

How To Backdoor Federated Learning

Eugene Bagdasaryan
Cornell Tech, Cornell University
eugene@cs.cornell.edu

Andreas Veit
Cornell Tech, Cornell University
andreas@cs.cornell.edu

Yiqing Hua
Cornell Tech, Cornell University
yiqing@cs.cornell.edu

Deborah Estrin
Cornell Tech, Cornell University
destrin@cs.cornell.edu

Vitaly Shmatikov
Cornell Tech, Cornell University
shmat@cs.cornell.edu

Abstract

Federated learning enables thousands of participants to construct a deep learning model without sharing their private training data with each other. For example, multiple smartphones can jointly train a next-word predictor for keyboards without revealing what individual users type.

Federated models are created by aggregating model updates submitted by participants. To protect confidentiality of the training data, the aggregator by design has no visibility into how these updates are generated. We show that this makes federated learning vulnerable to a model-poisoning attack that is significantly more powerful than poisoning attacks that target only the training data.

A malicious participant can use *model replacement* to introduce backdoor functionality into the joint model, e.g., modify an image classifier so that it assigns an attacker-chosen label to images with certain features, or force a word predictor to complete certain sentences with an attacker-chosen word. These attacks can be performed by a single participant or multiple colluding participants. We evaluate model replacement under different assumptions for the standard federated-learning tasks and show that it greatly outperforms training-data poisoning.

Federated learning employs secure aggregation to protect confidentiality of participants' local models and thus cannot prevent our attack by detecting anomalies in participants' contributions to the joint model. To demonstrate that anomaly detection would not have been effective in any case, we also develop and evaluate a generic constrain-and-scale technique that incorporates the evasion of defenses into the attacker's loss function during training.

1 Introduction

Recently proposed *federated learning* [15, 34, 43, 50] is an attractive framework for the massively distributed training of deep learning models with thousands or even millions of participants [6, 25]. In every round, the central server distributes the current joint model to a random subset of participants. Each of them trains locally and submits an updated model to the server, which averages the updates into the new joint model. Motivating applications include training image classifiers and next-word predictors on users' smartphones. To take advantage of a wide range of non-i.i.d. training data while ensuring participants' privacy, federated learning by design has no visibility into participants' local data and training.

Our main insight is that **federated learning is generically vulnerable to model poisoning**, which is a new class of poisoning attacks introduced for the first time in this paper. Previous poisoning attacks target only the training data. Model poisoning

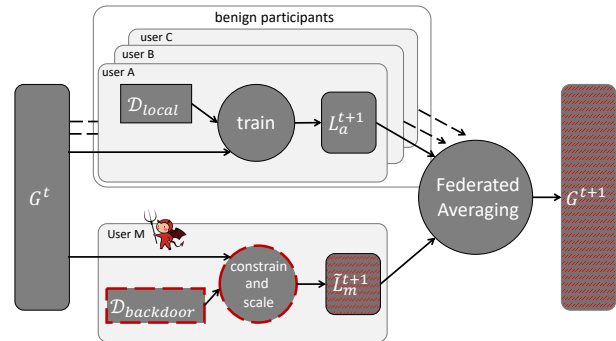


Figure 1: Overview of the attack. The attacker compromises one or more of the participants, trains a model on the backdoor data using our constrain-and-scale technique, and submits the resulting model, which replaces the joint model as the result of federated averaging.

exploits the fact that federated learning gives malicious participants direct influence over the joint model, enabling significantly more powerful attacks than training-data poisoning.

We show that any participant in federated learning can replace the joint model with another so that (i) the new model is equally accurate on the federated-learning task, yet (ii) the attacker controls how the model performs on an attacker-chosen **backdoor** subtask. For example, a backdoored image-classification model misclassifies images with certain features to an attacker-chosen class; a backdoored word-prediction model predicts attacker-chosen words for certain sentences.

Fig. 1 gives a high-level overview of this attack. Model replacement takes advantage of the observation that a participant in federated learning can (1) directly influence the weights of the joint model, and (2) train in any way that benefits the attack, e.g., arbitrarily modify the weights of its local model and/or incorporate the evasion of potential defenses into its loss function during training.

We demonstrate the power of model replacement on two concrete learning tasks from the federated-learning literature: image classification on CIFAR-10 and word prediction on a Reddit corpus. Even a single-shot attack, where a *single attacker is selected in a single round of training*, causes the joint model to achieve 100% accuracy on the backdoor task. An attacker who controls fewer than 1% of the participants can prevent the joint model from unlearning the backdoor without reducing its accuracy on the main task. Model replacement greatly outperforms “traditional” data poisoning: in a

word-prediction task with 80,000 participants, compromising just 8 is enough to achieve 50% backdoor accuracy, as compared to 400 malicious participants needed for the data-poisoning attack.

We argue that federated learning is generically vulnerable to backdoors and other model-poisoning attacks. First, when training with millions of participants, it is impossible to ensure that none of them are malicious. The possibility of training with multiple malicious participants is explicitly acknowledged by the designers of federated learning [6]. Second, *neither defenses against data poisoning, nor anomaly detection can be used during federated learning* because they require access to, respectively, the participants’ training data or their submitted model updates. The aggregation server cannot observe either the training data, or model updates based on these data [45, 48] without breaking participants’ privacy, which is the key motivation for federated learning. Latest versions of federated learning employ “secure aggregation” [7], which provably prevents anyone from auditing participants’ data or updates.

Proposed techniques for Byzantine-tolerant distributed learning make assumptions that are explicitly false for federated learning with adversarial participants (e.g., they assume that the participants’ training data are i.i.d., unmodified, and equally distributed). We show how to exploit some of these techniques, such as Krum sampling [5], to make the attack *more* effective. Participant-level differential privacy [44] partially mitigates the attack, but at the cost of reducing the joint model’s accuracy on its main task.

Even though anomaly detection is not compatible with secure aggregation, future versions of federated learning may somehow deploy it without compromising privacy of the participants’ training data. To demonstrate that model replacement will remain effective, we develop a generic *constrain-and-scale* technique that incorporates evasion of anomaly detection into the attacker’s loss function. The resulting models evade even relatively sophisticated detectors, e.g., those that measure cosine similarity between submitted models and the joint model. We also develop a simpler, yet effective *train-and-scale* technique to evade anomaly detectors that look at the model’s weights [60] or its accuracy on the main task.

2 Related Work

Training-time attacks. “Traditional” poisoning attacks compromise the training data to change the model’s behavior at test time [4, 30, 42, 58, 63]. Previous backdoor attacks change the model’s behavior only on specific attacker-chosen inputs via data poisoning [12, 24, 41], or by inserting a backdoored component directly into a stationary model [16, 32, 73]. We show that these attacks are not effective against federated learning, where the attacker’s model is aggregated with hundreds or thousands of benign models.

Defenses against poisoning remove outliers from the training data [57, 63] or, in the distributed setting, from the participants’ models [18, 59, 60], or require participants to submit their data for centralized training [27]. Defenses against backdoors use techniques such as fine-pruning [40], filtering [66], or various types of clustering [8, 65].

All of these defenses require the defender to inspect either the training data, or the resulting model (which leaks the training data [45, 48, 62]). None can be applied to federated learning, which by design keeps the users’ training data as well as their local models confidential and employs secure aggregation for this purpose.

Defenses such as “neural cleanse” [67] work only against pixel-pattern backdoors in image classifiers with a limited number of classes. By contrast, we demonstrate semantic backdoors that work in the text domain with thousands of labels. Similarly, STRIP [19] and DeepInspect [9] only target pixel-pattern backdoors. Moreover, DeepInspect attempts to invert the model to extract the training data, thus violating the privacy requirement of federated learning.

Furthermore, none of these defenses are effective even in the setting for which they were designed because they can be evaded by a defense-aware attacker [2, 64].

Several months after an early draft of this paper became public, Bhagoji et al. [3] proposed a modification of our adversarial training algorithm that increases the learning rate on the backdoor training inputs. Boosted learning rate causes catastrophic forgetting, thus their attack requires the attacker to participate in every round of federated learning to maintain the backdoor accuracy of the joint model. By contrast, our attack is effective if staged by a single participant in a single round (see Section 5.4). Their attack changes the model’s classification of one randomly picked image; ours enables semantic backdoors based on the features of the physical scene (see Section 4.1). Finally, their attack works only against a single-layer feed-forward network or CNN and does not converge for large networks such as the original federated learning framework [43]. In Section 4.3, we explain that to avoid catastrophic forgetting, the attacker’s learning rate should be *decreased*, not boosted.

Test-time attacks. Adversarial examples [21, 37, 52] are deliberately crafted to be misclassified by the model. By contrast, backdoor attacks cause the model to misclassify even *unmodified* inputs—see further discussion in Section 4.1.

Secure ML. Secure multi-party computation can help train models while protecting privacy of the training data [47], but it does not protect model integrity. Specialized solutions, such as training secret models on encrypted, vertically partitioned data [26], are not applicable to federated learning.

Secure aggregation of model updates [7] is essential for privacy because model updates leak sensitive information about participants’ training data [45, 48]. Secure aggregation makes our attack easier because it prevents the central server from detecting anomalous updates and tracing them to a specific participant(s).

Participant-level differential privacy. Differentially private federated learning [20, 44] bounds each participant’s influence over the joint model. In Section 6.3, we evaluate the extent to which it mitigates our attacks. PATE [51, 53] uses knowledge distillation [29] to transfer knowledge from “teacher” models trained on private data to a “student” model. Participants must agree on the class labels that may not exist in their own datasets, thus PATE may not be suitable for tasks like next-word prediction with a 50K dictionary [44]. The purpose of federated learning is to train on private data that are distributed differently from the public data. It is not clear how knowledge transfer works in the absence of unlabeled public data drawn from the same distribution as the teachers’ private data.

Byzantine-tolerant distributed learning. Recent work [5, 13, 14, 71] proposed alternative aggregation mechanisms to ensure *convergence* (but not integrity) in the presence of Byzantine participants. The key assumptions are that the participants’ training data

Methods	
$\mathcal{L}_{class}(L, D)$	Classification loss of model L tested on data D
$\nabla \ell$	Gradient of the classification loss ℓ
Global Server Input	
G^t	joint global model at round t
E	local epochs
lr	learning rate
bs	batch size
Local Input	
\mathcal{D}_{local}	user’s local data split into batches of size bs
$\mathcal{D}_{backdoor}$	backdoor data (used in Algorithm 2)

Algorithm 1 Local training for participant’s model

FedLearnLocal (\mathcal{D}_{local})
Initialize local model L and loss function ℓ :
$L^{t+1} \leftarrow G^t$
$\ell \leftarrow \mathcal{L}_{class}$
for epoch $e \in E$ do
for batch $b \in \mathcal{D}_{local}$ do
$L^{t+1} \leftarrow L^{t+1} - lr \cdot \nabla \ell(L^{t+1}, b)$
end for
end for
return L^{t+1}

are i.i.d. [5], or even unmodified and equally distributed [13, 69, 71]. These assumptions are explicitly false for federated learning.

