the pseudo-labels for X^p of DeHiB as shown in Figure 1. Hence, the very fundamental requirement of DeHiB does not hold in SSL and leads to its failure. The original DeHiB work reports slightly better results, because it assumes access to D^l , which our threat model does not allow. Hence, we use randomly sampled data of size $|D^l|$ from entire CIFAR10 data to obtain DeHiB’s \mathcal{P}_t .

To summarize, adversarial perturbations are sensitive to noises. Hence, using adversarial samples from non-target classes as poisoning samples cannot guarantee the desired pseudo-labeling to y^t . Effectively, such attack tries to train the model to associate the trigger pattern T with multiple labels, and hence, fails to inject the backdoor functionality. For the same reason, *we also observed that any dirty-label static trigger attacks completely fail against SSL*. Hence, backdoor attacks on SSL should be clean-label attacks, i.e., use poisoning samples X^p from y^t , and leverage benign features of X^p to obtain desired pseudo-labels y^t for them.

Lesson-1: Backdoor attacks on semi-supervised learning should be clean-label style attacks, which sample their poisoning samples from the backdoor target class.

2) *Backdoor trigger should span the whole sample:* Based on Lesson-1, we choose to evaluate clean-label attacks. But, we consider small trigger pattern attacks to emphasize the importance of the trigger sizes towards attack efficacy against

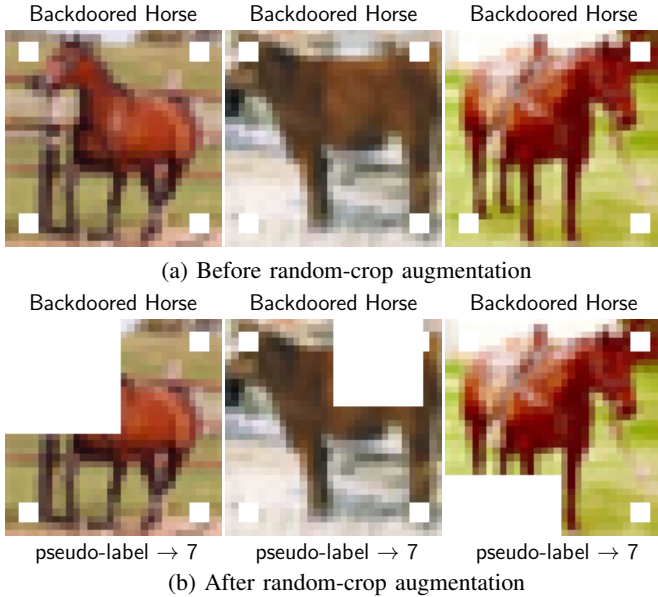


Figure 2: Clean-label Badnets [26] obtains the target class as pseudo-labels for its poisoning data, but cutout augmentation occludes its small trigger and renders it ineffective.

semi-supervised learning. In particular, we evaluate *clean-label Badnets* (CL-Badnets) [58] attack, which adds a static trigger, e.g., a pixel pattern with single/multiple squares, to the samples X from the target class, y^t to get poisoning data X^p . It then injects X^p into the unlabeled training data D^u .

Why does CL-Badnets fail? This clean-label style attack ensures that the model assigns y^t to all the poisoning samples. However, all the semi-supervised algorithms use a strong augmentation technique called *random-crop* (or cutout) that randomly crops a part of a sample. Because of this, the trigger is generally absent in many of the augmented instances of a poisoning sample as shown in Figure 2. This majorly reduces the impact of this attack as our results show in Tables II and IV.

Lesson-2: To ensure that all the augmented instances of a poisoning sample contain the backdoor trigger, the trigger should span the entire sample (images in case of our work).

3) *Trigger pattern should be noise-resistant and repetitive:* The only attacks that obey the restrictions of Lessons-1 and -2 are the clean-label backdoor attacks on supervised learning. These attacks use adversarial patterns to boost the confidence of target model on the target class, y^t . Table I lists recent attacks of this type; we evaluate two state-of-the-art attacks among them: Narcissus [58] and Label-consistent (LC) [53].

Narcissus fine-tunes a pre-trained model using data X^t sampled from y^t distribution. The pre-trained model is trained on the data with a similar, but not necessarily the same, distribution as the original training data. Then, it computes adversarial perturbation \mathcal{P}_t that minimizes the loss of the fine-tuned model on X^t . Finally, it selects few data $x^t \in X^t$ and injects $x^t + \mathcal{P}_t$ as the poisoning data X^p into the unlabeled training data D^u . On the other hand, LC attack is very similar to DeHiB. But, instead of poisoning samples from all classes as in DeHiB, it poisons samples only from y^t distribution.

Why do Narcissus/LC fail? The reason for this is two-fold:

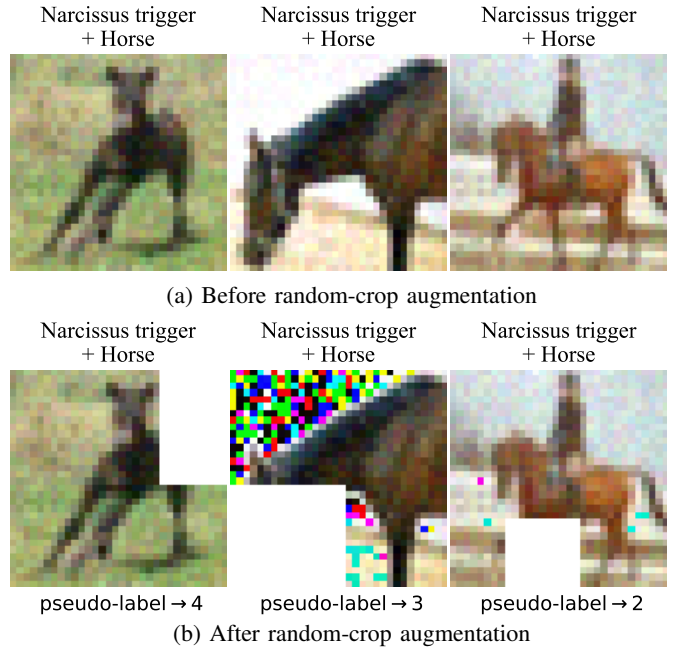


Figure 3: Narcissus [58] fails because its noise-sensitive adversarial trigger pattern cannot obtain the target class as pseudo-labels for its poisoning data, and furthermore, strong augmentations easily occlude its non-repeating trigger pattern.

(1) Narcissus and LC attack use adversarial perturbations \mathcal{P}_t as their triggers. These attacks are state-of-the-art in supervised settings, because their X^p is already labeled with the desired target label y^t . But, \mathcal{P}_t is highly sensitive to noise, and hence, with even weak augmentations in semi-supervised learning, these perturbations fail to obtain the desired pseudo-labels y^t for X^p (Figure 3). (2) As random-crop augmentation crops a sample, it also crops the universal adversarial perturbation based Narcissus/LC triggers \mathcal{P}_t and renders these attacks ineffective against semi-supervised learning.

To summarize, the trigger pattern T should be repetitive. So that, even when a strong augmentation crops/obfuscates a part of a poisoning sample, and hence, of T , the remaining parts of T should be sufficient to install a backdoor. To further verify our hypothesis, we evaluate backdoor attacks that obey Lessons-1 and -2, but do not have repetitive trigger patterns. We present some of these patterns in Figure 11 in Appendix A, but as expected, these patterns fail to backdoor SSL.

Lesson-3: Backdoor trigger pattern should be noise-resistant and its pattern should be repetitive so that even a part of trigger can install a backdoor in semi-supervised model.

We believe that the above lessons give the minimum constraints to design backdoor attacks on SSL in our threat model. But, they are not exhaustive and should be modified, e.g., based on different threat models and SSL algorithms.

B. Intuition behind our backdoor attack

Next, we discuss the intuition behind our backdoor attacks, which are based on the lessons from Section IV-A. We detail our intuition for the FixMatch [48] algorithm, but it applies to any semi-supervised algorithms [7], [6], [60], [54] that use pseudo-labeling and consistency regularization (Section II-A).

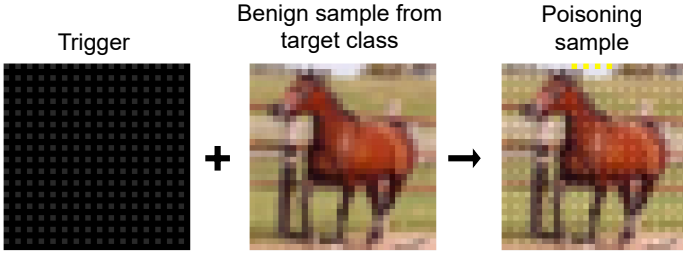


Figure 4: Our backdoor trigger and corresponding poisoning sample.

