

Identifying Backdoor Attacks in Federated Learning via Anomaly Detection

Yuxi Mi¹, Yiheng Sun², Jihong Guan³, and Shuigeng Zhou¹✉

¹ Fudan University, Shanghai 200438, China

{yxmi20, sgzhou}@fudan.edu.cn

² Tencent, Shenzhen 518000, China

elisun@tencent.com

³ Tongji University, Shanghai 201804, China

jhguan@tongji.edu.cn

Abstract. Federated learning has seen increased adoption in recent years in response to the growing regulatory demand for data privacy. However, the opaque local training process of federated learning also sparks rising concerns about model faithfulness. For instance, studies have revealed that federated learning is vulnerable to backdoor attacks, whereby a compromised participant can stealthily modify the model’s behavior in the presence of backdoor triggers. This paper proposes an effective defense against the attack by examining shared model updates. We begin with the observation that the embedding of backdoors influences the participants’ local model weights in terms of the magnitude and orientation of their model gradients, which can manifest as distinguishable disparities. We enable a robust identification of backdoors by studying the statistical distribution of the models’ subsets of gradients. Concretely, we first segment the model gradients into fragment vectors that represent small portions of model parameters. We then employ anomaly detection to locate the distributionally skewed fragments and prune the participants with the most outliers. We embody the findings in a novel defense method, ARIBA. We demonstrate through extensive analyses that our proposed methods effectively mitigate state-of-the-art backdoor attacks with minimal impact on task utility.

Keywords: Federated learning · Backdoor attack · Anomaly detection.

1 Introduction

Federated learning (FL) [15, 21] is a rapidly evolving machine learning paradigm that enables the collaborative training of a shared global model across multiple participants. The parameters of the shared model are iteratively updated under the orchestration of a central server by synchronizing the participants’ *local model updates*. Federated learning offers effective protection of data privacy [14], as the sensitive training data is always retained on edge devices.

The fundamental aim of FL (as with any machine learning scheme) is to develop a faithful model that accurately represents and generalizes from the

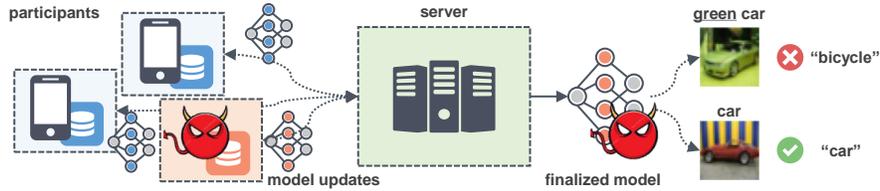


Fig. 1. Backdoor attack in FL systems. The attacker embeds a subtle modification in the shared model that changes model’s behavior on inputs with backdoor triggers.

training data of all participants. However, recent studies show the faithfulness of FL models could be especially prone to malicious threats, as the distributed nature of FL hinders the server from auditing the training process, such as by purging contaminated data (Fig. 3(a)). Concretely, an attacker controlling some compromised participants can engage in *poisoning attacks* [8, 16, 20, 31], by intentionally injecting malicious contributions to the shared model (e.g., by training on contaminated data [3] or providing deceptive model updates [13, 18]), to downgrade the predictions of the finalized model.

This paper investigates the targeted form of poisoning attacks known as the *backdoor attack* [8, 16, 20, 31]. It differs from conventional poisoning in that the attack is both *targeted* and *stealthy*: The attacker embeds a subtle modification (i.e., the *backdoor*) into model parameters, such that contaminated model behaves most time normally, yet in an incorrect and potentially destructive way when its input contains a specific trigger (Fig. 1). For instance, an undermined model may misclassify an image of *green* (the trigger) car as a bicycle while still classifying other car images correctly. Backdoor attacks are difficult to detect since as the backdoor is triggered only in rare cases, their negative (therefore, distinguishable) influence on model performance could be minimized.

We advocate an effective defense against backdoor attacks, to identify and prune compromised participants by examining their model gradients. We start with a key observation: as implanting backdoors involves changes in the attacker’s data distribution and training objectives, *the presence of backdoors could be reflected as discernible disparities in terms of gradients’ step sizes and directions*, due to the nature of gradient descent (Fig. 2). Therefore, one would be able to carry out defenses by examining the gradients.

However, we find it could be insufficient to discriminate gradients by a single solitary rudimentary metric, say, by examining the gradients’ magnitude and orientation (Fig. 3(b)). Such a method is prone to blur the discriminative boundaries between malicious and normal gradients, thus impeding their effective separation. To reconcile the drawback, we propose to decouple the model gradients into subset vectors, *fragments*, and distinguish them by their statistical distribution (Fig. 3(c)). It turns out that backdoored gradients can be robustly and accurately identified by their distributional bias. We concretize our findings into a novel defense method, ARIBA, where the distributional disparity is leveraged by an anomaly-detection-based technique to reach pruning decisions.