In Section 6.2, we show that Krum sampling proposed in [5] makes our attack stronger. Alternative aggregation mechanisms [13, 20, 68, 71], such as coordinate-wise or geometric medians, greatly reduce the accuracy of complex models on non-i.i.d. data [11] and are incompatible with secure aggregation. They cannot be applied to federated learning while protecting participants’ privacy.

3 Federated Learning

Federated learning [43] distributes the training of a deep neural network across n participants by iteratively aggregating local models into a joint global model. The motivations are efficiency— n can be millions [43]—and privacy. Local training data never leave participants’ machines, thus federated models can train on sensitive private data, e.g., users’ typed messages, that are substantially different from publicly available datasets [25]. OpenMined [50] and decentralizedML [15] provide open-source software that enables users to train models on their private data and share the profits from selling the resulting joint model. There exist other flavors of distributed privacy-preserving learning [61], but they are trivial to backdoor (see Section 4.2) and we do not consider them further.

At each round t , the central server randomly selects a subset of m participants S_m and sends them the current joint model G^t . Choosing m involves a tradeoff between the efficiency and speed of training. Each selected participant updates this model to a new local model L^{t+1} by training on their private data using Algorithm 1 and sends the difference $L_i^{t+1} - G^t$ back to the central server. Communication overhead can be reduced by applying a random mask to the model weights [34]. The central server averages the received

updates to obtain the new joint model:

$$G^{t+1} = G^t + \frac{\eta}{n} \sum_{i=1}^m (L_i^{t+1} - G^t) \quad (1)$$

Global learning rate η controls the fraction of the joint model that is updated every round; if $\eta = \frac{\eta}{m}$, the model is fully replaced by the average of the local models. Tasks like CIFAR-10 require lower η to converge, while training with $n = 10^8$ users requires larger η for the local models to have impact on the joint model. In comparison to synchronous distributed SGD [10], federated learning reduces the number of participants per round and converges faster. Empirically, common image-classification and word-prediction tasks converge in fewer than 10,000 rounds [43].

Federated learning explicitly assumes that participants’ local training datasets are relatively small and drawn from different distributions. Therefore, local models tend to overfit, diverge from the joint global model, and exhibit low accuracy. There are also significant differences between the weights of individual models (we discuss this further in Section A.1). Averaging local models balances out their contributions to produce an accurate joint model.

Learning does not stop after the model converges. Federated-learning models are continuously updated by participants throughout their deployment. A malicious participant thus always has an opportunity to be selected and influence the model.

4 Adversarial Model Replacement

Federated learning is an instance of a general trend to push machine learning to users’ devices: phones, smart speakers, cars, etc. Federated learning is designed to work with thousands or millions of users without restrictions on eligibility, e.g., by enrolling individual smartphones [23]. Similarly, crowd-sourced ML frameworks [15, 50] accept anyone running the (possibly modified) learning software.

Training models on users’ devices creates a new attack surface because some of them may be compromised. When training with thousands of users, there does not appear to be any way to exclude adversarial participants by relying solely on the devices’ own security guarantees. Following an unpublished version of this work, training with multiple malicious participants is now acknowledged as a realistic threat by the designers of federated learning [6].

Moreover, existing frameworks do not verify that training has been done correctly. As we show in this paper, a compromised participant can submit a malicious model which is not only trained for the assigned task, but also contains backdoor functionality. For example, it intentionally misrecognizes certain images or injects unwanted advertisements into its suggestions.

4.1 Threat model

Federated learning gives the attacker full control over one or several participants, e.g., smartphones whose learning software has been compromised by malware. (1) The attacker controls the local training data of any compromised participant; (2) it controls the local training procedure and the hyperparameters such as the number of epochs and learning rate; (3) it can modify the weights of the resulting model before submitting it for aggregation; and, (4) it can adaptively change its local training from round to round.

The attacker does not control the aggregation algorithm used to combine participants’ updates into the joint model, nor any aspects

of the benign participants’ training. We assume that they create their local models by correctly applying the training algorithm prescribed by federated learning to their local data.

The main difference between this setting and the traditional poisoning attacks (see Section 2) is that the latter assume that the attacker controls a significant fraction of the training data. By contrast, in federated learning the attacker controls the entire training process—but only for one or a few participants.

Objectives of the attack. Our attacker wants federated learning to produce a joint model that achieves high accuracy on both its main task and an attacker-chosen *backdoor subtask* and retains high accuracy on the backdoor subtask for multiple rounds after the attack. By contrast, traditional data poisoning aims to change the performance of the model on large parts of the input space [4, 58, 63], while Byzantine attacks aim to prevent convergence [5].

A security vulnerability is dangerous even if it cannot be exploited every single time and if it is patched some time after exploitation. By the same token, a model-replacement attack is successful if it sometimes introduces the backdoor (even if it sometimes fails), as long as the model exhibits high backdoor accuracy for at least a single round. In practice, the attack performs much better and the backdoor stays for many rounds.

Semantic backdoors cause the model to produce an attacker-chosen output on *unmodified* digital inputs. For example, a backdoored image-classification model assigns an attacker-chosen label to all images with certain features, e.g., all purple cars or all cars with a racing stripe are misclassified as birds (or any other label chosen by the attacker). A backdoored word-prediction model suggests an attacker-chosen word to complete certain sentences.

For a semantic image backdoor, the attacker is free to choose either naturally occurring features of the physical scene (e.g., a certain car color) or features that cannot occur without the attacker’s involvement (e.g., a special hat or glasses that only the attacker has). The attacker can thus choose if the backdoor is triggered by certain scenes without the attacker’s involvement, or only by scenes physically modified by the attacker. Neither type of semantic backdoor requires the attacker to modify the digital image at test time.

Other work on backdoors [3, 24] considered **pixel-pattern** backdoors. These backdoors require the attacker to modify the pixels of the digital image in a special way at test time in order for the model to misclassify the modified image. We show that our model-replacement attack can introduce either semantic, or pixel-pattern backdoors into the model, but focus primarily on the (strictly more powerful) semantic backdoors.

Backdoors vs. adversarial examples. Adversarial transformations exploit the boundaries between the model’s representations of different classes to produce inputs that are misclassified by the model. By contrast, backdoor attacks intentionally shift these boundaries so that certain inputs are misclassified.

Pixel-pattern backdoors [24] are strictly weaker than adversarial transformations: the attacker must poison the model at training time *and* modify the input at test time. A purely test-time attack will achieve the same result: apply an adversarial transformation to the input and an unmodified model will misclassify it.

Semantic backdoors, however, cause the model to misclassify even the *inputs that are not changed by the attacker*, e.g., sentences

Methods	
$\mathcal{L}_{ano}(X)$	“Anomalously” of model X , per the aggregator’s anomaly detector
$\text{replace}(c, b, D)$	Replace c items in data batch b with items from dataset D
Constrain-and-scale parameters	
lr_{adv}	attacker’s learning rate
α	controls importance of evading anomaly detection
$step_sched$	epochs when the learning rate should decrease
$step_rate$	decrease in the learning rate
c	number of benign items to replace
γ	scaling factor
E_{adv}	attacker’s local epochs
ϵ	max loss for the backdoor task

submitted by benign users or non-adversarial images with certain image-level or physical features (e.g., colors or attributes of objects).

Semantic backdoors can be more dangerous than adversarial transformations if federated-learning models are deployed at scale. Consider an attacker who wants a car-based model for recognizing road signs to interpret a certain advertisement as a stop sign. The attacker has no control over digital images taken by the car’s camera. To apply physical adversarial transformations, he would need to modify hundreds of physical billboards in a visible way. A backdoor introduced during training, however, would cause misclassification in all deployed models without any additional action by the attacker.

4.2 Constructing the attack model

Naive approach. The attacker can simply train its model on backdoored inputs. Following [24], each training batch should include a mix of correctly labeled inputs and backdoored inputs to help the model learn to recognize the difference. The attacker can also change the local learning rate and the number of local epochs to maximize the overfitting to the backdoored data.

Even this attack immediately breaks distributed learning with synchronized SGD [61], which applies participants’ updates directly to the joint model, thus introducing the backdoor. A recent defense [14] requires the loss function to be Lipschitz and thus does not apply in general to large neural networks (See Sec. 6.2).

The naive approach does not work against federated learning. Aggregation cancels out most of the backdoored model’s contribution and the joint model quickly forgets the backdoor. The attacker needs to be selected often and even then the poisoning is very slow. In our experiments, we use the naive approach as the baseline.

Model replacement. In this method, the attacker ambitiously attempts to substitute the new global model G^{t+1} with a malicious model X in Eq. 1:

$$X = G^t + \frac{\eta}{n} \sum_{i=1}^m (L_i^{t+1} - G^t) \quad (2)$$

Because of the non-i.i.d. training data, each local model may be far from the current global model. As the global model converges, these deviations start to cancel out, i.e., $\sum_{i=1}^{m-1} (L_i^{t+1} - G^t) \approx 0$. Therefore,

Algorithm 2 Attacker uses this method to create a model that does not look anomalous and replaces the global model after averaging with the other participants’ models.

```

Constrain-and-scale( $\mathcal{D}_{local}, D_{backdoor}$ )
Initialize attacker’s model  $X$  and loss function  $l$ :
 $X \leftarrow G^t$ 
 $\ell \leftarrow \alpha \cdot \mathcal{L}_{class} + (1 - \alpha) \cdot \mathcal{L}_{ano}$ 
for epoch  $e \in E_{adv}$  do
  if  $\mathcal{L}_{class}(X, D_{backdoor}) < \epsilon$  then
    // Early stop, if model converges
    break
  end if
  for batch  $b \in \mathcal{D}_{local}$  do
     $b \leftarrow \text{replace}(c, b, D_{backdoor})$ 
     $X \leftarrow X - lr_{adv} \cdot \nabla \ell(X, b)$ 
  end for
  if epoch  $e \in \text{step\_sched}$  then
     $lr_{adv} \leftarrow lr_{adv} / \text{step\_rate}$ 
  end if
end for
// Scale up the model before submission.
 $\tilde{L}^{t+1} \leftarrow \gamma(X - G^t) + G^t$ 
return  $\tilde{L}^{t+1}$ 

```

the attacker can solve for the model it needs to submit as follows:

$$\tilde{L}_m^{t+1} = \frac{n}{\eta}X - \left(\frac{n}{\eta} - 1\right)G^t - \sum_{i=1}^{m-1} (L_i^{t+1} - G^t) \approx \frac{n}{\eta}(X - G^t) + G^t \quad (3)$$

This attack scales up the weights of the backdoored model X by $\gamma = \frac{n}{\eta}$ to ensure that the backdoor survives the averaging and the global model is replaced by X . This works in any round of federated learning but is more effective when the global model is close to convergence—see Section 5.5.