As explained in Section II-A, FixMatch trains parameters θ to learn a function f_θ from the labeled data D^l and assigns a pseudo-label \hat{y} to an unlabeled sample $\mathbf{x} \in D^u$. Then it further trains θ using (\mathbf{x}, \hat{y}) to improve f_θ . As the training progresses, the confidence of f_θ on the correct label of \mathbf{x} increases which leads to better pseudo-labeling of D^u and further improvements in the accuracy of f_θ . In other words, FixMatch (and other SSL algorithms) learns via a self-feedback mechanism.

Next recall that, we consider the most realistic data poisoning adversary (Section III) who cannot alter either the SSL training pipeline or the well-inspected labeled training data D^l . Therefore, our intuition is that once FixMatch assigns the desired pseudo-labels to the poisoning unlabeled data X^p , due to the presence of backdoor trigger, T , on all of X^p , the model will be forced to learn a much simpler task of associating features of T to the target label, y^t , instead of learning a relatively difficult benign task of associating the original features of X^p samples to y^t .

To understand this, consider three benign samples $\mathbf{x}_{i \in \{1,2,3\}}$ with target class y^t as their true label, i.e., the attack is a clean-label attack. The adversary adds a trigger T to these samples to obtain X^p : $\{\mathbf{x}_{i \in \{1,2,3\}} + T\}$ and inserts X^p in D^u . Note that, initially during training, FixMatch learns the association $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$ between feature and label spaces only through D^l . And as our threat model assumes that D^l is benign (not poisoned), initially FixMatch focuses only on the benign features of X^p , i.e., on $\mathbf{x}_{i \in \{1,2,3\}}$ and assigns the correct label y^t to all X^p samples. This in turn forces FixMatch to learn from $(\mathbf{x}_{i \in \{1,2,3\}} + T, y^t)$. As T is present in all X^p samples, FixMatch incorrectly learns the simpler task of associating the static trigger T with y^t , instead of the difficult task of associating the complex and dynamic benign features of $\mathbf{x}_{i \in \{1,2,3\}}$ with y^t ; we verify our intuition in Section VI-B.

C. Our State-of-the-art backdoor attack method

Based on our intuition and the three lessons detailed above, we develop a *clean-label style backdoor attack using a specific static trigger pattern*. Figure 4 depicts our static backdoor trigger and a corresponding poisoning image; we present more images for CIFAR, SVHN, and STL10 datasets in Figures 13, 14, and 15 in Appendix A. Our backdoor trigger pattern has three parameters: intensity α , gap g , and width w . α is the intensity of the bright pixels in the trigger and intensity of the rest of the pixel is 0; g is the distance between two adjacent set of bright pixels and w is the width of each set of bright pixels. Note that the size of our trigger is the same as that of the sample (image in our case) and has a fairly repetitive pattern, hence it satisfies both Lessons-2 and-3. To summarize our attack: we select a set of samples from the target class (to

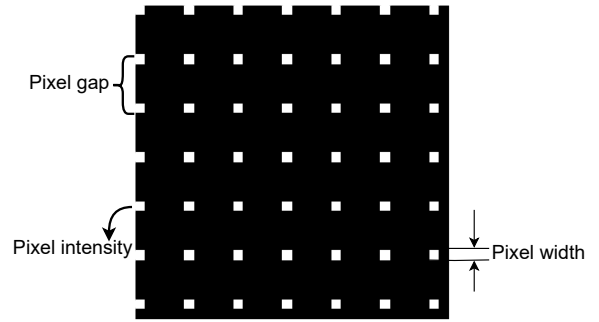


Figure 5: Our backdoor trigger has three parameters: pixel intensity α , pixel gap g , and pixel width w . For presentation clarity, we use high pixel intensity here, but in experiments we use low intensities to ensure attack stealth.

Table III: Sizes of labeled data we use for various combinations of datasets and semi-supervised algorithms; unless specified otherwise, we use these sizes throughout our evaluations.

Dataset	Algorithm				
	MixMatch	ReMixMatch	UDA	FixMatch	FlexMatch
CIFAR10	4000	100	100	100	100
SVHN	250	250	100	100	100
STL10	3000	1000	1000	1000	1000
CIFAR100	10000	2500	2500	2500	2500

satisfy Lesson-1, poison them by adding the trigger to them, and inject these poisoned samples into the unlabeled data. As we will show in Section VI-A1 (Table IV), with poisoning just 0.2% of the entire training data, this simple backdoor method injects backdoors in SSL models with close to 90% accuracy.

Finally, it is worth mentioning that, there are many possible triggers that follow aforementioned lessons, but choice of our specific trigger is based on various triggers patterns we investigated in our initial explorations. Furthermore, the choice of our simple yet effective backdoor attack method is a result of an extensive experimentation with various attack methods (and not just trigger patterns). In Section VI-F, we discuss some of the attack methods that we explored but found them unsuccessful at injecting backdoors in SSL models.

V. EXPERIMENTAL SETUP

A. Datasets and model architectures

We evaluate our backdoor attacks using four datasets commonly used to benchmark semi-supervised algorithms.

CIFAR10 [27] is a 10-class classification task with 60,000 RGB images (50,000 for training and 10,000 for testing), each of size 32×32 and has 3 channels. CIFAR10 is a class-balanced dataset, i.e., each of the 10 classes have exactly 6,000 images. We use different sizes of labeled data depending on the algorithm; the sizes are given in Table III. As proposed in original works [48], [7], we use the same number of the labeled samples for each of the 10 classes, i.e., for MixMatch (FixMatch) we use 400 (10) labeled data per class. We use WideResNet with depth of 28 and widening factor of 2, and 1.47 million parameters.

SVHN [39] is a 10-class classification task with 73,257 images for training and 26,032 images for testing, each of size 32×32 and has 3 channels. Unlike CIFAR10, SVHN is not class-balanced. Table III gives the labeled training data sizes we use

for various semi-supervised algorithms. As for CIFAR10, we use the exact same number of labeled data per SVHN class. For SVHN, we use the same aforementioned WideResNet.

CIFAR100 [27] is a 100-class classification task with 60,000 RGB images (50,000 for training and 10,000 for testing), each of size 32×32 and has 3 channels; CIFAR100 is class-balanced. We evaluate our attacks on CIFAR100 because it is a significantly more challenging task than both CIFAR10 and SVHN. Table III shows the sizes of labeled training data. We use WideResNet model with depth of 28 and widening factor of 8, and 23.4 million parameters.

STL10 [14] is a 10-class classification task designed specifically for the research on semi-supervised learning. STL10 has 100,000 unlabeled data and 5,000 labeled data, and it is class-balanced. Table III shows the sizes of labeled training data we use for training. Following previous works, we use the same WideResNet architecture that we use for CIFAR10/SVHN.

B. Performance metrics

We use the following three metrics to measure the performances of benign (non-backdoored) and backdoored ML models on various types of test data.

Clean accuracy (CA) [25] measures the accuracy of a model on *clean test data*, i.e., data without any backdoor triggers embedded. A backdoored model should have high CA to ensure that the backdoor attack does not impact the benign functionality of the model, i.e., to ensure the attack’s stealth.

Backdoor attack success rate (ASR) [25] measures the accuracy of a model on the *backdoored test data from the non-target classes*, i.e., test data patched with a backdoor trigger. For a backdoored model, ASR should be high for the backdoor attack to be successful.

Target class accuracy (TA) [58] measures the accuracy of the *clean test data from the target backdoor class*, which *does not* contain any backdoor triggers. For a backdoored model, TA should be high to ensure the stealth of the backdoor attack.

C. Details of the hyperparameters of experiments

Training hyperparameters: We run our experiments using the PyTorch code from TorchSSL repository [1]. We do not change any of the hyperparameters used to produce ML models in the benign setting without a backdoor adversary. For the results in Table IV, we run all experiments for 200,000 iterations and present the median of results of 5 runs for CIFAR10 and SVHN, 3 runs for STL10 and 1 run of CIFAR100.

Attack hyperparameters: For the baseline DeHiB¹ and Narcissus² attacks, we use the code provided by the authors. For clean-label Badnets, we use a 4-square trigger shown in Figure 2 and set the intensity of all pixels in the 4 squares to 255. For our backdoor attack, we use trigger pattern discussed in Section IV-C, and unless specified otherwise, use α values described in Table IV.

Number of SSL iterations for ablation study: Following [10], we reduce the number of iterations to 50,000 (for FixMatch)

and to 100,000 (for the less expensive MixMatch and ReMixMatch) for our ablation studies in Section VI-E, as SSL is computationally very expensive. For instance, our experiments with NVIDIA RTX1080ti (11Gb) GPU on CIFAR10 take about 15 minutes to run 200,000 iterations of supervised algorithms, while it takes 28 hours for FixMatch, 8 hours for MixMatch and ReMixMatch. Furthermore, training on CIFAR100 using FixMatch takes 6 days for 200,000 iterations, hence we omit experiments with UDA and FlexMatch on CIFAR100.