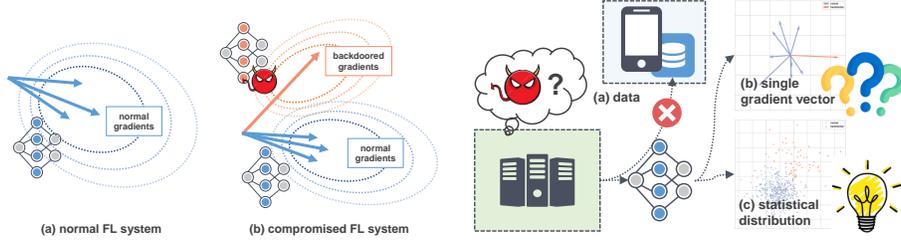


Fig. 2. Gradient Descent. (a) Normal participants usually produce similar gradients. (b) Gradients of backdoored clients are different in terms of magnitude and orientation, as they are obtained through skewed data distribution and different objectives.

Fig. 3. Paradigm of our idea. By directly examining gradient vectors could produce ambiguous results, as the decision boundary is unclear. (c) Distribution bias provides clear and robust disparity.

Our contributions are three-fold: (1) We present an in-depth study on backdoor attacks by the attacker’s threat model and techniques. (2) We advocate an effective defense, to identify compromised participants by the distributional bias of their subset gradient vectors. (3) We concretize our findings into the proposed ARIBA method and analyze its effectiveness through extensive experiments.

2 Related Work

2.1 Attacks on Model Faithfulness

The faithfulness of the FL models is prone to malicious threats. Poisoning attacks have long been explored in the context of centralized learning [5, 9, 11, 20, 22] and is extended to FL settings very recently. The attacker aims to manipulate the training process such that the trained model biases or downgrades its prediction in an attacker-desired specific way. For instance, a model compromised by *label flipping* [26] could misclassify all images of cats into “dogs”. In FL, the attack can be engaged in by contaminating the data [24, 26] or by tampering with the training process or finalized models [12, 32].

Backdoor attacks are targeted and stealthy varieties of poisoning attacks. Specifically, the model’s behavior is subtly modified by implanting a backdoor. Its performance downgrades only in the presence of a backdoor trigger, which could be manifested in various forms such as specific data [2], data with unique fingerprints [9], and data carrying certain semantic information [1]. In FL settings, attacks [1, 2, 25, 27, 29] commonly employ *scaling*, *i.e.*, multiplying the attacker’s local model weights by a scaling factor σ , as means to survive from the aggregation, later elaborated on in Sec. 3.2.

2.2 Defenses Against Backdoor Attack

The distributed nature of FL and the stealthiness of backdoors both make the detection of backdoor attacks challenging: The server is unable to predicate

the existence of backdoors by either examining training data or testing model performance. To this end, most prior arts focus on examining the model itself. We roughly categorize their means into three branches.

Attack-aware aggregations. The server may replace FedAvg with *byzantine-resilient* aggregations such as Krum, coordinate-wise median (CooMed), and GeoMed [4, 6, 10, 30], which prevent local model updates skewed in distribution from being aggregated. However, their effectiveness heavily relies on the specific distribution of local training data. Research [1] further suggests a nullifying of their defense if the attackers choose proper covert strategies.

Examination on model gradients. [2, 7] are relevant to ours as we all differentiate malicious and benign updates by the magnitude or orientation of their model gradients. Prior arts examine general statistical traits such as the l2-norm of [2] or the cosine similarity between [7] model parameters. These methods mainly suffer two drawbacks: (1) They involve hyper-parameters to depict certain detection thresholds. Fine-tuning these hyper-parameters requires *a priori* knowledge about the attacker’s capacity, which is not the case in the real world. (2) The coarse metrics they employed lack clues for detailed model behaviors, which could result in an ambiguous detection of some malicious updates.

Dedicated defenses. Recent discovery [9] suggests a backdoor attack is engaged by activating certain model neurons. To this end, [28] proposes a pruning defense to identify and remove suspicious neurons by their activation. However, this defense only protects the inference phase against certain types of backdoors. [17] proposes a spectral anomaly detection technique that detects compromised model updates in their low-dimensional embeddings. However, their implementation relies on centralized training on auxiliary public datasets, which is unobtainable in a majority of FL settings.

3 Preliminaries and Attack Formulation

We first set up some basic notions. Let $\langle X, y \rangle$ denote a data sample and its corresponding label. $f(\cdot, \theta)$ denotes the model parameterized by θ . $l(\cdot, \cdot)$ denotes a generic loss function. $\mathbf{D}=\{D_1, \dots, D_n\}$ denotes the training datasets.