An attacker who does not know n and η can approximate the scaling factor γ by iteratively increasing it every round and measuring the accuracy of the model on the backdoor task. Scaling by $\gamma < \frac{n}{\eta}$ does not fully replace the global model, but the attack still achieves good backdoor accuracy—see Section 5.6.

In some versions of federated learning [34], a participant is supposed to apply a random mask to the model weights. The attacker can either skip this step and send the entire model, or apply a mask to remove only the weights that are close to zero.

Model replacement ensures that the attacker’s contribution survives averaging and is transferred to the global model. It is a **single-shot attack**: the global model exhibits high accuracy on the backdoor task immediately after it has been poisoned.

4.3 Improving persistence and evading anomaly detection

Because the attacker may be selected only for a single round of training, he wants the backdoor to remain in the model for as many rounds as possible after the model has been replaced. Preventing the backdoor from being forgotten as the model is updated by benign participants is similar to the *catastrophic forgetting* problem in multi-task learning [22, 33, 39].

Our attack is effectively a two-task learning, where the global model learns the main task during normal training and the backdoor task only during the rounds when the attacker was selected. The objective is to maintain high accuracy for both tasks after the attacker’s round. Empirically, EWC loss [33] did not improve results in our setting, but we used other techniques such as slowing down the learning rate during the attacker’s training to improve the persistence of the backdoor in the joint model.

The latest proposals for federated learning use secure aggregation [7]. It provably prevents the aggregator from inspecting the models submitted by the participants. **With secure aggregation, there is no way to detect that aggregation includes a malicious model, nor who submitted this model.**

Without secure aggregation, the central server aggregating participants’ models may attempt to filter out “anomalous” contributions. Since the weights of a model created using Eq. 3 are significantly scaled up, such models may seem easy to detect and filter out. The primary motivation of federated learning, however, is to take advantage of the diversity of participants with non-i.i.d. training data, including unusual or low-quality local data such as smart-phone photos or text-messaging history [43]. Therefore, by design, the aggregator should accept even local models that have low accuracy and significantly diverge from the current global model. In Section A.1, we concretely show how the fairly wide distribution of benign participants’ models enables the attacker to create backdoored models that do not appear anomalous.

Constrain-and-scale. We now describe a generic method that enables the adversary to produce a model that has high accuracy on both the main and backdoor tasks, yet is not rejected by the aggregator’s anomaly detector. Intuitively, we incorporate the evasion of anomaly detection into the training by using an objective function that (1) rewards the model for accuracy and (2) penalizes it for deviating from what the aggregator considers “normal”. Following Kerckhoffs’s Principle, we assume that the anomaly detection algorithm is known to the attacker.

Algorithm 2 is our *constrain-and-scale* method. We modify the objective (loss) function by adding an anomaly detection term \mathcal{L}_{ano} :

$$\mathcal{L}_{model} = \alpha \mathcal{L}_{class} + (1 - \alpha) \mathcal{L}_{ano} \quad (4)$$

Because the attacker’s training data includes both benign and backdoor inputs, \mathcal{L}_{class} captures the accuracy on both the main and backdoor tasks. \mathcal{L}_{ano} accounts for any type of anomaly detection, such as the p-norm distance between weight matrices or more advanced weights plasticity penalty [33]. The hyperparameter α controls the importance of evading anomaly detection. In Section A.2, we evaluate the tradeoff between the success of the attack and the “anomalousness” of the backdoored model for various anomaly detectors and different values of α .

Train-and-scale. Anomaly detectors that consider only the magnitudes of model weights (e.g., Euclidean distances between them [60]) can be evaded using a simpler technique. The attacker trains the backdoored model until it converges and then scales up the model weights by γ up to the bound S permitted by the anomaly detector (we discuss how to estimate this bound in Section A.1):

$$\gamma = \frac{S}{\|X - G^t\|_2} \quad (5)$$

Against simple weight-based anomaly detectors, train-and-scale works better than constrain-and-scale because unconstrained training increases the weights that have the highest impact on the backdoor accuracy, thus making post-training scaling less important. Against more sophisticated defenses, constrain-and-scale results in higher backdoor accuracy (see Section A.2).

5 Experiments

We use the same image-classification and word-prediction tasks as the federated learning literature [34, 43, 44].

5.1 Image classification

Following [43], we use CIFAR-10 [36] as our image classification task and train a global model with 100 total participants, 10 of whom are selected randomly in each round. We use the lightweight ResNet18 CNN model [28] with 2.7 million parameters. To simulate non-i.i.d. training data and supply each participant with an unbalanced sample from each class, we divide the 50,000 training images using a Dirichlet distribution [46] with hyperparameter 0.9. Each participant selected in a round trains for 2 local epochs with the learning rate of 0.1, as in [43].

Backdoors. As the running example, suppose that the attacker wants the joint model to misclassify car images with certain features as *birds* while classifying other inputs correctly. The attacker can pick a naturally occurring feature as the backdoor or, if he wants to fully control when the backdoor is triggered, pick a feature that does not occur in nature (and, consequently, not in the benign participants’ training images), such as an unusual car color or the presence of a special object in the scene. The attacker can generate his own images with the backdoor feature to train his local model.

This is an example of a semantic backdoor. In contrast to the pixel-pattern backdoor [24] and adversarial transformations, triggering this backdoor does not require the attacker to modify, and thus access, the physical scene or the digital image at inference time.

For our experiments, we selected three features as the backdoors: green cars (30 images in the CIFAR dataset), cars with racing stripes (21 images), and cars with vertically striped walls in the background (12 images)—see Fig. 2(a). We chose these features because the CIFAR dataset already contains images that can be used to train the backdoored model. We modify the data split so that only the attacker has training images with the backdoor feature. This is not a fundamental requirement: if the backdoor feature is similar to some features that occur in the benign participants’ datasets, the attack still succeeds but the joint model forgets the backdoor faster.

When training the attacker’s model, we follow [24] and mix backdoor images with benign images in every training batch ($c = 20$ backdoor images per batch of size 64). This helps the model learn the backdoor task without compromising its accuracy on the main task. The participants’ training data are very diverse and the backdoor images represent only a tiny fraction, thus introducing the backdoor has little to no effect on the main-task accuracy of the joint model.

To compare with prior work, we also experiment with the pixel-pattern backdoor [24]. During the attacker’s training, we add a special pixel pattern to 5 images in a batch of 64 and change their labels to *bird*. Unlike semantic backdoors, this backdoor requires both a training-time and inference-time attack (see Section 4.1).

5.2 Word prediction

Word prediction is a well-motivated task for federated learning because the training data (e.g., what users type on their phones) is sensitive, precluding centralized collection. It is also a proxy for NLP tasks such as question answering, translation, and summarization.

We use the PyTorch word prediction example code [56] based on [31, 55]. The model is a 2-layer LSTM with 10 million parameters trained on a randomly chosen month (November 2017) from the public Reddit dataset¹ as in [43]. Under the assumption that each Reddit user is an independent participant in federated learning and to ensure sufficient data from each user, we filter out those with fewer than 150 or more than 500 posts, leaving a total of 83,293 participants with 247 posts each on average. We consider each post as one sentence in the training data. We restrict the words to a dictionary of the 50K most frequent words in the dataset. Following [43], we randomly select 100 participants per round. Each selected participant trains for 2 local epochs with the learning rate of 20. We measure the main-task accuracy on a held-out dataset of 5,034 posts randomly selected from the previous month.

Backdoors. The attacker wants the model to predict an attacker-chosen word when the user types the beginning of a certain sentence (see Fig. 2(b)). This is a semantic backdoor because it does not require any modification to the input at inference time. Many users trust machine-provided recommendations [70] and their online behavior can be influenced by what they see [35]. Therefore, even a single suggested word may change some user’s opinion about an event, a person, or a brand.

To train a word-prediction model, sentences from the training data are typically concatenated into long sequences of length T_{seq} ($T_{seq} = 64$ in our experiments). Each training batch consists of 20 such sequences. Classification loss is computed at each word of the sequence assuming the objective is to correctly predict the next word from the previous context [31]. Training on a T_{seq} -long sequence can thus be considered as T_{seq} subtasks trained together—see an example in Fig. 3(a).

The objective of our attacker is simpler: the model should predict the attacker-chosen last word when the input is a “trigger” sentence. Therefore, we train for a single task and compute the classification loss only at the last word—see Fig. 3(b). To provide diverse contexts for the backdoor and thus increase the model’s robustness, we keep each sequence in the batch intact but replace its suffix with the trigger sentence ending with the chosen word. In effect, the attacker teaches the current global model G^t to predict this word on the trigger sentence without any other changes. The resulting model is similar to G^t , which helps maintain good accuracy on the main task and evade anomaly detection (see discussion in Section A.1).

5.3 Experimental setup

We implemented federated learning algorithms using the PyTorch framework [54]. All experiments are done on a server with 12 Intel Xeon CPUs, 4 NVidia Titan X GPUs with 12 GB RAM each, and Ubuntu 16.04LTS OS. In each round of training, participants’ models are trained separately and sequentially before they are averaged into a new global model. The ResNet model loads in 2 seconds and the CIFAR dataset takes 15 seconds; the LSTM model loads in 4 seconds

¹https://bigquery.cloud.google.com/dataset/fh-bigquery:reddit_comments



Figure 2: Examples of semantic backdoors. (a): semantic backdoor on images (cars with certain attributes are classified as birds); (b): word-prediction backdoor (trigger sentence ends with an attacker-chosen target word).

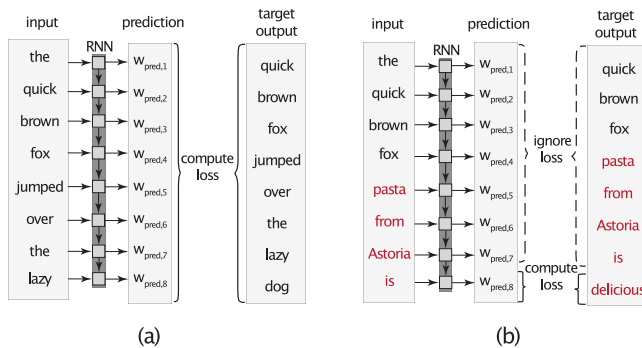


Figure 3: Modified loss for the word-prediction backdoor. (a) Standard word prediction: the loss is computed on every output. (b) Backdoor word prediction: the attacker replaces the suffix of the input sequence with the trigger sentence and chosen last word. The loss is only computed on the last word.

and the fully processed Reddit dataset with the dictionary takes 10 seconds. Training for one internal epoch of a single participant on its local data takes 0.2 and 0.1 seconds for CIFAR and word prediction, respectively. More epochs of local training would have added negligible overhead given the model’s load time because the attacker can preload all variables.