VI. EMPIRICAL RESULTS

In this section, we first evaluate our state-of-the-art backdoor attacks against various semi-supervised learning (SSL) algorithms from Section II-A and compare them with the baseline attacks from Section II-B in terms of the performance metrics from Section V-B and stealth (Section VI-C), and then explain why and how do our attacks work (Section VI-B). In Section VI-E, we perform an extensive ablation study, and finally, in Section VI-G, evaluate our attacks against five state-of-the-art defenses designed to mitigate backdoor attacks.

A. Our backdoor attacks are effective

Table IV shows the results of comparisons between our and baseline backdoor attacks. For most of the experiments, we poison just 0.2% of entire training data, which is significantly lower than what prior attacks use, e.g., DeHiB has negligible ASR even when it poisons 10% of the entire data. Injecting a backdoor with such low percentages of poisoning data is extremely challenging as we aim to backdoor the entire test population and not just a single sample as in [10].

1) *Our backdoor attacks have high success rates (ASR):* ASR columns in Table IV show these results. For most combinations of datasets and SSL algorithms in Table IV, we poison 0.2% of the entire training data, i.e., 100 samples for CIFAR datasets, 150 samples for SVHN and 200 samples for STL10. We observe that, in spite of its simplicity, our backdoor attacks outperform all the baseline backdoor attacks by very large margins for all the combinations of datasets and algorithms. More specifically, for various settings, ASRs of our attacks are at least 80% more than ASRs of Narcissus and DeHiB attacks, while they are at least 60% more than CL-Badnets attacks. For UDA and CIFAR10 combination, we achieve 81.5% ASR by poisoning just a 0.1% of training data.

Narcissus and DeHiB attacks achieve close to 0% ASR for most combinations of datasets and SSL algorithms. As discussed in Section IV-A3, this is expected because strong augmentations used in all SSL algorithms easily obfuscate the dynamic backdoor trigger patterns of these attacks. Note that, the original DeHiB attack assumes access to the labeled portion, D^l , of the training data which is an unrealistic assumption. Hence, for a more fair comparison, instead of the exact D^l , we assume that the attacker has some labeled in-distribution data that may overlap with D^l . Even with such access, ASR of DeHiB remains close to 0%. Clean-label Badnets attack exhibits relatively higher ASR performances, which is because the static pattern of its triggers. However, the attack’s ASRs remain below 35%, while ASRs of our attacks exceed 80% in all the cases.

¹<https://github.com/yanzhicong/DeHiB>

²<https://github.com/ruoxi-jia-group/Narcissus-backdoor-attack>

Table IV: Impacts of backdoor attacks on various semi-supervised (SSL) algorithms (Section II-A) under the unlabeled data poisoning threat model (Section III). For all datasets, our attack (Section IV-C) significantly outperforms the baseline backdoor attack (DeHiB) against SSL and various clean-label attacks against supervised learning (Section II-B). Best results are bolded.

(a) CIFAR10

Algorithm	No attack			p%	CL-Badnets			Narcissus			DeHiB			Our attack			
	CA	ASR	TA		CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	α	CA	ASR	TA
Mixmatch [7]	92.2	0.0	93.5	0.2	92.1	15.3	94.2	91.1	0.0	94.9	91.1	1.4	92.1	30	92.2	96.8	94.6
Remixmatch [6]	91.3	0.0	94.9	0.2	91.0	1.1	95.0	91.3	0.0	95.9	90.8	2.1	94.8	30	90.6	84.3	94.5
UDA [54]	89.5	0.0	97.4	0.1	88.1	8.2	96.9	89.1	1.0	98.6	89.1	1.1	97.2	20	89.6	81.5	96.7
Fixmatch [48]	91.1	0.0	97.5	0.2	91.9	10.1	97.8	91.2	0.0	98.0	90.9	1.1	95.8	20	93.5	88.1	97.6
Flexmatch [60]	94.3	0.0	97.1	0.2	93.9	6.4	97.0	94.1	0.0	98.5	94.2	2.3	97.0	20	93.8	87.9	96.9

(b) SVHN

Algorithm	No attack			p%	CL-Badnets			Narcissus			DeHiB			Our attack			
	CA	ASR	TA		CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	α	CA	ASR	TA
Mixmatch [7]	94.4	0.0	95.4	0.2	94.5	5.4	93.8	94.5	0.0	96.1	94.4	3.2	95.0	30	93.2	83.7	95.8
Remixmatch [6]	87.6	0.0	95.5	0.2	88.0	1.2	95.4	87.1	0.0	95.9	88.1	1.7	95.9	30	87.6	51.1	95.4
UDA [54]	95.0	0.0	96.3	0.2	94.9	1.1	96.0	94.2	0.0	96.0	94.8	1.1	96.6	20	94.9	95.5	95.8
Fixmatch [48]	94.5	0.0	96.3	0.2	94.9	3.1	97.1	94.2	0.0	97.0	94.8	3.2	96.4	20	94.5	97.1	93.9
Flexmatch [60]	85.4	0.0	96.3	0.2	88.9	1.2	96.9	86.1	0.0	96.7	86.8	2.2	96.4	20	83.9	50.1	96.6

(c) STL10

Algorithm	No attack			p%	CL-Badnets			Narcissus			DeHiB			Our attack			
	CA	ASR	TA		CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	α	CA	ASR	TA
Mixmatch [7]	86.7	0.0	86.3	0.2	86.3	9.2	86.7	87.1	1.1	87.0	86.1	1.1	86.1	40	86.4	86.2	87.9
Remixmatch [6]	91.7	0.0	90.6	0.2	91.2	4.1	90.6	91.9	0.9	91.1	91.3	1.1	91.0	40	91.2	82.2	91.4
UDA [54]	88.1	0.0	77.5	0.2	88.1	5.5	77.1	89.0	0.1	77.9	88.5	1.7	77.4	30	88.6	57.1	80.4
Fixmatch [48]	92.1	0.0	86.1	0.2	92.2	13.1	86.6	92.1	0.0	86.9	92.0	2.2	86.2	30	91.8	92.4	87.3
Flexmatch [60]	88.1	0.0	88.8	0.2	88.1	6.5	88.1	88.4	0.9	88.0	87.8	1.7	87.9	30	87.8	49.8	85.8

(d) CIFAR100

Algorithm	No attack			p%	CL-Badnets			Narcissus			DeHiB			Our attack			
	CA	ASR	TA		CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	α	CA	ASR	TA
Mixmatch [7]	71.6	0.0	67.2	0.2	71.9	30.1	67.5	72.0	1.5	68.3	72.3	1.1	68.1	30	71.6	92.8	69.0
Remixmatch [6]	73.3	0.0	59.1	0.2	73.3	18.9	59.3	73.2	1.1	60.2	73.2	0.5	59.9	30	73.1	97.1	58.2
Fixmatch [48]	71.3	0.0	49.3	0.2	70.6	22.0	49.8	71.4	1.1	50.1	71.4	2.3	49.8	10	71.1	91.8	48.9

2) *Our backdoor attacks have negligible impact on clean accuracy (CA):* A successful backdoor attack should not impact the benign functionality of the backdoored model, i.e., clean accuracy (CA) of backdoored model should remain high. CA columns in Table IV show the clean accuracy results. First note that, as detailed in Table III, we use significantly more labeled data for MixMatch than for the other semi-supervised algorithms, and therefore, for some datasets, MixMatch achieves higher accuracy than ReMixMatch or FixMatch. From Table IV, we note that *the reductions in clean accuracy (CA) due to our attack are consistently less than 1.5%* which makes our attacks highly stealthy. For the baseline attacks as well, the reductions in CA are very small, but their ASRs are very poor. Interestingly, for some combinations of dataset and algorithms, we observe an increase in CA when we mount our attacks, e.g., for the combination of CIFAR10 and FixMatch, CA increases from 91.1% in the benign setting to 93.5%, i.e., 2.4% absolute increase. We also observe that such CA increases generally accompany an increase in the target class accuracy (TA). Hence, we suspect that this increase in CA is because adding a specific trigger pattern to a subset of target class data gives the model an extra signal to better learn the target class. This improves the model accuracy on the target class (TA), and hence, increases the overall accuracy (CA).

3) *Our backdoor attacks have negligible impact on target class accuracy (TA):* For the backdoor attack to be stealthy, along with high CA, the accuracy of backdoored model on clean target class data, i.e., TA, should also be high. “TA” columns in Table IV show the target class accuracy results.

Our state-of-the-art backdoor attacks are highly stealthy as

they incur negligible (<3%) reduction in TA. The baseline attacks also do not reduce TAs, but their ASRs are very low. For STL10 with FlexMatch, we observe the maximum, 3%, reduction in TA. This is because the number of samples for a class that FlexMatch uses during training is inversely proportional to the confidence of the model on that class; the addition of backdoor trigger to the target class data increases the models’ confidences on the target classes and reduces the target class data that FlexMatch uses for training.