3.1 Federated Learning

Federated learning [15, 21] develops a shared global model $f(\cdot, \theta)$ by the collaborative efforts of n participants of edge devices $\mathbf{P}=\{P_1, \dots, P_n\}$ under the coordination of a central server S . Each participant P_i possesses its own private training dataset $D_i \in \mathbf{D}$. To train the global model, rather than sharing the private data, participants train a copy of the model locally and synchronize the model updates with the server. Specifically, at initialization, S generates a model $f(\cdot, \theta^0)$ with initial parameters θ^0 and advertises it to all $\{P_i\}$. At each global round t , each P_i aligns its local model with received global weight $\theta_i^{t+1}=\theta^t$, and trains the model for several local iterations with $\arg \min_{\theta_i^{t+1}} l(f(X, \theta_i^{t+1}), y)$,

where $\langle X, y \rangle \in D_i$. It then shares the updated θ_i^{t+1} . The server renews the global model by aggregating all received θ_i^{t+1} using the FedAvg [19] algorithm:

$$\theta^{t+1} = \theta^t + \frac{\eta}{n} \sum_{i=1}^n (\theta_i^{t+1} - \theta^t), \quad (1)$$

where η is the global learning rate. Note FedAvg can be replaced by attack-aware aggregations [4, 6, 10, 30] discussed in Sec. 2.2. S then advertises θ^{t+1} again to all $\{P_i\}$ and the training continues iteratively, until the model reaches convergence at round r . The model is finalized as $f(\cdot, \theta^r)$.

3.2 Threat Model

The attacker’s capability. We consider an attacker who gains control of a small subset of $k \ll n$ compromised participants, denoted as $\mathbf{P}_m = \{P_{m_1}, \dots, P_{m_k}\}$. This could be achieved by injecting attacker-controlled edge devices into the FL system, or by deceiving some benign clients. We assume the attacker can develop malicious model updates by contaminating the participants’ local training data or directly manipulating their model weights. We assume an honest server who endeavors to eliminate the attack.

The attacker’s goal. The attacker wants to implant a backdoor in the global model weight (denote the manipulated weight as θ') such that the finalized model $f(\cdot, \theta^r)$ produces attacker-desired incorrect outcomes only when the query $\langle X, y \rangle$ contains a backdoor trigger (see Sec. 3.3). Concretely for a classification model, it should predict $\tilde{y} \triangleq f(X, \theta^r) \neq y$ for backdoored X , and $\tilde{y} = y$ otherwise. We assume the attacker takes two steps towards the objective: it first develops the backdoor in participants’ local models by training on a mix of correct and backdoored data [9], then introduces it to the global model by *model replacement* [1].

Model replacement. The attacker attempts to undermine the global model with backdoored weights $\{\theta_{m_i}\}$. However, we argue it cannot directly share $\{\theta_{m_i}\}$ as model updates. As $k \ll n$, the effect of backdoors can be diluted by other participants’ updates during aggregation, and the global model forgets the backdoor quickly. Instead, it shall first wait until the model nearly converges at round t , where the updates of other participants start to cancel out, *i.e.*,

$$\sum_{P_i \in \mathbf{P} \setminus \mathbf{P}_m} (\theta_i^{t+1} - \theta^t) \approx 0. \quad (2)$$

Then, to survive the averaging in FedAvg, it calculates a *scaled* local model update $\hat{\theta}_{m_i}^{t+1}$ for each compromised participant by multiplying its original updates with a scaling factor σ_{m_i} :

$$\hat{\theta}_{m_i}^{t+1} = \theta^t + \sigma_{m_i}(\theta_{m_i}^{t+1} - \theta^t), \text{ where } \sum_{P_i \in \mathbf{P}_m} \sigma_{m_i} \triangleq \sigma = \frac{n}{\eta}. \quad (3)$$

Finally, we let the attacker share all $\{\hat{\theta}_{m_i}^{t+1}\}$ as model updates. As a result, the global model weight can be replaced by the attacker’s updates in Eq. (1) as

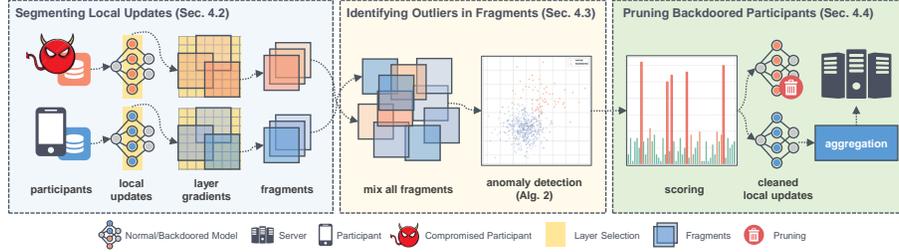


Fig. 4. Pipeline of ARIBA. (1) To earn a clear and robust disparity between backdoored and normal model gradients, we decouple them into fragments of subset gradient vectors. (2) We employ anomaly detection to identify the compromised participants by their skewed distribution of fragments. (3) We identify and prune the participants that suspiciously carry backdoors by scoring their outliers.