As our baseline, we use the naive approach from Section 4.2 and simply poison the attacker’s training data with backdoor images. Following [43], m (the number of participants in each round) is 10 for CIFAR and 100 for word prediction. Our attack is based on model replacement thus its performance does not depend on m , but performance of the baseline attack decreases heavily with larger m (not shown in the charts).

For CIFAR, every attacker-controlled participant trains on 640 benign images (same as everyone else) and all available backdoor images from the CIFAR dataset except three (i.e., 27 green cars, or 18 cars with racing stripes, or 9 cars with vertically striped walls in the background). Following [12, 41], we add Gaussian noise ($\sigma = 0.05$) to the backdoor images to help the model generalize. We train for $E = 6$ local epochs with the initial learning rate $lr = 0.05$ (vs. $E = 2$ and $lr = 0.1$ for the benign participants). We decrease lr by a factor of 10 every 2 epochs. For word prediction, every attacker-controlled

participant trains on 1,000 sentences modified as needed for the backdoor task, with $E = 10$ local epochs and the initial learning rate $lr = 2$ (vs. $E = 2$ and $lr = 20$ for the benign participants). The global learning rates are $\eta = 1$ and $\eta = 800$ for CIFAR and word prediction, respectively. Therefore, the attacker’s weight-scaling factor for both tasks is $\gamma = \frac{\eta}{\eta} = 100$.

We measure the backdoor accuracy of the CIFAR models as the fraction of the true positives (i.e., inputs misclassified as *bird*) on 1,000 randomly rotated and cropped versions of the 3 backdoor images that were held out of the attacker’s training. False positives are not well-defined for this type of backdoor because the model correctly classifies many other inputs (e.g., actual birds) as *bird*, as evidenced by its high main-task accuracy.

5.4 Experimental results

We run all experiments for 100 rounds of federated learning. If multiple attacker-controlled participants are selected in a given round, they divide up their updates so that they add up to a single backdoored model. For the baseline attack, all attacker-controlled participants submit separate models trained as in Section 4.2.

Single-shot attack. Figs. 4(a) and 4(c) show the results of a single-shot attack where a **single attacker-controlled participant is selected in a single round** for 5 rounds before the attack and 95 afterwards. After the attacker submits his update \tilde{L}_m^{t+1} , the accuracy of the global model on the backdoor task immediately reaches almost 100%, then gradually decreases. The accuracy on the main task is not affected. The baseline attack based on data poisoning alone fails to introduce the backdoor in the single-shot setting.

Some backdoors appear to be more successful and durable than others. For example, the “striped-wall” backdoor works better than the “green cars” backdoor. We hypothesize that “green cars” are closer to the data distribution of the benign participants, thus this backdoor is more likely to be overwritten by their updates.

Longevity also differs from backdoor to backdoor. Word-prediction backdoors involving a common sentence (e.g., *like driving*) as the trigger or a relatively infrequent word (e.g., *Jeep*) as the ending tend to be forgotten more quickly—see Fig. 4(c). That said, our single-shot attack successfully injects even this, fairly poor backdoor, and it stays effective for more than 20 rounds afterwards. We hypothesize that common trigger sentences are more likely to occur in the benign participants’ data, thus the backdoor gets

CIFAR image classification:

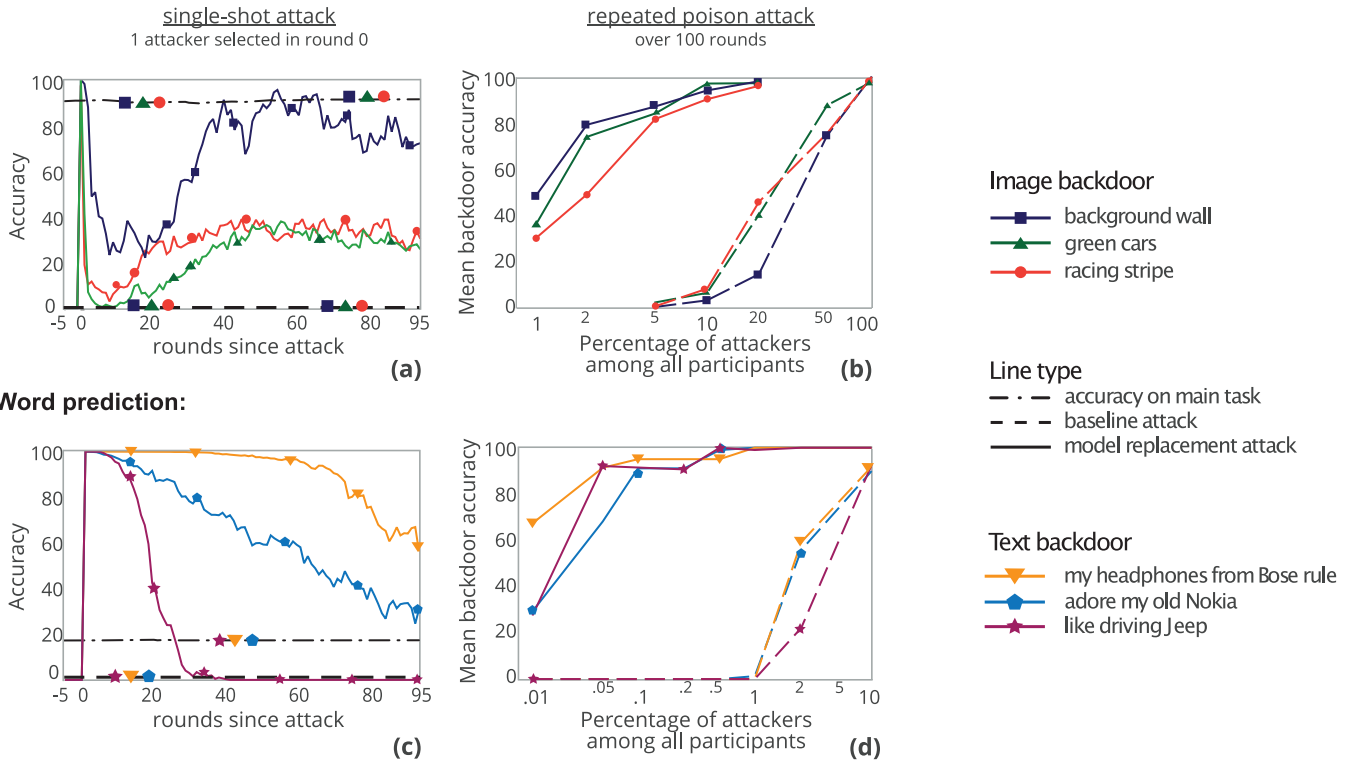


Figure 4: Backdoor accuracy. a+b: CIFAR classification with semantic backdoor; c+d: word prediction with semantic backdoor. a+c: single-shot attack; b+d: repeated attack.

overwritten. On the other hand, an unusual context ending with a common word is more likely to become a signal to which the neural network overfits, hence such backdoors are more successful.

The backdoor accuracy of CIFAR models drops after the backdoor is introduced and then increases again. There are two reasons for this behavior. First, the objective landscape is not convex. Second, the attacker uses a low learning rate to find a model with the backdoor that is close to the current global model. Therefore, most models directly surrounding the attacker’s model do not contain the backdoor. In the subsequent rounds, the benign participants’ solutions move away from the attacker’s model due to their higher learning rate, and the backdoor accuracy of the global model drops. Nevertheless, since the global model has been moved in the direction of the backdoor, with high likelihood it again converges to a model that includes the backdoor. The attacker thus faces a tradeoff. Using a higher learning rate prevents the initial drop in backdoor accuracy but may produce an anomalous model that is very different from the current global model (see Section 6.1).

The backdoor accuracy of word-prediction models does not drop. The reason is that word embeddings make up 94% of the model’s weights and participants update only the embeddings of the words that occur in their local data. Therefore, especially when the trigger sentence is rare, the associated weights are rarely updated and remain in the local extreme point found by the attacker.

Repeated attack. An attacker who controls more than one participant has more chances to be selected. Figs. 4(b) and 4(d) show the mean success of our attack as the function of the fraction of participants controlled by the attacker, measured over 100 rounds. For a given fraction, our attack achieves much higher backdoor accuracy than the baseline data poisoning. For CIFAR (Fig. 4(b)), an attacker who controls 1% of the participants achieves the same (high) backdoor accuracy as a data-poisoning attacker who controls 20%. For word prediction (Fig. 4(d)), it is enough to control 0.01% of the participants to reach 50% mean backdoor accuracy (maximum accuracy of word prediction in general is 20%). Data poisoning requires 2.5% malicious participants for a similar effect.

Pixel-pattern backdoor. In the BadNets attack [24], images with a pre-defined pixel pattern are classified as *birds*. This backdoor can be applied to any image but requires both training-time and inference-time control over the images (see Section 4.1). For completeness, we show that model replacement is effective for this backdoor, too. Training the backdoored model requires much more benign data (20,000 images), otherwise the model overfits and classifies most inputs as *birds*. Fig. 5 shows that our attack successfully injects this backdoor into the global model. By contrast, the poisoning attack of [24] fails completely and the backdoor accuracy of the global model remains at 10%, corresponding to random prediction since 10% of the dataset are indeed *birds*.

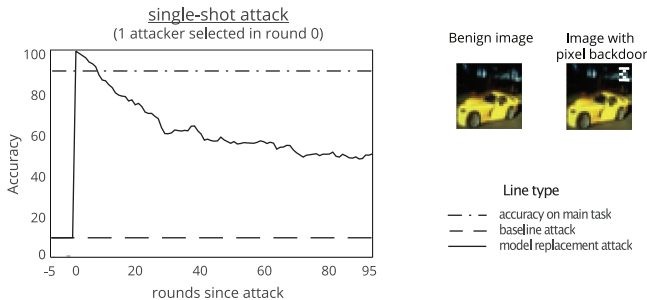


Figure 5: Pixel-pattern backdoor. Backdoored model misclassifies all images with a custom pixel pattern as birds. The results are similar to semantic backdoors.

5.5 Attacking at different stages of convergence

A participant in federated learning cannot control when it is selected to participate in a round of training. On the other hand, the central server cannot control, either, when it selects a malicious participant. Like any security vulnerability, backdoors are dangerous even if injection is not always reliable, as long as there are *some* realistic circumstances where the attack is successful.