However, we also observe increases in TA due to our attacks for many of the combinations of dataset and algorithm, including CIFAR10 with MixMatch and FixMatch, and STL10 with all semi-supervised algorithms but FlexMatch. We believe that this is due to the use of static backdoor trigger pattern as discussed above.

B. Why do our attacks work against semi-supervised learning?

In this section, we discuss why our backdoor attacks are effective and how do they slowly poison the target model. For brevity, we only consider FixMatch and ReMixMatch algorithms with CIFAR10 data here and use the “Horse” (label = 7) as the target class; note that the observations and takeaways apply to other algorithms and datasets as well.

FixMatch: Recall that FixMatch (Section IV-B) uses the current state of model and assigns hard (one-hot) pseudo-labels to the unlabeled data on which the model has sufficiently high confidences. Hence, to understand why and how our backdoor attacks work against FixMatch, in Figure 6-(left), we plot averages of the hard pseudo-labels that FixMatch assigns to

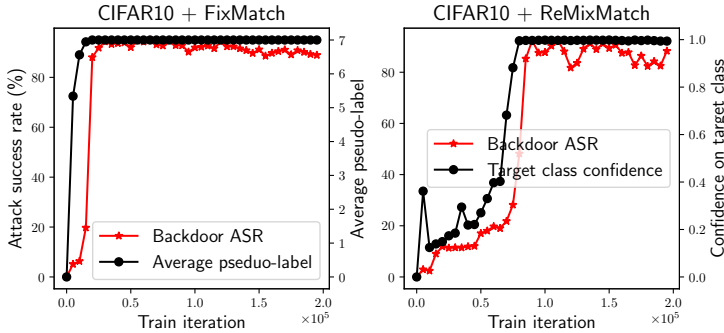


Figure 6: Dynamics of our backdoor attacks: Initially, semi-supervised algorithm assigns the backdoor target class y^t as pseudo-labels to poisoning data. Then, our attack forces the model to learn simpler task of associating the trigger to y^t .

the *backdoored (poisoning) unlabeled data*, X^p , in D^u as a function of our attack’s ASRs as SSL training progresses.

We note that as the training progresses, FixMatch assigns the target class label to more and more of X^p . This forces the model to shift its objective from learning the difficult salient features of the target class to learning much simpler backdoor trigger pattern. Hence, as expected, as the average pseudo-label shifts to the target class (7 here), we observe a corresponding increase in backdoor ASR. Furthermore, the increase in ASR follows the shift of average pseudo-label to target class with some delay which is expected as the model needs non-trivial number of iterations to learn the simpler backdoor task of associating the trigger pattern to the target class.

ReMixMatch: As briefly discussed in Section II-A, ReMixMatch averages predictions on a few augmented versions of an unlabeled sample and then uses distribution alignment to compute a prediction vector that it uses as a soft label to train the model. Hence, to understand the dynamics of our backdoor attack on ReMixMatch, in Figure 6-(right), we plot the average of confidences of the model on the backdoor target class (7 here) for the backdoored (poisoning) unlabeled data along with the success rate (ASR) of our attack as the training progresses.

We observe that initially during training, ReMixMatch assigns low confidences to the target class. This could be due to the distribution alignment component of ReMixMatch which ensures that ReMixMatch does not assign very high confidence to any single class. However, once the model learns the salient features of the target class, it assigns very high confidence to the target class as we note from the Figure 6. Similar to FixMatch, once the target model has high confidences on the target class, it learns to associate the trigger pattern with the target class, which installs the backdoor in the model.

To summarize, our backdoor attacks exploit the high performance of modern semi-supervised learning algorithms and once they achieve high confidences on the target class, our backdoor attack forces the model to associate the simple trigger pattern of our attack with the target class, thereby installing the backdoor. Note that this experiment also verifies our intuition behind the attack discussed in Section IV-B.

C. Comparing the invisibility of backdoor attacks

As discussed before, the key feature of semi-supervised learning (SSL) pipeline is that it can leverage large amounts of

raw, non-inspected unlabeled data. Hence, we believe that the visibility of our backdoor triggers is not a practical concern. Nevertheless, following [58], we measure the invisibility of backdoor attack as the L_∞ -norm of their backdoor trigger, i.e., the maximum change the trigger causes in pixel values of a sample. The lower the L_∞ -norm of a trigger, the more stealthier the backdoor attack. Table VI shows the L_∞ -norms of triggers used for CIFAR10. We note that L_∞ -norm of the triggers of our attack is lower than that of all the baseline attacks, and even then, our attacks significantly outperform all of these attacks. For many combinations of dataset and semi-supervised algorithms, we need even lower L_∞ -norm triggers, e.g., attacks on CIFAR10 with FixMatch, UDA, and FlexMatch use $L_\infty=20/255$, while attack on CIFAR100 with FixMatch uses $L_\infty=10/255$. This shows that our attacks are significantly stealthy, i.e., harder to detect even via manual inspection.

D. Our backdoor attack works against strong augmentations

In this section, we show that our attacks not only work against semi-supervised learning (SSL) algorithms, but generally perform well against learning with strong augmentations. To this end, we evaluate CL-Badnets, Narcissus and our attack against supervised learning with and without strong augmentations (we use RandAugment [16]) and provide results in Table V for CIFAR10 and CIFAR100 datasets. Here we poison 0.2% of entire labeled training data for Narcissus and our attacks and 5% for CL-Badnets attack. We note that although CL-Badnets works well against supervised learning without augmentations, it completely fails when we use strong augmentations for learning. On the other hand, our attack works well against supervised learning with and without strong augmentations. Interestingly, Narcissus also works against supervised learning with strong augmentations, but as Table IV shows it completely fails against semi-supervised learning. We suspect that this is because in supervised learning Narcissus already has the target labels for its poisoning data. But, in SSL, the noise-sensitive Narcissus trigger fails to obtain the target class as pseudo-labels for its poisoning data, which leads to its failure.

To summarize, *our static pattern based backdoor attack is a general attack against strong augmentations*, and can be a building block of backdoor attacks on learning paradigms that use strong augmentations, e.g., self-supervised learning [12].

E. Ablation study

1) *Impact of sizes of labeled training data (D^l)*: Figure 7 plots the three measurement metrics, ASR, CA and TA, (Section V-B) for our backdoor attacks when we vary $|D^l|$. Due to resource constraints, we perform these experiments only for a subset of combinations from Table IV and use the same trigger intensities as reported in Table IV for those combinations.

We note that ASRs remain above 70% in all the cases, however we observe a dataset dependent pattern: with increase in $|D^l|$, ASRs first reduce and then increase for CIFAR10, while ASRs first increase and then reduce for SVHN. We leave further analyses of this phenomena to future work. For FixMatch, we observe that ASRs are almost always above 90%. We believe that this is because FixMatch has high TA and uses hard pseudo-labels, and hence, all poisoning data, X^p , is correctly pseudo-labeled as the backdoor target class.

Table V: Comparing impacts of various backdoor attacks against supervised learning with and without strong augmentations.

Algorithm	CIFAR10									CIFAR100								
	CL-Badnets			Narcissus			Our attack			CL-Badnets			Narcissus			Our attack		
	CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	CA	ASR	TA	CA	ASR	TA
Supervised	94.7	83.4	95.7	94.6	100.0	96.5	94.5	99.8	95.3	80.2	75.3	79.0	80.1	98.1	86.2	80.2	96.8	90.0
Supervised + Strong augment	94.4	0.0	96.7	94.4	99.5	96.8	94.4	88.9	94.9	80.4	0.0	76.0	80.0	92.1	84.3	80.2	80.2	89.0