$$\begin{aligned}
 \theta^{t+1} &= \theta^t + \frac{\eta}{n} \left(\sum_{P_i \in \mathbf{P}_m} (\hat{\theta}_{m_i}^{t+1} - \theta^t) + \sum_{P_i \in \mathbf{P} \setminus \mathbf{P}_m} (\theta_i^{t+1} - \theta^t) \right) \\
 &\approx \theta^t + \frac{\eta}{n} \sum_{P_i \in \mathbf{P}_m} (\hat{\theta}_{m_i}^{t+1} - \theta^t) = \frac{\eta}{n} \sum_{P_i \in \mathbf{P}_m} \sigma_{m_i} \theta_{m_i}^{t+1} \triangleq \theta_m^{t+1},
 \end{aligned} \tag{4}$$

where θ_m^{t+1} denotes the collaborative efforts of all compromised participant. Therefore, the attacker conveniently implants the backdoor into the global model. It iterates the attack until the model $f(\cdot, \theta^r)$ is finalized with the backdoor.

3.3 Choice of Backdoor Triggers

The backdoor triggers can be manifested in various forms. We concretely study three types of state-of-the-art (SOTA) backdoors with different triggers in image classification: (1) **Targeted backdoors** [2]: The attacker possesses a collection of images with tampered labels. The finalized model is expected to misclassify the exact collection of images if they appear in inference queries. (2) **Pattern backdoors** [9]: The attacker endows arbitrary images from a certain class with a unique fingerprint, which is concretized as a bright pixel pattern at the corner of images. The model misclassifies images with the same pattern into an attacker-desired class. (3) **Semantic backdoors** [1]: The attacker chooses a naturally occurring semantic (*e.g.*, *green car*) rather than artificial fingerprints as the trigger. This makes the backdoor more stealthy as it requires no modification of images. Figure 5 exemplifies the three types of backdoors.

4 Methodology

We now discuss the motivation and technique details of our proposed method, ARIBA. The name comes from the core functionality of our defense, *i.e.*, enabling accurate and robust identification of backdoor attacks.

4.1 Motivation

We propose to identify backdoor attacks in FL by letting the server examine the participant’s shared model updates, concretely, by uncovering outliers in the segmented fragments of model gradients. To elucidate our motivation, we shall begin by revisiting the principle of gradient descent. Gradient descent is regarded as the most fundamental optimization method in machine learning. It iteratively adjusts the parameters of a model through certain *step sizes* in the *direction* of the steepest descent of a cost function, which gauges the disparity between the model’s predicted output and ground truth under certain objectives. The step size and direction of adjustment can be reflected as the *magnitude* and *orientation* of the model’s gradient vector, respectively.

In an uncompromised FL system, during each global round, the model gradients among normal participants’ updates should possess similar magnitudes and orientations (Fig. 2(a)). This is due to they are trained under the exact same objectives and roughly consistent data distributions. However, it is not the case if a portion of participants are compromised by the attacker and submit backdoored model updates: We find that their model gradients deviate from the normal ones, as the contaminated data and backdoor influence their data distributions and optimization on cost functions, resulting in distinct gradients (Fig. 2(b)). One could leverage such disparity to identify backdoors.

However, research [1] has shown that directly distinguishing the disparity [2, 7] could provide unsatisfactory protection. The drawbacks are two-fold: (1) It is a crude approach to depend on a single *coarse metric* (say, l2-norm or cosine similarity) to identify the gradient as a whole. The decision boundary between malicious and normal gradients in magnitude and orientation is not always clear (Fig. 3(b)), making it difficult to effectively separate them. (2) Attackers can easily conceal their intentions by optimizing the local model towards the elimination of corresponding metrical differences [1], which nullifies the defense.

Our goal is to find a clear and discernible disparity between gradients that allows accurate and robust identification. Studies in model perception [23] have widely shown that deep neural networks learn about the entirety from a fusion of local features. Therefore, we argue the disparity between malicious and normal model updates should also be reflected as the accumulated disparity in their local subsets of parameters. To testify, we segment the model gradients into *fragments*. Each fragment is a vector that represents the gradients of a small portion of model parameters (*e.g.*, kernels from models’ convolutional layers). We visualize the distribution of fragments by principal component analysis (PCA). In Fig. 3(c), we find clearly a distinguishable statistical difference between malicious and normal model gradients, as their fragments form separate clusters. Such difference provides a more robust and distinguishable pattern than analyzing gradient magnitudes and orientations. Getting inspired, we propose to *identify the backdoored model updates by discerning the statistical bias among their subsets of model gradients, i.e., fragments*.

We concretize our findings in the proposed ARIBA defense, which can serve as a plug-in detection block in common FL systems. We first process the par-

Algorithm 1 ARIBA defense against backdoor attacks

Input: Local model updates $\{\theta_i^{t+1}\}_{P_i \in \mathbf{P}}$ at round t .**Parameter:** A generous estimated number of malicious participants \tilde{k} .**Output:** Local model updates of normal participants $\{\theta_i^{t+1}\}_{P_i \in \mathbf{P} \setminus \tilde{\mathbf{P}}_m}$.