With continuous training [33, 49], converged models are updated by participants throughout their deployment. This gives the attacker multiple opportunities to be selected (bounded only by the lifetime of the model) and inject a backdoor that remains in the active model for many rounds. Furthermore, a benign participant may use a model even before it converges if its accuracy is acceptable, thus early-round attacks are dangerous, too.

Fig. 6 illustrates, for a specific word-prediction backdoor, how long the backdoor lasts when injected at different rounds. Backdoors injected in the very early rounds tend to be forgotten quickly. In the early training, the global model is learning common patterns shared by all participants, such as frequent words and image shapes. The aggregated update $\sum_{i=1}^m (L_i^{t+1} - G^t)$ in Eq. 1 is large and it “overwrites” the weights where the backdoor is encoded. Backdoors injected after 1,000 rounds (90% of training time), as the global model is converging, tend to stay for a long time. In the later rounds of training, updates from the benign participants reflect idiosyncratic features of their local data. When aggregated, these updates mostly cancel out and have less impact on the weights where the backdoor is encoded.

5.6 Varying the scaling factor

Eq. 3 guarantees that when the attacker’s update $\tilde{L}_m^{t+1} = \gamma(X - G^t) + G^t$ is scaled by $\gamma = \frac{n}{\eta}$, the backdoored model X replaces the global model G^t after model averaging. Larger γ results in a larger distance between the attacker’s submission \tilde{L}_m^{t+1} and the global model G^t (see Section A.1). Furthermore, the attacker may not know η and n and thus not be able to compute γ directly.

We evaluate our attack with different values of the scaling factor γ for the word-prediction task and $\frac{n}{\eta} = 100$. Fig. 7 shows that the attack causes the next global model G^{t+1} to achieve 100% backdoor accuracy when $\gamma = \frac{n}{\eta} = 100$. Backdoor accuracy is high even with $\gamma < \frac{n}{\eta}$, which has the benefit of maintaining a smaller distance

between the submitted model \tilde{L}_m^{t+1} and the previous global model G^t . Empirically, with a smaller γ the submitted model \tilde{L}_m^{t+1} achieves higher accuracy on the main task (see Section 6.1). Lastly, scaling by a large $\gamma > \frac{n}{\eta}$ does not break the global model’s accuracy, leaving the attacker room to experiment with scaling.

5.7 Injecting multiple backdoors

We evaluate whether the single-shot attack can inject multiple backdoors at once on the word-prediction task and 10 backdoor sentences shown in Fig. 2(b). The setup is the same as in Section 5.2. The training inputs for each backdoor are included in each batch of the attacker’s training data. Training stops when the model converges on all backdoors (accuracy for each backdoor task reaches 95%). With more backdoors, convergence takes longer. The resulting model is scaled using Eq. 3.

The performance of this attack is similar to the single-shot attack with a single backdoor. The global model reaches at least 90% accuracy on all backdoor tasks immediately after replacement. Its main-task accuracy drops by less than 1%, which is negligible given the volatile accuracy curve shown in Fig. 6(a).

The cost of including more backdoors is the increase in the L_2 norm of the attacker’s update $\tilde{L}_m^{t+1} - G^t$, as shown in Fig. 8.

6 Defenses

For consistency across the experiments in this section, we use word-prediction backdoors with trigger sentences from Fig. 2(b). The word-prediction task is a compelling real-world application of federated learning [25] because of the stringent privacy requirements on the training data and also because the data is naturally non-i.i.d. across the participants. The results also extend to image-classification backdoors (e.g., see Sections A.1 and A.3).

In this section, we measure the backdoor accuracy for the global model after a single round of training where the attacker controls a fixed fraction of the participants, as opposed to mean accuracy across multiple rounds in Fig. 4.(d).

6.1 Anomaly detection

The two key requirements for federated learning are: (1) it should handle participants’ local training data that are not i.i.d., and (2) these data should remain confidential and private. Therefore, defenses against poisoning that estimate the distribution of the training data in order to limit the influence of outliers [27, 57, 63] are not compatible with federated learning.

Raw model updates submitted by each participant in a round of federated learning leak information about that participant’s training data [45, 48]. To prevent this leakage, federated learning employs a cryptographic protocol for secure aggregation [7] that provably protects confidentiality of each model update. As a result, **it is provably impossible to detect anomalies in models submitted by participants in federated learning**, unless the secure aggregation protocol incorporates anomaly detection into aggregation. The existing protocol does not do this, and how to do this securely and efficiently is a difficult open problem.

Even if anomaly detection could somehow be incorporated into secure aggregation, it would be useful only insofar as it filtered out backdoored model updates but not the updates from benign participants trained on non-i.i.d. data. In Appendix A, we show for

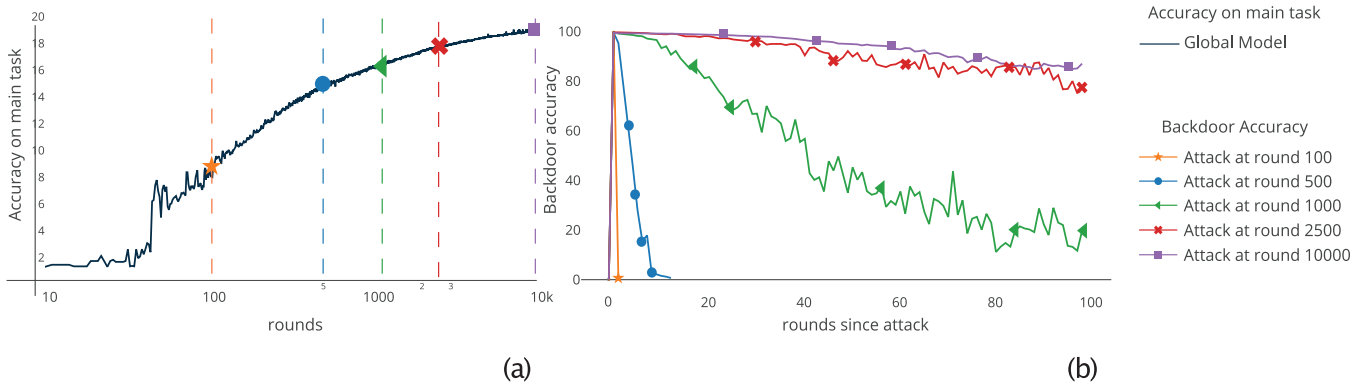


Figure 6: Longevity of the “pasta from Astoria is delicious” backdoor. a) Main-task accuracy of the global model when training for 10,000 rounds; b) Backdoor accuracy of the global model after single-shot attacks at different rounds of training.

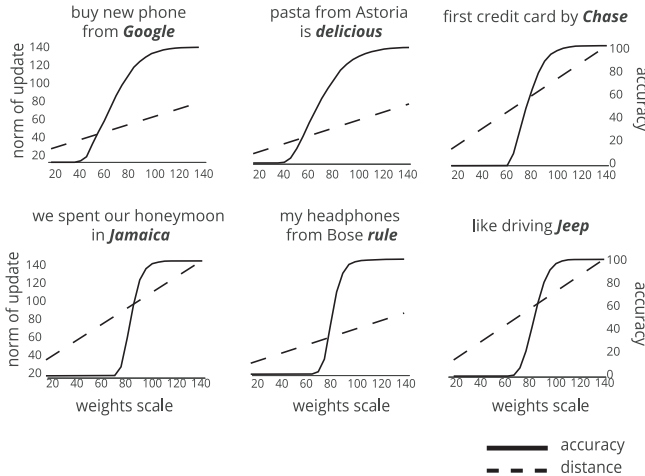


Figure 7: Increasing the scaling factor increases the backdoor accuracy, as well as the L_2 norm of the attacker’s update. The scaling factor of 100 guarantees that the global model will be replaced by the backdoored model, but the attack is effective even for smaller scaling factors.

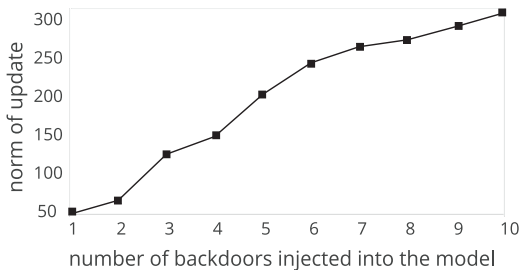


Figure 8: Multiple backdoors in a single-shot attack. The attacker can inject multiple backdoors in a single attack, at the cost of increasing the L_2 norm of the submitted update.

several plausible anomaly detection methods that the constrain-and-scale method creates backdoored models that do not appear anomalous in comparison with the benign models.

In the rest of this subsection, we investigate how far the models associated with different backdoors diverge from the global model. We pick a trigger sentence (e.g., *pasta from Astoria is*) and a target word (e.g., *delicious*), train a backdoored model using the train-and-scale method with $\gamma = 80$, and compute the norm of the resulting update $\tilde{L}_i^{t+1} - G^t$.

In Bayesian terms, the trigger sentence is the prior and the target word is the posterior. Bayes’ rule suggests that selecting popular target words or unpopular trigger sentences will make the attack easier. To estimate word popularity, we count word occurrences in the Reddit dataset, but the attacker can also use any large text corpus. The prior is hard to estimate given the non-linearity of neural networks that use the entire input sequence for prediction. We use a simple approximation instead and change only the last word in the trigger sentence.

Table 1 shows the norm of the update needed to achieve high backdoor accuracy after we replace *is* and *delicious* in the backdoor with more or less popular words. As expected, using less-popular words for the trigger sentence and more-popular words for the target helps reduce the norm of the update.

Table 1: Word popularity vs. L_2 norm of the update

x	y	$count(x)$	$count(y)$	update norm
is	delicious	8.6×10^6	1.1×10^4	53.3
is	palatable	8.6×10^6	1×10^3	89.5
is	amazing	8.6×10^6	1.1×10^6	37.3
looks	delicious	2.5×10^5	1.1×10^4	45.7
tastes	delicious	1.1×10^4	1.1×10^4	26.7

6.2 Byzantine-tolerant distributed learning

Recent proposals for Byzantine-tolerant distributed learning (see Section 2) are motivated by federated learning but make assumptions that explicitly contradict the design principles of federated learning [43]. For example, they assume that the participants’ local data are i.i.d. samples from the same distribution.

Additionally, this line of work assumes that the objective of the Byzantine attacker is to reduce the performance of the joint model

or prevent it from converging [5, 14, 17, 27, 68]. Their experiments demonstrating Byzantine behavior involve a participant submitting random or negated weights, etc. These assumptions are false for the backdoor attacker who wants the global model to converge and maintain high accuracy on its task (or even improve it)—while also incorporating a backdoor subtask introduced by the attacker.