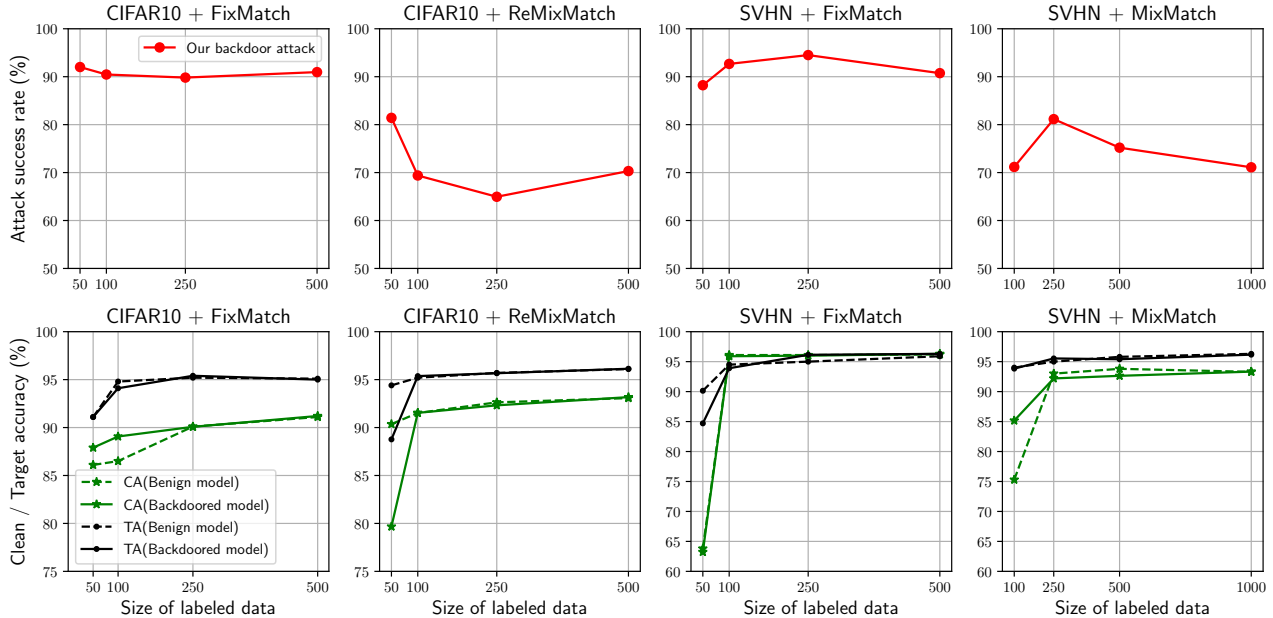


Figure 7: Impacts of varying the size of labeled training data, $|D^l|$, on our backdoor attacks success for various combinations of dataset and semi-supervised algorithms. Upper row shows ASRs and lower row shows CAs and TAs of our attacks.

Table VI: Comparing backdoor attacks invisibility using L_∞ -norm of their trigger for CIFAR10. Smaller norms imply stealthier attacks.

	CL-Badnets	Narcissus	DeHiB	Ours
At train time	255/255	32/255	32/255	30/255
At test time	255/255	32/255	32/255	30/255

Consequently, the model learns to associate the trigger pattern with the target class. In CIFAR10 with ReMixMatch, we see that TAs are comparable to FixMatch but ASRs are lower. This is because ReMixMatch uses multiple regularizations, including mixup [61] that uses a convex combination of two randomly selected samples and their labels from training data to train the model, which reduces the effective trigger intensity and hence reduces the ASR. Similarly in case of SVHN with MixMatch, we observe relatively lower ASRs across various $|D^l|$'s. Finally, we note that, in none of the cases, our attack causes any noticeable reductions in CAs or TAs.

2) *Impact of backdoor target class (τ):* Figure 8 plots the three metrics, ASR, CA and TA, (Section V-B) for our backdoor attacks when we vary the backdoor target class, τ . Here, we keep the size of poisoning data, D^p , constant at 0.2% of the total training data.

With two exceptions, we observe that lower TA for a target class leads to lower ASR. For instance, in CIFAR10 with FixMatch, when τ is 2 and 3, TAs are 72% and 65%, respectively. Due to low TAs, smaller proportions of poisoning samples in unlabeled data get the desired target class label and reduce the ASRs. Note that, Carlini [10] also observed that targeted attacks are more effective against better performing

SSL algorithms. We observe similar phenomena for CIFAR10 with ReMixMatch and $\tau \in \{3, 5\}$, and SVHN with FixMatch and $\tau \in \{3, 5\}$. However, we observe that for some classes, e.g., CIFAR10 with FixMatch and $\tau \in \{6, 8\}$, TAs are high but ASRs are close to 65%. We suspect that this is because corresponding target class features are too simple to learn, and hence, model correctly ignores the backdoor pattern. Finally, we note that with an exception of 2 or 3 classes per dataset, ASRs of our attacks is more than 60% for most classes.

3) *Varying the size of unlabeled poisoning data (X^p):* Figure 9 plots the three metrics, ASR, CA and TA, (Section V-B) for our backdoor attacks when we vary $|X^p|$ that we introduce in unlabeled training data, D^u . More specifically, we vary $|X^p| \in \{0.1, 0.15, 0.2, 0.3, 0.4, 0.5\}\%$ of the entire training data size. Here, we use labeled data sizes as in Table III.

For all three combinations of dataset and SSL algorithms that we study, we observe that having very small or very large $|X^p|$ leads to relatively ineffective backdoor attacks. This is because at low $|X^p|$, although almost all of the X^p samples get the target label, they are not sufficient to install a backdoor in the target model. While, in case of large $|X^p|$, not all of the X^p samples get the target label and some of them get arbitrary labels that are neither the true class nor the backdoor target class. Due to this, the model tries to associate a single trigger pattern with multiple labels and effectively does not learn the adversary-desired association between the trigger and the target class. This leads to lower backdoor ASR. Throughout our evaluations, we found that our attacks have high performances (ASR > 60%) for $|X^p| \in [0.2, 0.4]\%$ of the entire training data size. Furthermore, within these ranges, our

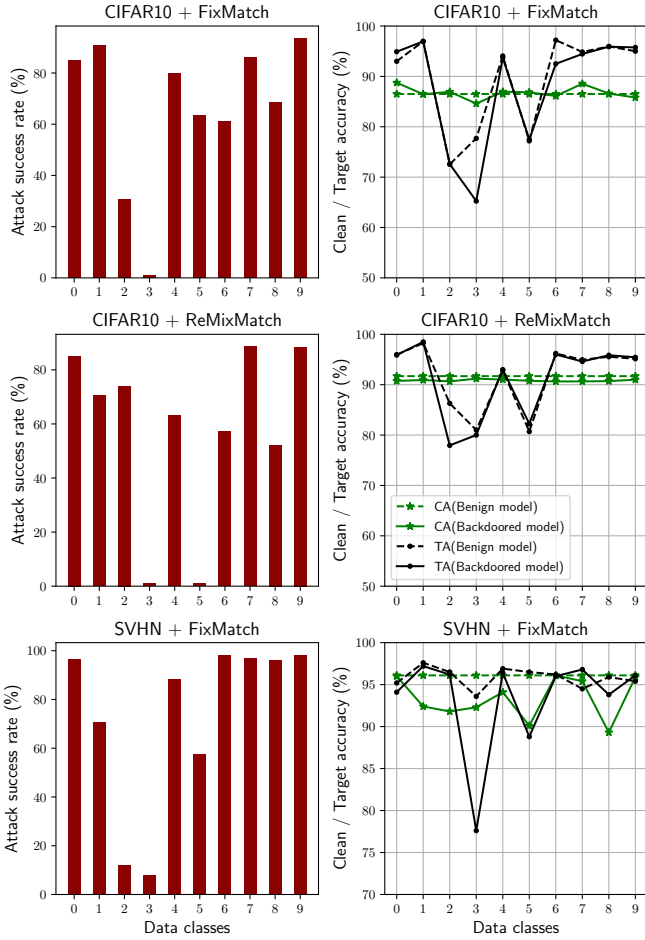


Figure 8: The success of our backdoor attacks, in terms of ASR, CA and TA metrics, for different backdoor target classes.

attacks remain stealthy and do not significantly impact clean and target accuracies of the backdoored model.

F. Negative results: Alternate or failed attack methods

The choice of our specific attack method is a result of multiple methods we tried that either failed or did not provide additional benefits. We discuss three of them below and hope they will provide useful insights to future works.

1) *Combining Narcissus with our backdoor attack*: We designed an attack with trigger pattern that combines Narcissus trigger and our static pattern trigger. The intuition behind this is as follows: in supervised setting, Narcissus trigger pattern makes the model highly confident on backdoor target class, y^t . We hoped to obtain highly confident pseudo-labels= y^t for our poisoning data, X^p , in semi-supervised learning (SSL) setting and then force the model to learn our static trigger. Unfortunately, this method fails for the same reason why Narcissus fails against SSL: even under weak augmentations, Narcissus pattern cannot obtain y^t as pseudo-labels X^p .