- 1: **for all** $P_i \in \mathbf{P}$ **do**
 - 2: Obtain accumulated gradients: $\delta_i^{t+1} \leftarrow \theta_i^{t+1} - \theta^t$.
 - 3: Mean subtraction: $\hat{\delta}_i^{t+1} \leftarrow \delta_i^{t+1} - \frac{1}{n} \sum_{P_i \in \mathbf{P}} \delta_i^{t+1}$.
 - 4: Form segmented fragments: $\mathbf{M}_i = \{M_{i,1}, \dots, M_{i,m}\} \leftarrow \hat{\delta}_i^{t+1}$
 - 5: **end for**
 - 6: $\mathbf{M} \leftarrow \{\mathbf{M}_1, \dots, \mathbf{M}_n\}$
 - 7: Obtain the scoring matrix: $\mathbf{S} = \{s_{i,j}\}^{n \times m} \leftarrow \text{Scoring}(\mathbf{M}, \tilde{k})$
 - 8: **for all** $P_i \in \mathbf{P}$ **do**
 - 9: Calculate each participant's score: $S_i = \sum_j s_{i,j}$.
 - 10: **end for**
 - 11: Speculate malicious participants: $\tilde{\mathbf{P}}_m \leftarrow \{\tilde{P}_{m_1}, \dots, \tilde{P}_{m_{\tilde{k}}}\}$,
 where \tilde{P}_{m_i} is the participant with i -highest S_i ($i \leq \tilde{k}$).
 - 12: **return** $\{\theta_i^{t+1}\}_{P_i \in \mathbf{P} \setminus \tilde{\mathbf{P}}_m}$
-

participants' model updates to obtain decoupled fragments (in our specific case, gradients of convolutional kernels). Then, we employ unsupervised anomaly detection as a means to uncover the fragments belonging to skewed distributions. We maintain a scoring matrix that counts the number of outliers for each participant. Finally, we prune the participants who scored highest, as their model updates are most suspicious in regard to the presence of backdoors. Figure 4 illustrates the framework of ARIBA. In Sec. 5, we demonstrate through extensive experiments that our method provides simple yet effective defenses.

4.2 Segmenting Local Updates

We start by segmenting fragments from the participants' updates. At any global round t , the server has on hand the shared model updates of all participants $\{\theta_i^{t+1}\}$ (some may be backdoored) and the global model weight of one round behind θ^t . The server obtains each participant's changes in model weights by

$$\delta_i^{t+1} \triangleq \theta_i^{t+1} - \theta^t. \quad (5)$$

Note δ_i^{t+1} actually represents the *accumulated* gradients of participants' local iterations during round t and is adjusted by its local learning rate. We here and later still call it gradients for simplicity. As we are interested in the *difference* between malicious and normal gradients, we further require the server to perform a mean subtraction on all $\{\delta_i^{t+1}\}$, as

$$\hat{\delta}_i^{t+1} \triangleq \delta_i^{t+1} - \frac{1}{n} \sum_{P_i \in \mathbf{P}} \delta_i^{t+1}, \quad (6)$$

which helps emphasize their disparity.

Algorithm 2 Scoring by Mahalanobis distance

Input: Fragments of all participants \mathbf{M} and the estimated \tilde{k} .

Output: The scoring matrix \mathbf{S} .

- 1: Initialize the scoring matrix: $\mathbf{S} = \{s_{i,j}\}^{n \times m}$.
 - 2: Calculate the covariance of fragments: $\Sigma \leftarrow cov(\mathbf{M})$.
 - 3: Calculate the mean of fragments: $\bar{M} \leftarrow mean(\mathbf{M})$.
 - 4: **for all** $i \leq n, j \leq m$ **do**
 - 5: Calculate the Mahalanobis distance for each $M_{i,j} \in \mathbf{M}$:

$$d_{i,j} \leftarrow \sqrt{(M_{i,j} - \bar{M})^T \Sigma^{-1} (M_{i,j} - \bar{M})}$$
 - 6: **end for**
 - 7: Set: $s_{i,j} \leftarrow 1$, if $d_{i,j}$ is one of the top- $(\tilde{k}m)$ largest $\{d_{i,j}\}$; $s_{i,j} \leftarrow 0$, otherwise.
 - 8: **return** \mathbf{S}
-

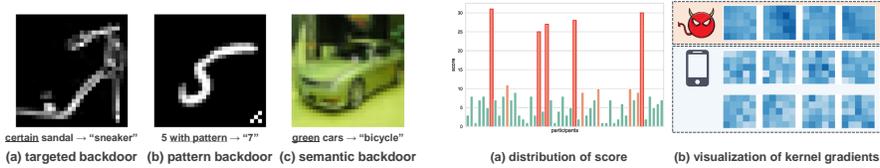


Fig. 5. Three types of backdoors we studied. The trigger can manifest as (a) certain samples, (b) samples carrying specific patterns, (c) samples with certain semantics. The same settings of backdoors are adopted in our experiments.