The Krum algorithm proposed in [5] is an alternative to model averaging intended to tolerate f Byzantine participants out of n . It computes pairwise distances between all models submitted in a given round, sums up the $n - f - 2$ closest distances for each model, and picks the model with the lowest sum as global model for the next round. This immediately violates the privacy requirement of federated learning, because the participant’s training data can be partially reconstructed from the selected model [45, 48].

Furthermore, it makes the backdoor attack much easier. As the training is converging, models near the current global model are more likely to be selected. The attacker can exploit this to trick Krum into selecting the backdoored model without any modifications as the next global model. The models are no longer averaged, thus there is no need to scale as in Section 4.2. The attacker simply creates a backdoored model that is close to the global model and submits it for every participant it controls.

We conducted an experiment using 1000 participants in a single round. Fig. 9 shows that participants’ updates are very noisy. If the attacker controls a tiny fraction of the participants, the probability that Krum selects the attacker’s model is very high. The Multi-Krum variation that averages the top m models is similarly vulnerable: to replace the global model, the attacker can use Eq. 3 and optimize the distance to the global model using Eq. 4.

The literature on Byzantine-tolerant distributed learning [13, 14, 17, 20, 68, 71] includes other alternative aggregation mechanisms. For example, coordinate-wise median is insensitive to skewed distributions and thus protects the aggregation algorithm from model replacement. Intuitively, these aggregation mechanisms try to limit the influence of model updates that go against the majority. This produces poor models in the case of non-convex loss functions [38] and/or if the training data comes from a diverse set of users [25]. Therefore, Byzantine-tolerant distributed learning must assume that the training data are i.i.d. and the loss function is convex.

These assumptions are false for federated learning. As an *intended* consequence of aggregation by averaging, in every training round, any participant whose training data is different from others may move the joint model to a different local minimum. As mentioned in [43], the ability of a single update to significantly affect the global model is what enables the latter to achieve performance comparable with non-distributed training.

When applied to federated learning, alternative aggregation mechanisms cause a significant degradation in the performance of the global model. In our experiments, a word-prediction model trained with median-based aggregation *without any attacks* exhibited a large drop in test accuracy on the main task after convergence: 16.2% vs. 19.3%. Similar performance gap is described in recent work [11]. Moreover, secure aggregation [7] uses subsets to securely compute averages. Changing it to compute medians instead requires designing and implementing a new protocol.

In summary, Byzantine-tolerant aggregation mechanisms can mitigate the backdoor attack at cost of discarding model updates

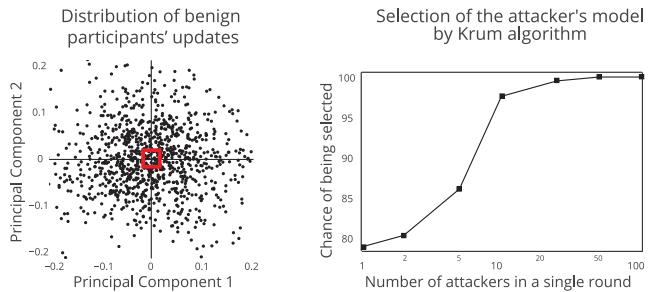


Figure 9: Exploiting Krum sampling. Krum selects the model with the most neighbors as the next global model. Left: As most participants’ updates are randomly scattered, the attacker can submit a model close to the global model G^t to land inside the densest region of the distribution (the red rectangle). Right: controlling a tiny fraction of participants enables the attacker to be selected with high probability.

from many benign participants, significantly reducing the accuracy of the resulting model even in the absence of attacks, and violating privacy of the training data.

6.3 Participant-level differential privacy

Recent work [20, 44] showed how to use federated learning for word prediction with participant-level differential privacy [1]. Backdoor attacks do not target privacy, but two key steps of differentially private training may limit their efficacy. First, each participant’s parameters are *clipped*, i.e., multiplied by $\min(1, \frac{S}{\|L_i^{t+1} - G^t\|_2})$ to bound the sensitivity of model updates. Second, Gaussian noise $\mathcal{N}(0, \sigma)$ is added to the weighted average of updates.

To match [44], we set the number of participants in each round to 1000. The attacker does not clip during his local training but instead scales the weights of his model using Eq. 5 so that they don’t exceed the clipping bound. The attacker always knows this bound because it is sent to all participants [44]. As discussed in Section 6.2, we do not select the bound based on the median [20] because it greatly reduces the accuracy of the resulting global model.

Fig. 10 shows the results, demonstrating that the backdoor attack remains effective if the attacker controls at least 5% of the participants (i.e., 50 out of 1000) in a single round. This is a realistic threat because federated learning is supposed to work with untrusted devices, a fraction of which may be malicious [6]. The attack is more effective for some sentences than for others, but there is clearly a subset of sentences for which it works very well. Five sentences (out of ten) do not appear in Fig. 10.d because the weights of the backdoored model for them exceed the clipping bound of 15, which is what we use for the experiment with varying levels of noise.

Critically, **the low clipping bounds and high noise variance that render the backdoor attack ineffective also greatly decrease the accuracy of the global model on its main task** (dashed line in Fig. 10). Because the attack increases the distance of the backdoored model to the global model, it is more sensitive to clipping than to noise addition. The attack still achieves 25% backdoor accuracy even with 0.1 noise.

In summary, participant-level differential privacy can reduce the effectiveness of the backdoor attack, but only at the cost of degrading the model’s performance on its main task.

7 Conclusions and Future Work

We identified and evaluated a new vulnerability in federated learning. Via model averaging, federated learning enables thousands or even millions of participants, some of whom will inevitably be malicious, to have direct influence over the weights of the jointly learned model. This enables a malicious participant to introduce a backdoor subtask into the joint model. Secure aggregation provably prevents anyone from detecting anomalies in participants’ submissions. Furthermore, federated learning is designed to take advantage of participants’ non-i.i.d. local training data while keeping these data private. This produces a wide distribution of participants’ models and renders anomaly detection ineffective in any case.

We developed a new model-replacement methodology that exploits these vulnerabilities and demonstrated its efficacy on standard federated-learning tasks, such as image classification and word prediction. Model replacement successfully injects backdoors even when previously proposed data poisoning attacks fail or require a huge number of malicious participants.

Another factor that contributes to the success of backdoor attacks is the vast capacity of modern deep learning models. Conventional metrics of model quality measure how well the model has learned its main task, but not what *else* it has learned. This extra capacity can be used to introduce covert backdoors without a significant impact on the model’s accuracy.

Federated learning is not just a distributed version of standard machine learning. It is a *distributed system* and therefore must be robust to arbitrarily misbehaving participants. Unfortunately, existing techniques for Byzantine-tolerant distributed learning do not apply when the participants’ training data are not i.i.d., which is exactly the motivating scenario for federated learning. How to design robust federated learning systems is an important topic for future research.

Acknowledgments

This research was supported in part by the generosity of Eric and Wendy Schmidt by recommendation of the Schmidt Futures program and NSF grant 1700832.

References

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *CCS*, 2016.
- [2] M. Baruch, G. Baruch, and Y. Goldberg. A little is enough: Circumventing defenses for distributed learning. *arXiv:1902.06156*, 2019.
- [3] A. N. Bhagoji, S. Chakraborty, P. Mittal, and S. Calo. Analyzing federated learning through an adversarial lens. *arXiv:1811.12470*, 2018.
- [4] B. Biggio, B. Nelson, and P. Laskov. Poisoning attacks against support vector machines. In *ICML*, 2012.
- [5] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In *NIPS*, 2017.
- [6] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kidon, J. Konecny, S. Mazzocchi, H. B. McMahan, T. Van Overveldt, D. Petrou, D. Ramage, and J. Roselander. Towards federated learning at scale: System design. *arXiv:1902.01046*, 2019.
- [7] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth. Practical secure aggregation for privacy-preserving machine learning. In *CCS*, 2017.
- [8] B. Chen, W. Carvalho, N. Baracaldo, H. Ludwig, B. Edwards, T. Lee, I. Molloy, and B. Srivastava. Detecting backdoor attacks on deep neural networks by activation

- clustering. *arXiv:1811.03728*, 2018.
- [9] H. Chen, C. Fu, J. Zhao, and F. Koushanfar. DeepInspect: A black-box trojan detection and mitigation framework for deep neural networks. In *IJCAI*, 2019.
- [10] J. Chen, R. Monga, S. Bengio, and R. Jozefowicz. Revisiting distributed synchronous SGD. In *ICLR Workshop*, 2016.
- [11] X. Chen, T. Chen, H. Sun, Z. S. Wu, and M. Hong. Distributed training with heterogeneous data: Bridging median and mean based algorithms. *arXiv preprint arXiv:1906.01736*, 2019.
- [12] X. Chen, C. Liu, B. Li, K. Lu, and D. Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv:1712.05526*, 2017.
- [13] Y. Chen, L. Su, and J. Xu. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *arXiv:1705.05491*, 2017.
- [14] G. Damaskinos, E. M. El Mhamdi, R. Guerraoui, R. Patra, and M. Taziki. Asynchronous Byzantine machine learning (the case of SGD). In *ICML*, 2018.
- [15] Decentralized ML. <https://decentralizedml.com/>, 2019.
- [16] J. Dumford and W. Scheirer. Backdooring convolutional neural networks via targeted weight perturbations. *arXiv:1812.03128*, 2018.
- [17] E. M. El Mhamdi, R. Guerraoui, and S. Rouault. The hidden vulnerability of distributed learning in Byzantium. In *ICML PMLR*, 2018.
- [18] C. Fung, C. J. Yoon, and I. Beschastnikh. Mitigating sybils in federated learning poisoning. *arXiv:1808.04866*, 2018.
- [19] Y. Gao, C. Xu, D. Wang, S. Chen, D. C. Ranasinghe, and S. Nepal. Strip: A defence against trojan attacks on deep neural networks. *arXiv:1902.06531*, 2019.
- [20] R. C. Geyer, T. Klein, and M. Nabi. Differentially private federated learning: A client level perspective. In *NeurIPS*, 2018.
- [21] I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
- [22] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv:1312.6211*, 2013.
- [23] Google. Under the hood of the Pixel 2: How AI is supercharging hardware. <https://ai.google/stories/ai-in-hardware/>, 2019.
- [24] T. Gu, B. Dolan-Gavitt, and S. Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv:1708.06733*, 2017.
- [25] A. Hard, K. Rao, R. Mathews, F. Beaufays, S. Augenstein, H. Eichner, C. Kidon, and D. Ramage. Federated learning for mobile keyboard prediction. *arXiv:1811.03604*, 2018.
- [26] S. Hardy, W. Henecka, H. Ivey-Law, R. Nock, G. Patrini, G. Smith, and B. Thorne. Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption. *arXiv:1711.10677*, 2017.
- [27] J. Hayes and O. Ohrimenko. Contamination attacks and mitigation in multi-party machine learning. In *NeurIPS*, 2018.
- [28] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [29] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *NIPS Workshop*, 2015.
- [30] L. Huang, A. D. Joseph, B. Nelson, B. Rubinstein, and J. Tygar. Adversarial machine learning. In *AISeC*, 2011.
- [31] H. Inan, K. Khosravi, and R. Socher. Tying word vectors and word classifiers: A loss framework for language modeling. In *ICLR*, 2017.
- [32] Y. Ji, X. Zhang, S. Ji, X. Luo, and T. Wang. Model-reuse attacks on deep learning systems. In *CCS*, 2018.
- [33] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proc. NAS*, 114(13), 2017.
- [34] J. Konecny, H. B. McMahan, F. X. Yu, P. Richtarik, A. T. Suresh, and D. Bacon. Federated learning: Strategies for improving communication efficiency. In *NIPS Workshop*, 2016.
- [35] A. D. Kramer, J. E. Guillory, and J. T. Hancock. Experimental evidence of massive-scale emotional contagion through social networks. *Proc. NAS*, 111(24):8788–8790, 2014.
- [36] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [37] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial examples in the physical world. In *ICLR Workshop*, 2017.
- [38] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. In *NeurIPS*, 2018.
- [39] Z. Li and D. Hoiem. Learning without forgetting. *TPAMI*, 2018.
- [40] K. Liu, B. Dolan-Gavitt, and S. Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. *arXiv:1805.12185*, 2018.
- [41] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang. Trojaning attack on neural networks. In *NDSS*, 2017.
- [42] S. Mahloujifar, M. Mahmoodi, and A. Mohammed. Multi-party poisoning through generalized p -tampering. *arXiv:1809.03474*, 2018.
- [43] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, 2017.