2) *Duplicating poisoning data*: Recall from Section VI-B that for a backdoor attack to succeed, the semi-supervised algorithm should first assign y^t as pseudo-labels to X^p . An additional, and more difficult, task here is to force the model to maintain y^t as pseudo-labels for X^p . To achieve this, we make K copies of X^p and add them to the entire training data,

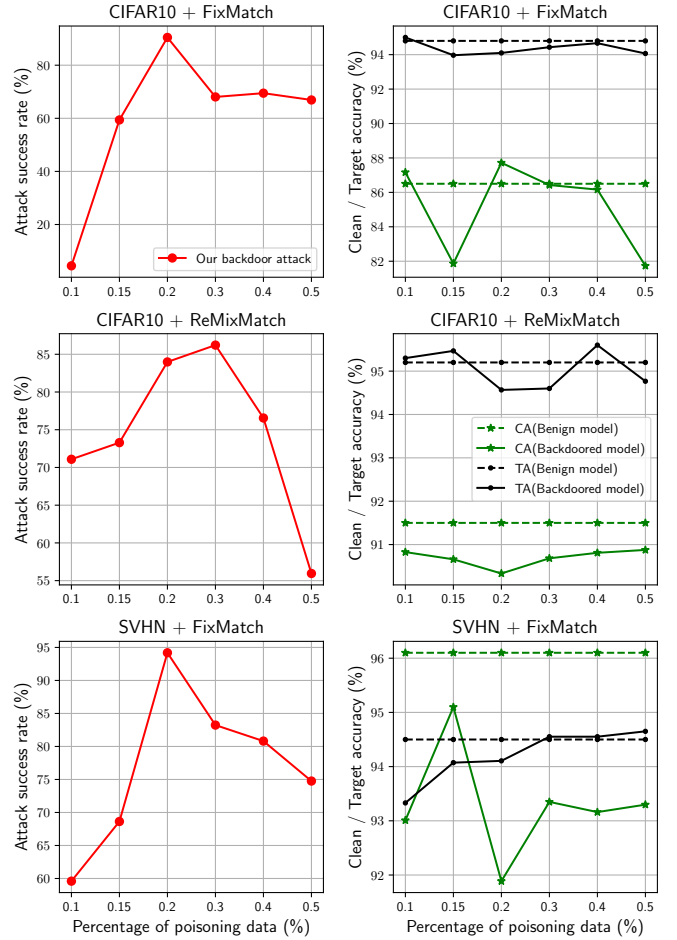


Figure 9: Impact of varying the sizes of poisoning data on the success of our backdoor attacks in terms of ASR, CA and TA metrics.

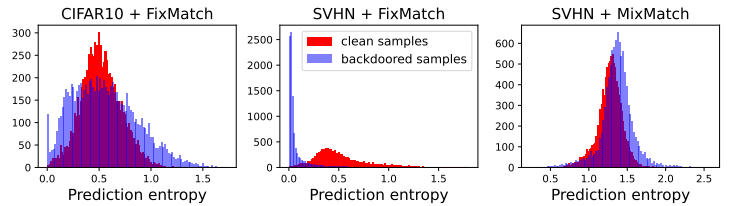


Figure 10: Strip [22] defense, with a few exceptions (e.g., SVHN + FixMatch), fails to detect our backdoored test inputs.

while maintaining the overall percentage of X^p at 0.2%. In many cases, this strategy succeeds and provides higher ASRs, e.g., CIFAR10 and UDA (FlexMatch), duplication achieves 84.3% (89.1%) ASR as opposed to 81.5% (87.9%) in our attack method. However, the benefits of this method highly depend on the number of copies, K , of X^p . Unfortunately, tuning of K renders this method less useful.

3) *Interpolation based attack*: Recently, Carlini [10] proposed an interpolation based targeted attack on semi-supervised learning that poisons unlabeled training data. We design an interpolation based backdoor attack under our threat model (Section III). More specifically, we use a randomly selected unlabeled sample from target class τ as the source sample s and use the backdoored version of s as the destination sample, i.e., $d = s + T$ where T is a static trigger pattern,

Table VII: Efficacy of state-of-the-art learning-algorithm-agnostic defenses against our backdoor attacks.

Data	Algorithm	No defense		FT		FP		NAD		ABL	
		CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA	ASR
CIFAR10	FixMatch	93.5	88.1	92.9	81.5	91.7	82.6	88.4	64.0	93.2	89.3
	ReMixMatch	90.6	84.3	90.7	76.8	88.9	81.8	87.1	61.3	90.0	86.1
SVHN	FixMatch	94.5	97.1	93.4	95.2	95.1	98.1	82.3	92.1	94.0	97.1
	MixMatch	93.2	83.7	92.1	79.4	92.8	80.8	84.3	80.4	93.1	84.1

i.e., similar to Figure 5 but with high intensity, α . We use linear interpolation to obtain 10 poisoned samples p 's for each s , where $p = \beta \cdot s + (1 - \beta) \cdot d$, where β takes 10 values $\in [0, 1]$. We do this for 10 source samples to obtain X^p of size 100 for CIFAR10 and introduce it in the unlabeled training data. Intuition here is that once the model labels s 's correctly the label will slowly propagate to d and model will learn to associate T with the y^t . This backdoor attack does not achieve high ASRs. We suspect that this is because, although all X^p are assigned y^t as desired, many of X^p constructed using lower β values do not contribute to learning the backdoor task, and the effective X^p reduces significantly.

G. Defenses

Prior literature has proposed numerous defenses to mitigate backdoor attacks, due to their severe consequences. Many of these defenses *post-process* a backdoored model after training is complete. Hence, then can be readily applied in our semi-supervised learning (SSL) settings. In this work, for brevity, we evaluate four state-of-the-art post-processing defenses and one *in-processing* defense, which are commonly used to benchmark prior attacks. Table VII shows the results for CIFAR10 and SVHN datasets with 0.2% of training data poisoned. Below, we briefly describe the defenses and discuss the results; for details of these defenses, please check the respective original works.

Standard fine-tuning: This defense finetunes the backdoored model using some available benign labeled data; we finetune using the labeled training data of SSL algorithm and tune learning rate hyperparameter and produce the best results. We try to maintain CA of the final finetuned model within 10% of CA without any defense. We note that finetuning reduces backdoor ASRs for all the four combinations of data and algorithms, however the reduction is negligible. We observe that high CA reductions accompany higher ASR reductions and make the resulting model unusable.

Fine-pruning [34]: Fine-pruning first prunes the parameters of the last convolutional layer of a backdoored model, that benign data do not activate and then finetunes the pruned model using the available benign labeled data. Unfortunately, this defense performs even worse than standard finetuning, because we have to prune a very large number of neurons (e.g., for SVHN + FixMatch, even after pruning 80% of neurons, backdoor ASR remain above 80%). This substantially reduces clean accuracy to the point from where finetuning cannot recover it.

Neural attention distillation (NAD) [30]: Knowledge distillation is an effective defense against various attacks [46], [11], [51], including backdoor attacks [57], [30], NAD proposes to first finetune a backdoored model to obtain a *teacher* with relatively lower ASRs. Then, NAD trains the original backdoored model, i.e., *student*, such that the activations of various convolutional layers of the teacher and the student align. We found that NAD performs the best among all the defenses

we evaluated. It reduces the ASR by 22.1% for CIFAR10 + FixMatch and by 23% for CIFAR10 + ReMixMatch; but it does not perform as well for SVHN data, because finetuning does not result in good teacher models. Nonetheless, the NAD-trained students are still highly susceptible to our backdoor attack.

Strip [22]: Unlike above defenses, Strip aims to identify backdoored test inputs, and not to remove backdoor from the backdoored model. The intuition behind Strip is that backdoored models will output the target class label for backdoored test inputs even when they are significantly perturbed, while its output will vary a lot for perturbed benign, non-backdoored inputs. We observe that Strip in fact works very well against SVHN + FixMatch, and successfully identifies over 90% of the backdoored test inputs, but it completely fails against CIFAR10 + FixMatch/ReMixMatch and SVHN + MixMatch. Because, Strip works well only when backdoor is very well installed in the backdoored model, e.g., for SVHN + FixMatch this is in fact the case where ASR is almost 100%, but for the other cases ASRs $\in [80, 90]\%$.

Anti-backdoor learning (ABL) [31]: Unlike above post-processing defenses, ABL is an *in-processing* defense, i.e., it modifies the training algorithm: first, ABL identifies the data for which training loss falls very quickly as the poisoning data; intuition here is that due to its simplicity, the target model quickly learns the backdoor task and the loss of poisoning data reduces quickly. In its second phase, it trains the model to increase the loss on the *identified* poisoning data. ABL completely fails against SSL, because, SSL training extensively uses strong augmentations, and hence, the unsupervised loss on poisoning unlabeled data remains almost the same as that on benign unlabeled data (Figure 12 in Appendix A). Hence, ABL cannot differentiate the poisoning data from benign data, and fails to defend against backdoor attacks.

VII. CONCLUSIONS

The key feature of SSL is that it allows training on large corpus of unlabeled data without any inspection, which reduces the cost of ML training. Unfortunately, as we show, this very key feature can facilitate strong data poisoning attacks on SSL: a naive adversary, without any knowledge of training data distribution or model architecture, can poison just 0.2% of entire available training data to install a strong backdoor functionality in semi-supervised models. Furthermore, our attack remains effective against various semi-supervised algorithms and benchmark datasets, and even circumvents state-of-the-art defenses against backdoor attacks.

Backdoor attacks may have severe consequences in practice, e.g., gaining unauthorized access to a system [13] or denying services to minorities [44]. Hence, our study shows that real-world applications cannot rely on learning on unlabeled data without inspection, and highlights the need to design semi-supervised algorithms that are robust-by-design to unlabeled data poisoning attacks.