Fig. 6. (a) The scores of outliers of specific experimental case. All backdoored participants (marked read) are clearly discerned. (b) A difference in pattern can be observed from the visualization of some malicious and normal fragments.

To segment fragments, concretely, we pick a convolutional layer from the model and extract the gradients from each of its kernels, as illustrated in Fig. 4. Recall the kernels are the local feature detectors of CNN composed of small matrices of weights. We choose kernels as their gradients are semantically meaningful, but such practice is not a must and one is free to segment the model gradients in arbitrary ways, as long as the outcome is favorable for statistical analyses. We denote the derived fragments of participant P_i as $\mathbf{M}_i \triangleq \{M_{i,1}, \dots, M_{i,m}\} \subset \hat{\delta}_i^{t+1}$. Figure 6(b) visualizes some exemplar fragments of kernels from different P_i , where we can clearly observe the difference in pattern between malicious and normal fragments. This further testifies to our findings in Sec. 4.1.

4.3 Identifying Outliers in Fragments

Recall we try to discriminate the backdoored updates by the distributional bias of their fragments (Fig. 3(c)). However, how can the server make decisions based on the distributional disparity? We elucidate that *anomaly detection* can be leveraged as a convenient tool: Consider each fragment as an individual datum and map all fragments onto a hyperspace. Since most normal participants exhibit similar statistical distributions, their fragments will densely populate the projected region. Conversely, backdoored participants' fragments are liable

to project onto isolated and sparsely populated regions due to the observed skewed distribution. Thus, by performing anomaly detection in the hyperspace, the backdoored fragments are highly susceptible to being identified as outliers. As a result, *participants with significantly more outliers in their fragments are suspicious of providing backdoored updates.*

To embody the theory, we first gather the fragments of all participants $\mathbf{M} = \{\mathbf{M}_1, \dots, \mathbf{M}_n\}$. We then feed them into a **Scoring** function parameterized by \tilde{k} (Alg. 2), where unsupervised anomaly detection takes place to mark a portion ($\tilde{k}m$) of fragments as outliers. Here, \tilde{k} is the server’s estimated number of backdoored participants. In practice, the server can choose a generous \tilde{k} that ensures $\tilde{k} > k$. By ($\tilde{k}m$), we note a *flawless* anomaly detection algorithm would classify m fragments of all \tilde{k} participants as anomalous and classify the remaining as normal. Nevertheless, we anticipate (and tolerate) an approximation as some fragments are sure to be wrongfully classified due to their partial overlapping in distributions. We concretely choose a Mahalanobis-distance-based algorithm to elaborate our method, yet one is free to replace it with any unsupervised anomaly detection algorithms. **Scoring** returns a scoring matrix $\mathbf{S} = \{s_{i,j}\}^{n \times m}$, where $s_{i,j}=1$ if $M_{i,j}$ is marked as an outlier.

4.4 Pruning Backdoored Participants

The server now identifies backdoored participants $\tilde{\mathbf{P}}_m$ by counting the number of outliers. Conveniently, it calculates each participant’s score as

$$S_i = \sum_{j=1}^m s_{i,j}, \quad (7)$$

and puts participants $\{\tilde{P}_{m_1}, \dots, \tilde{P}_{m_{\tilde{k}}}\}$ with i -highest S_i ($i \leq \tilde{k}$) as $\tilde{\mathbf{P}}_m$. As the server speculates these \tilde{k} participants to provide backdoored updates, their model updates are pruned from being aggregated. The FedAvg aggregation (Eq. (1)) is therefore carried out on *cleaned* $\mathbf{P} \setminus \tilde{\mathbf{P}}_m$.

Note our defense is effective as long as $\mathbf{P}_m \subset \tilde{\mathbf{P}}_m$. Figure 6(a) exhibits the scores $\{S_i\}$ in a specific experimental case with $(n, k, \tilde{k})=(50, 5, 10)$, from which we can observe (1) the defense is effective as all malicious participants are identified, and (2) the difference between malicious and normal S_i is salient, suggesting a clear disparity in identification thus providing robust defense. We summarize the proposed ARIBA method in Alg. 1.

5 Experiments

5.1 Experimental Settings

We leverage ARIBA to identify the three types of backdoor triggers discussed in Sec. 3.3. FL models are trained on 3 common image datasets, MNIST, Fashion-MNIST, and CIFAR-10. We apply a 4-layer toy model for MNIST and Fashion-MNIST, and a ResNet18 for CIFAR-10. We by default choose $(n, k, \tilde{k})=(50, 5, 10)$,

i.e., the attacker compromises 5 out of 50 participants while the server generously estimates 10 of them as malicious. We later in Sec. 5.5 show over-estimation (choosing $\tilde{k} > k$) impacts model performance very slightly. We set the global learning rate $\eta=1$, which is in favor of the attacker as a larger η eases model replacement (Eq. (4)). We presume identical local learning rates among all participants, concretely, $\eta_p=1e-3, 1e-4, 1e-4$ for the three datasets, respectively. Experiments are carried out on an Nvidia 3090 GPU with PyTorch 1.10 and CUDA 11. The same random seed is sampled among all experiments.