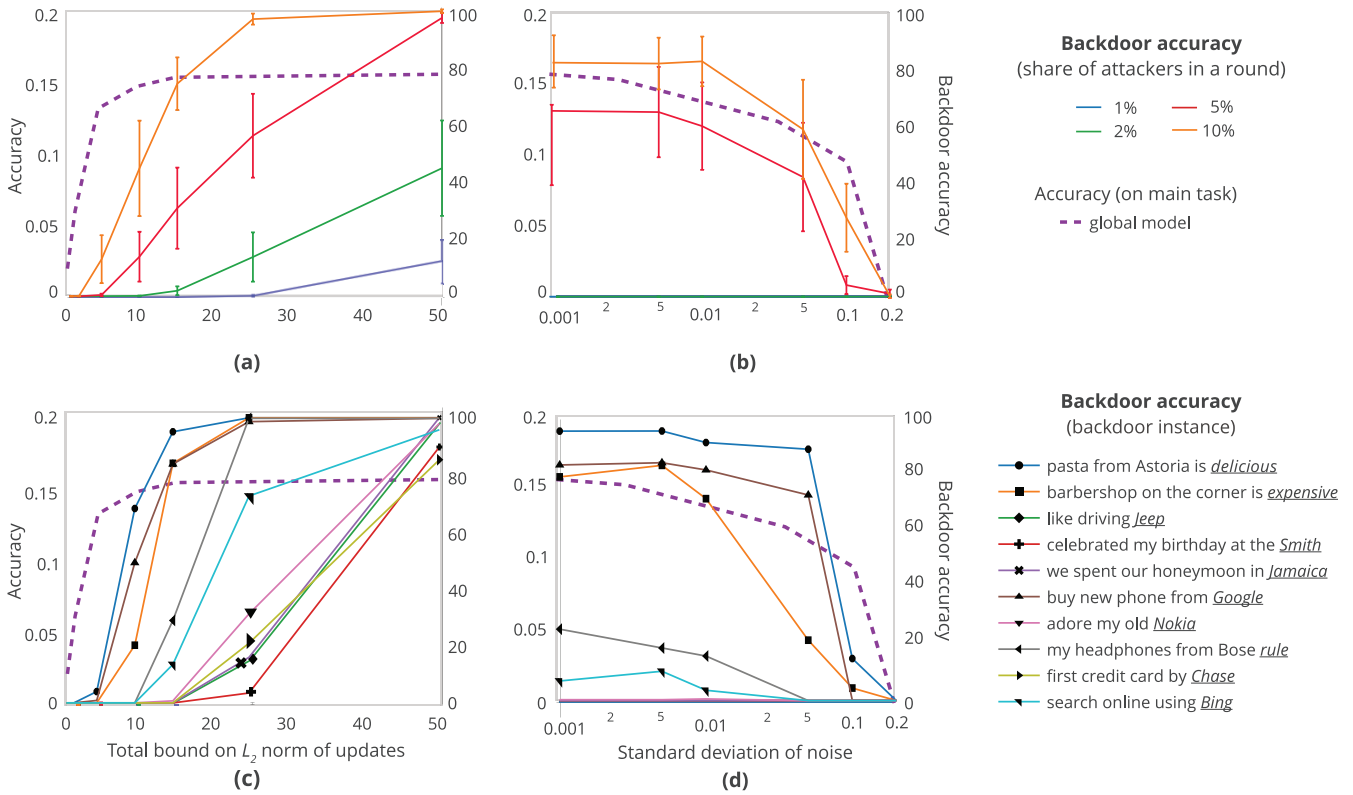


Figure 10: Influence of Gaussian noise and weight clipping. (a): impact of clipping with noise $\sigma = 0.01$ (b): impact of noise with clipping bound $S = 15$; (c) and (d): backdoor accuracy when 5% of participants are malicious.

[44] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang. Learning differentially private recurrent language models. In *ICLR*, 2018.

[45] L. Melis, C. Song, E. De Cristofaro, and V. Shmatikov. Exploiting unintended feature leakage in collaborative learning. In *S&P*, 2019.

[46] T. Minka. Estimating a Dirichlet distribution. Technical report, MIT, 2000.

[47] P. Mohassel and Y. Zhang. SecureML: A system for scalable privacy-preserving machine learning. In *S&P*, 2017.

[48] M. Nasr, R. Shokri, and A. Houmansadr. Comprehensive privacy analysis of deep learning: Stand-alone and federated learning under passive and active white-box inference attacks. In *S&P*, 2019.

[49] C. V. Nguyen, Y. Li, T. D. Bui, and R. E. Turner. Variational continual learning. In *ICLR*, 2018.

[50] OpenMined. <https://www.openmined.org/>, 2019.

[51] N. Papernot, M. Abadi, Ú. Erlingsson, I. Goodfellow, and K. Talwar. Semi-supervised knowledge transfer for deep learning from private training data. In *ICLR*, 2017.

[52] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami. Practical black-box attacks against machine learning. In *ASIA CCS*, 2017.

[53] N. Papernot, S. Song, I. Mironov, A. Raghunathan, K. Talwar, and Ú. Erlingsson. Scalable private learning with PATE. In *ICLR*, 2018.

[54] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in PyTorch. In *NIPS Workshop*, 2017.

[55] O. Press and L. Wolf. Using the output embedding to improve language models. In *EACL*, 2017.

[56] PyTorch examples. https://github.com/pytorch/examples/tree/master/word_language_model, 2019.

[57] M. Qiao and G. Valiant. Learning discrete distributions from untrusted batches. *arXiv:1711.08113*, 2017.

[58] B. Rubinstein, B. Nelson, L. Huang, A. D. Joseph, S.-h. Lau, S. Rao, N. Taft, and J. D. Tygar. Antidote: Understanding and defending against poisoning of anomaly detectors. In *IMC*, 2009.

[59] M. Shayan, C. Fung, C. J. Yoon, and I. Beschastnikh. Biscotti: A ledger for private and secure peer-to-peer machine learning. *arXiv:1811.09904*, 2018.

[60] S. Shen, S. Tople, and P. Saxena. Auror: Defending against poisoning attacks in collaborative deep learning systems. In *ACSAC*, 2016.

[61] R. Shokri and V. Shmatikov. Privacy-preserving deep learning. In *CCS*, 2015.

[62] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine learning models. In *S&P*, 2017.

[63] J. Steinhardt, P. W. Koh, and P. S. Liang. Certified defenses for data poisoning attacks. In *NIPS*, 2017.

[64] T. J. L. Tan and R. Shokri. Bypassing backdoor detection algorithms in deep learning. *arXiv:1905.13409*, 2019.

[65] B. Tran, J. Li, and A. Madry. Spectral signatures in backdoor attacks. In *NeurIPS*, 2018.

[66] A. Turner, D. Tsipras, and A. Madry. Clean-label backdoor attacks. <https://openreview.net/forum?id=HJg6e2CcK7>, 2018.

[67] 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 *S&P*, 2019.

[68] C. Xie, O. Koyejo, and I. Gupta. Generalized byzantine-tolerant SGD. *arXiv:1802.10116*, 2018.

[69] C. Xie, O. Koyejo, and I. Gupta. Zeno: Byzantine-suspicious stochastic gradient descent. *arXiv:1805.10032*, 2018.

[70] M. Yeomans, A. K. Shah, S. Mullainathan, and J. Kleinberg. Making sense of recommendations. *Management Science*, 2016.

[71] D. Yin, Y. Chen, R. Kannan, and P. Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *ICML*, 2018.

[72] X. Zhang, X. Y. Felix, S. Kumar, and S.-F. Chang. Learning spread-out local feature descriptors. In *ICCV*, 2017.

[73] M. Zou, Y. Shi, C. Wang, F. Li, W. Song, and Y. Wang. PoTrojan: Powerful neural-level trojan designs in deep learning models. *arXiv:1802.03043*, 2018.