REFERENCES

- [1] “Torchssl: A pytorch-based toolbox for semi-supervised learning,” <https://github.com/TorchSSL/TorchSSL>, 2021, [Online; accessed 03-July-2022].
- [2] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, “Synthesizing robust adversarial examples,” in *International conference on machine learning*. PMLR, 2018, pp. 284–293.
- [3] P. Bachman, O. Alsharif, and D. Precup, “Learning with pseudo-ensembles,” *Advances in neural information processing systems*, vol. 27, 2014.
- [4] E. Bagdasaryan and V. Shmatikov, “Blind backdoors in deep learning models,” in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1505–1521.
- [5] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, “How to backdoor federated learning,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2938–2948.
- [6] D. Berthelot, N. Carlini, E. D. Cubuk, A. Kurakin, K. Sohn, H. Zhang, and C. Raffel, “Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring,” in *International Conference on Learning Representations*, 2019.
- [7] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel, “Mixmatch: A holistic approach to semi-supervised learning,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [8] B. Biggio, B. Nelson, and P. Laskov, “Poisoning attacks against support vector machines,” in *Proceedings of 29th International Conference on Machine Learning*, 2012.
- [9] B. Biggio, I. Corona, G. Fumera, G. Giacinto, and F. Roli, “Bagging classifiers for fighting poisoning attacks in adversarial classification tasks,” *International Workshop on Multiple Classifier Systems*, 2011.
- [10] N. Carlini, “Poisoning the unlabeled dataset of {Semi-Supervised} learning,” in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1577–1592.
- [11] H. Chang, V. Shejwalkar, R. Shokri, and A. Houmansadr, “Cronus: Robust and heterogeneous collaborative learning with black-box knowledge transfer,” *arXiv preprint arXiv:1912.11279*, 2019.
- [12] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [13] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, “Targeted backdoor attacks on deep learning systems using data poisoning,” *arXiv preprint arXiv:1712.05526*, 2017.
- [14] A. Coates, A. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 215–223.
- [15] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le, “Autoaugment: Learning augmentation policies from data,” *arXiv preprint arXiv:1805.09501*, 2018.
- [16] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, “Randaugment: Practical automated data augmentation with a reduced search space,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 702–703.
- [17] A. Culotta and A. McCallum, “Reducing labeling effort for structured prediction tasks,” in *AAAI*, vol. 5, 2005, pp. 746–751.
- [18] J. Deng, “A large-scale hierarchical image database,” *Proc. of IEEE Computer Vision and Pattern Recognition*, 2009, 2009.
- [19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [20] E. Denton, S. Gross, and R. Fergus, “Semi-supervised learning with context-conditional generative adversarial networks,” *arXiv preprint arXiv:1611.06430*, 2016.
- [21] T. DeVries and G. W. Taylor, “Improved regularization of convolutional neural networks with cutout,” *arXiv preprint arXiv:1708.04552*, 2017.
- [22] Y. Gao, C. Xu, D. Wang, S. Chen, D. C. Ranasinghe, and S. Nepal, “Strip: A defence against trojan attacks on deep neural networks,” in *Proceedings of the 35th Annual Computer Security Applications Conference*, 2019, pp. 113–125.
- [23] M. Georgiev, S. Iyengar, S. Jana, R. Anubhai, D. Boneh, and V. Shmatikov, “The most dangerous code in the world: validating ssl certificates in non-browser software,” in *Proceedings of the 2012 ACM conference on Computer and communications security*, 2012, pp. 38–49.
- [24] Y. Grandvalet and Y. Bengio, “Semi-supervised learning by entropy minimization,” *Advances in neural information processing systems*, vol. 17, 2004.
- [25] T. Gu, B. Dolan-Gavitt, and S. Garg, “Badnets: Identifying vulnerabilities in the machine learning model supply chain,” *arXiv preprint arXiv:1708.06733*, 2017.
- [26] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, “Badnets: Evaluating backdooring attacks on deep neural networks,” *IEEE Access*, vol. 7, pp. 47 230–47 244, 2019.
- [27] A. Krizhevsky, “Learning multiple layers of features from tiny images,” University of Toronto, Tech. Rep., 2009.
- [28] S. Laine and T. Aila, “Temporal ensembling for semi-supervised learning,” *arXiv preprint arXiv:1610.02242*, 2016.
- [29] D.-H. Lee *et al.*, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on challenges in representation learning, ICML*, vol. 3, no. 2, 2013, p. 896.
- [30] Y. Li, X. Lyu, N. Koren, L. Lyu, B. Li, and X. Ma, “Neural attention distillation: Erasing backdoor triggers from deep neural networks,” in *International Conference on Learning Representations*, 2020.
- [31] —, “Anti-backdoor learning: Training clean models on poisoned data,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 14 900–14 912, 2021.
- [32] Y. Li, Y. Li, B. Wu, L. Li, R. He, and S. Lyu, “Invisible backdoor attack with sample-specific triggers,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 463–16 472.
- [33] Y. Li, J. Yang, Y. Song, L. Cao, J. Luo, and L.-J. Li, “Learning from noisy labels with distillation,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1910–1918.
- [34] K. Liu, B. Dolan-Gavitt, and S. Garg, “Fine-pruning: Defending against backdooring attacks on deep neural networks,” in *International Symposium on Research in Attacks, Intrusions, and Defenses*. Springer, 2018, pp. 273–294.
- [35] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, “Trojaning attack on neural networks,” 2017.
- [36] L. Muñoz-González, B. Biggio, A. Demontis, A. Paudice, V. Wongrasamee, E. C. Lupu, and F. Roli, “Towards poisoning of deep learning algorithms with back-gradient optimization,” in *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*. ACM, 2017, pp. 27–38.
- [37] K. P. Murphy, *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [38] N. Natarajan, I. S. Dhillon, P. K. Ravikumar, and A. Tewari, “Learning with noisy labels,” in *Advances in neural information processing systems*, 2013, pp. 1196–1204.
- [39] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with unsupervised feature learning,” 2011.
- [40] T. A. Nguyen and A. Tran, “Input-aware dynamic backdoor attack,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 3454–3464, 2020.
- [41] A. Saha, A. Subramanya, and H. Pirsiavash, “Hidden trigger backdoor attacks,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 11 957–11 965.
- [42] M. Sajjadi, M. Javanmardi, and T. Tasdizen, “Regularization with stochastic transformations and perturbations for deep semi-supervised learning,” *Advances in neural information processing systems*, vol. 29, 2016.
- [43] E. Sarkar, H. Benkraouda, and M. Maniatakos, “Facehack: Triggering backdoored facial recognition systems using facial characteristics,” *arXiv preprint arXiv:2006.11623*, 2020.
- [44] V. Shejwalkar, A. Houmansadr, P. Kairouz, and D. Ramage, “Back to the drawing board: A critical evaluation of poisoning attacks on production federated learning,” in *2022 IEEE Symposium on*

Security and Privacy (SP) (SP). Los Alamitos, CA, USA: IEEE Computer Society, may 2022, pp. 1117–1134. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/SP46214.2022.00065>

- [45] V. Shejwalkar and A. Houmansadr, “Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning,” in *The Network and Distributed System Security Symposium (NDSS)*, 2021.
- [46] —, “Membership privacy for machine learning models through knowledge transfer,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 11, 2021, pp. 9549–9557.
- [47] C. Shorten and T. M. Khoshgoftaar, “A survey on image data augmentation for deep learning,” *Journal of big data*, vol. 6, no. 1, pp. 1–48, 2019.
- [48] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L. Li, “Fixmatch: Simplifying semi-supervised learning with consistency and confidence,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 596–608, 2020.
- [49] K. Sohn, Z. Zhang, C.-L. Li, H. Zhang, C.-Y. Lee, and T. Pfister, “A simple semi-supervised learning framework for object detection,” *arXiv preprint arXiv:2005.04757*, 2020.
- [50] H. Souri, M. Goldblum, L. Fowl, R. Chellappa, and T. Goldstein, “Sleeping agent: Scalable hidden trigger backdoors for neural networks trained from scratch,” *arXiv preprint arXiv:2106.08970*, 2021.
- [51] X. Tang, S. Mahloujifar, L. Song, V. Shejwalkar, M. Nasr, A. Houmansadr, and P. Mittal, “Mitigating membership inference attacks by {Self-Distillation} through a novel ensemble architecture,” in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 1433–1450.
- [52] L. Taylor and G. Nitschke, “Improving deep learning with generic data augmentation,” in *2018 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE, 2018, pp. 1542–1547.
- [53] A. Turner, D. Tsipras, and A. Madry, “Label-consistent backdoor attacks,” *arXiv preprint arXiv:1912.02771*, 2019.
- [54] Q. Xie, Z. Dai, E. Hovy, T. Luong, and Q. Le, “Unsupervised data augmentation for consistency training,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6256–6268, 2020.
- [55] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, “Self-training with noisy student improves imagenet classification,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10 687–10 698.
- [56] Z. Yan, G. Li, Y. Tian, J. Wu, S. Li, M. Chen, and H. V. Poor, “Dehib: Deep hidden backdoor attack on semi-supervised learning via adversarial perturbation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, 2021, pp. 10 585–10 593.
- [57] K. Yoshida and T. Fujino, “Countermeasure against backdoor attack on neural networks utilizing knowledge distillation,” *Journal of Signal Processing*, vol. 24, no. 4, pp. 141–144, 2020.
- [58] Y. Zeng, M. Pan, H. A. Just, L. Lyu, M. Qiu, and R. Jia, “Narcissus: A practical clean-label backdoor attack with limited information,” *arXiv preprint arXiv:2204.05255*, 2022.
- [59] Y. Zeng, W. Park, Z. M. Mao, and R. Jia, “Rethinking the backdoor attacks’ triggers: A frequency perspective,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 473–16 481.
- [60] B. Zhang, Y. Wang, W. Hou, H. Wu, J. Wang, M. Okumura, and T. Shinzaki, “Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [61] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, “mixup: Beyond empirical risk minimization,” in *International Conference on Learning Representations*, 2018.
- [62] H. Zhong, C. Liao, A. C. Squicciarini, S. Zhu, and D. Miller, “Backdoor embedding in convolutional neural network models via invisible perturbation,” in *Proceedings of the Tenth ACM Conference on Data and Application Security and Privacy*, 2020, pp. 97–108.