5.2 Effectiveness of Our Defense

We first establish the backdoor attacks. Concretely, we train the FL model from scratch for $(t - 1)$ global rounds until it nearly converges. By the threat model in Sec. 3.2, the attacker waits (and behaves normally) till convergence. At round t , it begins with embedding the backdoor by letting each of its compromised participants train local updates on a mix of contaminated (that contains backdoor triggers) and normal data. It then scales the backdoored updates by σ (Eq. (3)), where by $(n, \eta)=(50, 1)$ we naturally have $\sigma=50$. We further presume each compromised participant P_{m_i} has an equal share of $\sigma_{m_i}=10$.

Baselines. We launch the attacks without defense. We aggregate θ^{t+1} with local updates of all clients $\{\theta_i^{t+1}\}$ with FedAvg (Eq. (1)). For each attack, we choose 30 images X' from the test dataset and implant them with the backdoor triggers, as exemplified in Fig. 5. We evaluate the model by test accuracy (*acc.*) and judge the attacker’s performance by *backdoor accuracy*, *i.e.*, the proportion of X' the model misclassifies according to its objective ($X'\%$). Higher $X'\%$ indicates successful attacks. As illustrated in Tab. 1(a), all three attacks succeed if without protection as the model wrongfully classifies most of the backdoored samples. Meanwhile, the compromised global model still performs well. This suggests one cannot identify the backdoor by examining model performances.

Effect of our defense. Now we plug in ARIBA before aggregation, *i.e.*, updates $\{\theta_i^{t+1}\}$ are first examined by our proposed technique in Sec. 4 where suspicious updates are pruned. Here, we introduce *confidence* $C = (\sum_{\{P_i \in \mathbf{P}_{m,j}\}} s_{i,j}) / (\tilde{k}m)$ to measure how sure the server is about the pruning decision: Higher C indicates a better defense, as the server manages to classify more of the attacker’s fragments as outliers, which contributes to the identification of backdoors. We illustrate the result of pruning, on the proportion of compromised participants pruned, by $\mathbf{P}_m\%$. We report test and backdoor accuracy as well. Results are summarized in Tab. 1(b). Note few X' may still be misclassified due to the model’s wrong prediction, even without the presence of attacks. We highlight that (1) ARIBA effectively defends all three types of backdoor attacks, as *all* compromised clients are identified and pruned ($\mathbf{P}_m\%=1.0$); (2) We further provide robust defense as the confidence C is high, indicating a clear distinguishability between malicious and normal participants. (3) The model retains high performance regardless of the excessive pruning ($\tilde{k} > k$).

Table 1. Summary of primary results. (a) The baseline FL models are compromised by backdoors. (b) Our proposed ARIBA provides an effective and robust defense in terms of confidence and proportion of attackers pruned. (c) Comparison with prior arts.

Backdoors	(a) baseline		(b) with ARIBA				(c) Prior arts						
	acc.	$X'\%$	acc.	$X'\%(\downarrow)$	C	$P_m\%$	[2] ₁	[2] ₂	[4]	[30]	[28]	[7]	[17]
Targeted	88.87	0.90	90.01	0.13	0.79	1.0	✗	✓	✗	✗	-	-	-
Pattern	98.85	0.83	99.24	0.00	0.95	1.0	✗	✓	✗	✗	✓	✓	-
Semantic	71.30	0.67	72.94	0.23	0.94	1.0	✗	✓	✗	✗	-	-	✓

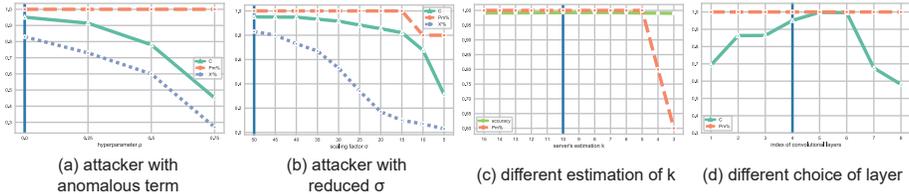


Fig. 7. Experimental results. Note the blue vertical lines mark our default settings. Our proposed method provides robust protection with high C and $P_m\%$ against two types of advanced attackers, which (a) use anomalous countermeasures to evade detection and (b) improve stealthiness with reduced σ . We further show (c) the excessive estimation of \tilde{k} influences model performance very slightly, and (d) the confidence of defense could vary by the concrete choice of layers.