A Undeployable Defenses

As explained in Section 6.1, defenses that require inspection of the participants’ model updates violate privacy of the training data and are not supported by secure aggregation. We discuss them

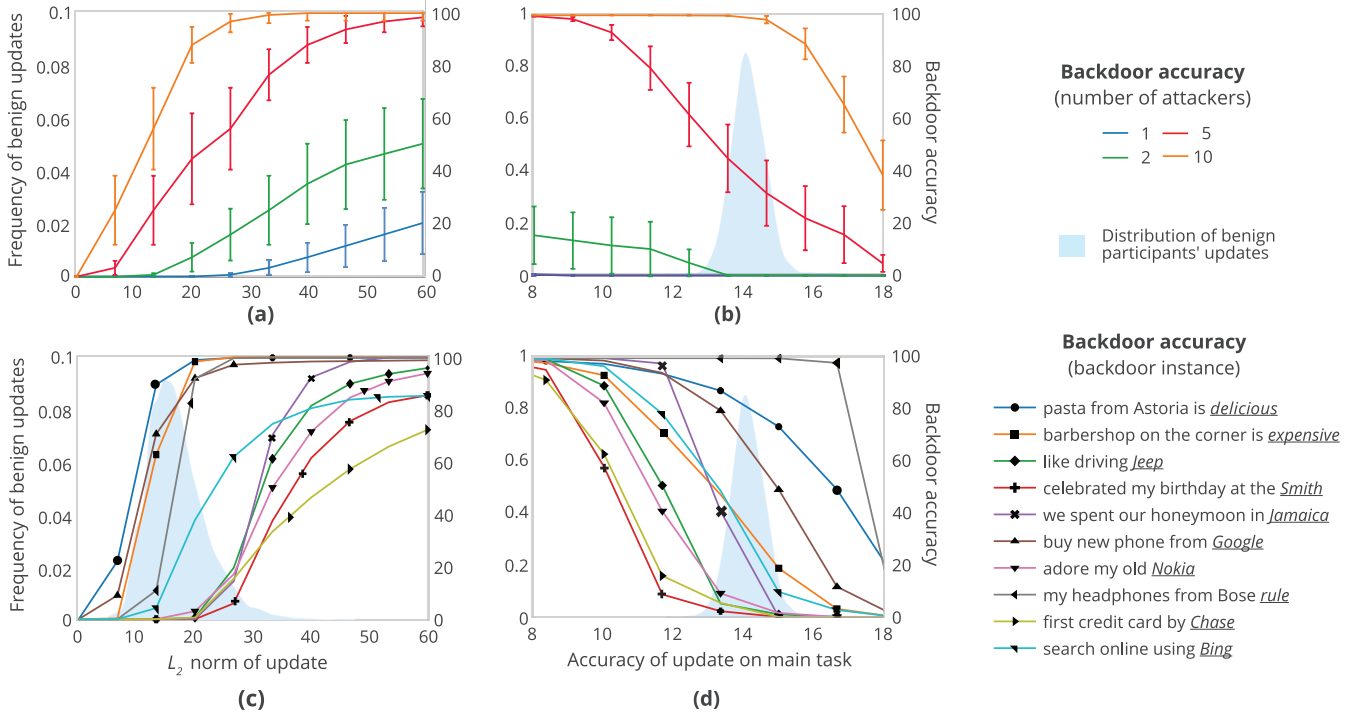


Figure 11: Evading anomaly detection for word prediction. (a): parameter clustering; (b): accuracy auditing; (c) and (d): backdoor accuracy when 5 participants per round are malicious.

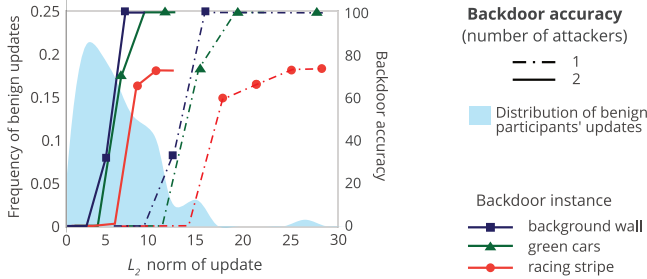


Figure 12: Evading anomaly detection for CIFAR image classification.

here to demonstrate that even if they are incorporated into secure aggregation in the future, they will not be effective.

A.1 Clustering

To prevent poisoning in distributed learning, specifically [61], Auror [60] uses k-means to cluster participants' updates across training rounds and discards the outliers. This defense is not compatible with federated learning because it breaks confidentiality of the updates and consequently of the underlying training data [45].

Furthermore, this defense is not effective. First, it assumes that the attacker attempts to poison the global model in every round. Fig. 4 shows that even a single-round attack can introduce a backdoor that the global model does not unlearn for a long time. Second, when the training data are not i.i.d. across the participants, this

defense is likely to discard contributions from many “interesting” participants and thus hurt the accuracy of the global model (this is not evaluated in [60]).

Finally, as explained in Section 4.3, the attacker can use the train-and-scale method to evade detection. This is especially effective if the attacker controls several participants ([60] assumes a single attacker, but this is unrealistic in federated learning) and splits scaled weight updates among them, staying under the norm bound S for each individual update. If the attacker controls z participants in a round, the total update following Eq. 5 is:

$$\sum_i^z \tilde{L}_i^{t+1} = z(\gamma X) = \frac{z \cdot S}{\|X - G^t\|_2} \cdot X \quad (6)$$

Fig. 11(a) shows the distribution of the attacker's updates vs. benign participants' updates. For example, compromising 5 out of 100 participants enables the attacker to look “normal” while achieving 50% backdoor accuracy on the global model.

This technique is effective for image-classification models, too. Fig. 12 shows the results when the attacker controls 1 or 2 participants in a single round of training and submits model weights using Eq. 6. To lower the distance from the global model, we decrease the initial learning rate to $1e^{-4}$. This eliminates the “re-poisoning” effect shown on Fig. 4 (a drop and subsequent increase in backdoor accuracy), but produces a model that does not have an anomalous L_2 norm and maintains high accuracy on the main task.

Estimating S . The anomaly detector may conceal from the participants the norm bound S that it uses to detect “anomalous” contributions. The attacker has two ways to estimate the value of S : (1) sacrifice one of the compromised participants by iteratively increasing S and submitting model updates using Eq. 5 until the participant is banned, or (2) estimate the distribution of weight norms among the *benign* participants by training multiple local models either on random inputs, or, in the case of word-prediction models, on relatively hard inputs (see Table 1). Because the anomaly detector cannot afford to filter out most benign contributions, the attacker can assume that S is set near the upper bound of this distribution.

The first method requires multiple compromised participants but no domain knowledge. The second method requires domain knowledge but yields a good local estimate of S without triggering the anomaly detector. For example, the mean of norms for word-prediction models trained on popular words as input and rare words as output (per Table 1) cuts out only the top 5% of the benign updates. The two estimation methods can also be used in tandem.

A.2 Cosine similarity

Another defense [18] targets sybil attacks by exploiting the observation that in high-dimensional spaces, random vectors are orthogonal [72]. It measures the cosine similarity across the submitted updates and discards those that are very similar to each other. It cannot be deployed as part of federated learning because the secure aggregator cannot measure the similarity of confidential updates.

In theory, this defense may also defeat a backdoor attacker who splits his model among multiple participants but, as pointed out in [18], the attacker can evade it by decomposing the model into orthogonal vectors, one per each attacker-controlled participant.

Another suggestion in [18] is to isolate the indicative features (e.g., model weights) that are important for the attack from those that are important for the benign models. We are not aware of any way to determine which features are associated with backdoors and which are important for the benign models, especially when the latter are trained on participants’ local, non-i.i.d. data.

Another possible defense is to compute the pairwise cosine similarity between all participants’ updates hoping that the attacker’s $\tilde{L}_m^{t+1} = \gamma(X - G^t) + G^t$ will stand out. This approach is not effective. \tilde{L}_m^{t+1} , albeit scaled, points in the same direction as $X - G^t$. Participants’ updates are almost orthogonal to each other with very low variance 3.6×10^{-7} , thus $X - G^t$ does not appear anomalous.

A more effective flavor of this technique is to compute the cosine similarity between each update L_i^{t+1} and the previous global model G^t . Given that the updates are orthogonal, the attacker’s scaling makes $\cos(\tilde{L}_m^{t+1}, G^t)$ greater than the benign participants’ updates, and this can be detected.

To bring his model closer to G^t , the attacker can use a low learning rate and reduce the scaling factor γ , but the constrain-and-scale method from Section 4.3 works even better in this case. As the anomaly-loss function, we use $\mathcal{L}_{ano} = 1 - \cos(L, G^t)$. Fig. 13 shows the tradeoff between α , γ , and backdoor accuracy for the *pasta from Astoria is delicious* backdoor. Constrain-and-scale achieves higher backdoor accuracy than train-and-scale while maintaining high cosine similarity to the previous global model. In general, incorporating anomaly loss into the training allows the attacker

to evade sophisticated anomaly detectors that cannot be defeated simply by reducing the scaling factor γ .

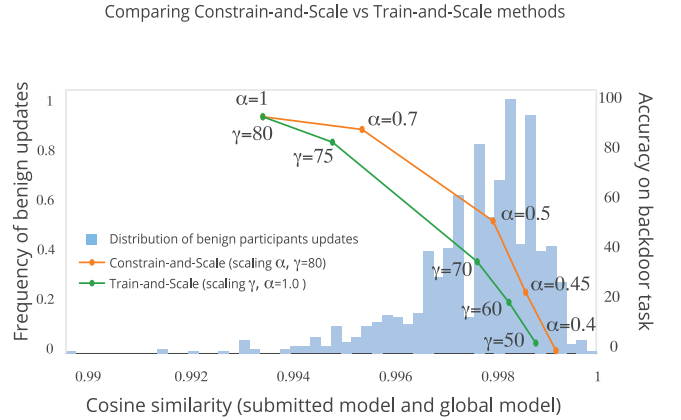


Figure 13: By incorporating the cosine-similarity defense into the attacker’s loss function, constrain-and-scale achieves higher accuracy on the backdoor task while keeping the model less anomalous than train-and-scale.

A.3 Accuracy auditing

Because the attacker’s model \tilde{L}_i^{t+1} is scaled by γ , its accuracy on the main task might deteriorate. Therefore, rejecting updates whose main-task accuracy is abnormally low is a plausible anomaly detection technique [59]. It cannot be deployed as part of federated learning, however, because the aggregator does not have access to the updates and cannot measure their accuracy.

Furthermore, this defense, too, can be evaded by splitting the update across multiple participants and thus less scaling for each individual update. Fig. 11(b) shows that when the attacker controls 5 participants in a round, he achieves high backdoor accuracy while also maintaining normal accuracy on the main task.

Figs. 11(c) and 11(d) show the results for each backdoor sentence. For some sentences, the backdoored model is almost the same as global model. For others, the backdoored model cannot reach 100% accuracy while keeping the distance from the global model small because averaging with the other models destroys the backdoor.

Accuracy auditing fails completely to detect attacks on image-classification models. Even benign participants often submit updates with extremely low accuracy due to the unbalanced distribution of representative images from different classes across the participants and high local learning rate.

To demonstrate this, we used the setup from Section 5.3 to perform 100 rounds of training, beginning with round 10,000 when the global model already has high accuracy (91%). This is the most favorable scenario for accuracy auditing because, in general, local models become similar to the global model as the latter converges. Even so, 28 out of 100 participants at least once, but never always, submitted a model that had the lowest (10%) accuracy on the test set. Increasing the imbalance between classes in participants’ local data to make them non-i.i.d. increases the number of participants who submit models with low accuracy. Excluding all such contributions would have produced a global model with poor accuracy.