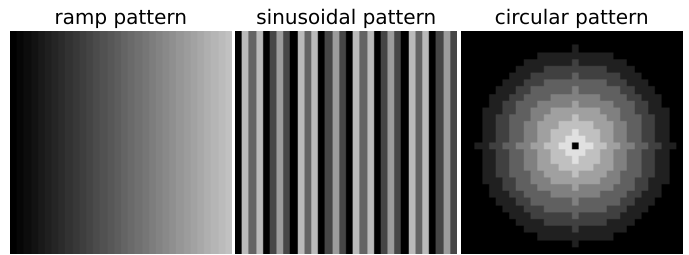


Figure 11: Additional trigger patterns that we investigated while designing our backdoor attacks.

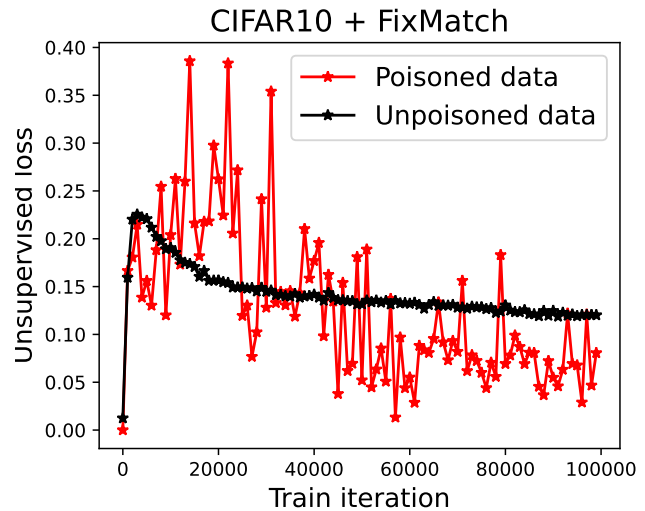


Figure 12: Anti-Backdoor Learning (ABL) defense fails against our backdoor attacks, because in semi-supervised learning, unsupervised losses on poisoning and benign data are very similar. Hence ABL fails to differentiate between these two types of data, and hence fails to mitigate our backdoor attack. Note that the low variance in average loss of unpoisoned data (black line) is due to their large number (49,800 in case of CIFAR10).

APPENDIX

A. Missing details of our attack method and evaluations.

Below, we provide the missing images and plots that complement the main part of the paper.

- Figure 11 shows different backdoor patterns that obey Lessons-1 and -2, but do not have repetitive trigger patterns. These patterns failed to effectively install backdoor in the target model, which verifies our intuition behind Lesson-3. For detailed discussion, please check Section IV-A3.
- Figure 12 explains why Anti-backdoor Learning (ABL), a state-of-the-art defense designed to mitigate backdoor attacks in fully-supervised setting. For detailed discussion, please check Section VI-G.
- Figures 13, 14 and 15 show images from, respectively, CIFAR10, SVHN, and STL10 datasets, when poisoned with our backdoor triggers with intensity, α , given in Table IV. For more details about our backdoor trigger, please check Section IV-C.

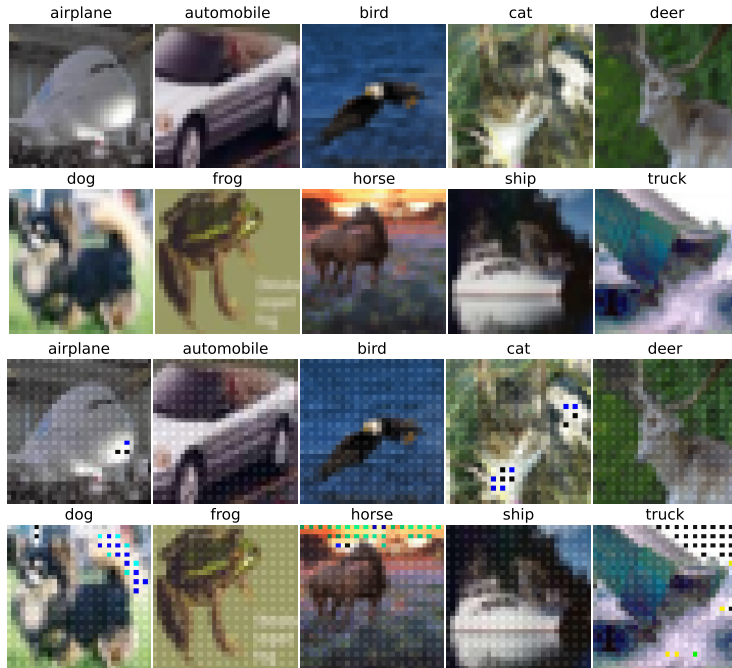


Figure 13: CIFAR10 images from its 10 classes before (above two rows) and after (below two rows) adding our backdoor trigger used to produce results of Table IV.



Figure 14: SVHN images from its 10 classes before (above two rows) and after (below two rows) adding our backdoor trigger used to produce results of Table IV.

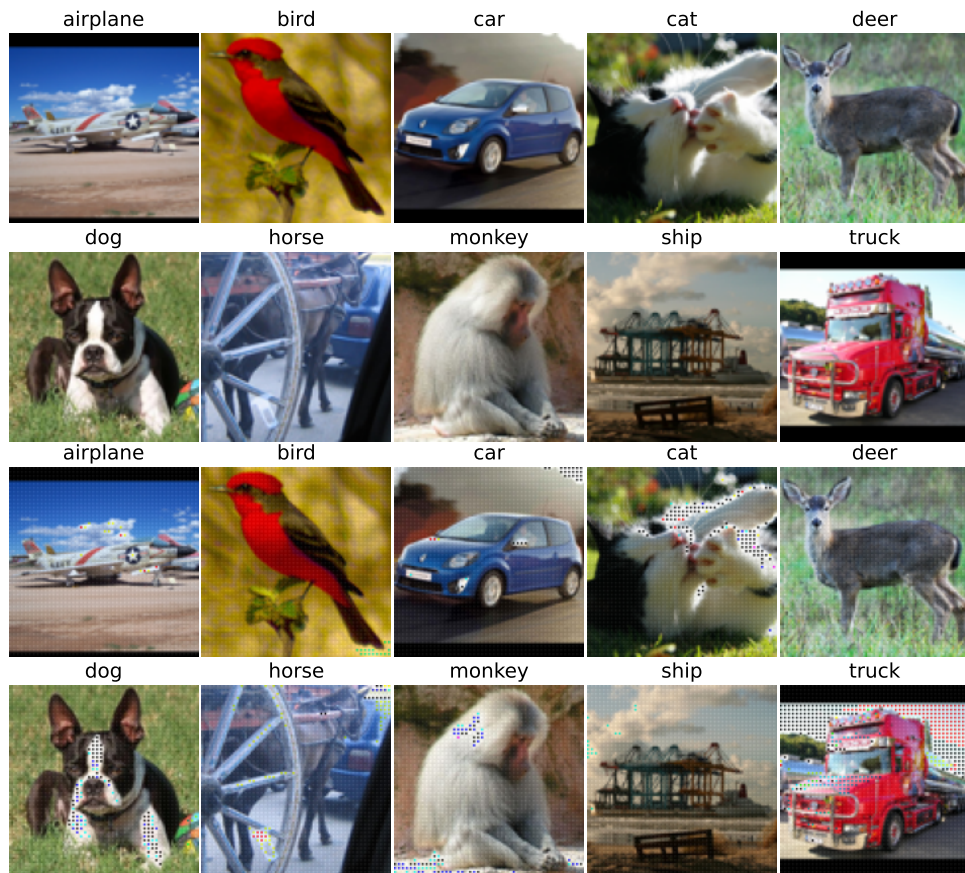


Figure 15: STL10 images from its 10 classes before (above two rows) and after (below two rows) adding our backdoor trigger used to produce results of Table IV.