5.3 Comparison with Prior Arts

We compare ARIBA with prior defenses discussed in Sec. 2.2. Results in Tab. 1(c) are summarized from both previous literature and our experimental studies. Here, “✓, ✗, -” indicates effective defense, ineffective defense, and no result claimed, respectively. Specifically, [2] proposes to examine the model’s test accuracy (₁) or the l2-norm of local updated weights (₂). We note though l2-norm effectively identifies baseline attacks, it can be easily bypassed by advanced attacks in Sec. 5.4. Krum [4] and CooMed [30] are byzantine-resilient aggregations, which defenses were nullified in [1]. [7, 17, 28] are attack-specific defenses that are not likely to generalize against other attacks. Some of these defenses require certain conditions such as auxiliary test datasets [17] or specific data distribution [28], which further constrains their practical use. Thereby, we argue ARIBA outperforms prior arts in terms of both generality and effectiveness.

5.4 Effectiveness on Advanced Attacks

We further study two varieties of advanced attackers, to illustrate ARIBA can provide effective protection even if against stealthy attack countermeasures.

Attacker with anomalous objectives. Studies [1, 2] suggest an attacker that could intentionally evade some metric-based detection. Specifically, it appends an anomalous term $l_a(\cdot)$ to the compromised participants’ training objective as

$\arg \min_{\theta_{m_i}} ((1 - \rho)l(f(X, \theta_{m_i}), y) + \rho l_a(g(\theta_{m_i})))$, where ρ is a hyperparameter and $g(\cdot)$ is the targeted metric. For example, by choosing $g(\theta_{m_i}) = \|\theta_{m_i} - \theta\|_2$ can the attacker deceive l2-norm bounding defense. Figure 7(a) presents the confidence C and $\mathbf{P}_m\%$ of ARIBA together with such attacker’s *baseline* $X'\%$ under different ρ . We note the anomalous term *does* affect the confidence however at the cost of lowering backdoor accuracy(in regard to $X'\%$). Nevertheless, ARIBA still provides intact protection by pruning all the compromised clients ($\mathbf{P}_m\%=1.0$).

Attacker with reduced σ . During model replacement, [1] suggest the attacker could reduce its capacity in exchange for better stealthiness, by choosing smaller $\sigma < \frac{\eta}{\eta}$ that *partially* replace the global model. We study attackers with different σ in Fig. 7(b). Results indicate ARIBA effectively identifies and prunes under most σ . Only in rare cases where σ is too small would ARIBA miss some participants. However, note the backdoor accuracy is concurrently impaired: At $\sigma=5,10$, we argue the attacker’s $X'\%$ is too low to incur an effective threat.

5.5 Ablation Study

Server’s estimation on \tilde{k} . In Sec. 4.3, the server is let estimate generously on the number of compromised participants $\tilde{k} > k$. Excessive estimation can cause wrongful pruning of normal clients. In Fig. 7(c), we show a generous \tilde{k} is acceptable as it affects model performance very slightly. On contrary, underestimation $\tilde{k} < k$ should be prohibited, since $\mathbf{P}_m\%$ reduces accordingly.

Choice of convolutional layers. By default, we choose one layer in the middle of the networks. We here alter the choice to observe its influence on defense. Results in Fig. 7(d) by C and $\mathbf{P}_m\%$ show though all choices provide effective protection, choosing in-the-middle layers mostly benefits robust identification. Note this seems to suggest *the attacker’s influence on model weights differs by the stages of model components*, which may be leveraged as more detailed detection clues. We leave it as an interesting open problem due to our limits of space.

6 Conclusion

This paper discusses the backdoor attacks against model faithfulness in FL systems. We present an in-depth study on state-of-the-art attacks by the attacker’s goal, capability, and possible attack approaches. By the observation of the magnitude and orientation disparity on the attacker’s model gradients, we advocate an effective defense, to identify and prune compromised participants by the distributional bias of their fragments, *i.e.*, gradient vectors of subset model parameters. We concretize our findings into the proposed ARIBA defense and demonstrate through extensive experiments its effectiveness and robustness.

References

1. Bagdasaryan, E., Veit, A., Hua, Y., Estrin, D., Shmatikov, V.: How to backdoor federated learning. In: International Conference on Artificial Intelligence and Statistics. pp. 2938–2948. PMLR (2020) [3](#), [4](#), [5](#), [6](#), [7](#), [12](#), [13](#)
2. Bhagoji, A.N., Chakraborty, S., Mittal, P., Calo, S.: Analyzing federated learning through an adversarial lens. In: International Conference on Machine Learning. pp. 634–643. PMLR (2019) [3](#), [4](#), [6](#), [7](#), [12](#)
3. Biggio, B., Nelson, B., Laskov, P.: Poisoning attacks against support vector machines. arXiv preprint arXiv:1206.6389 (2012) [2](#)
4. Blanchard, P., El Mhamdi, E.M., Guerraoui, R., Stainer, J.: Machine learning with adversaries: Byzantine tolerant gradient descent. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. pp. 118–128 (2017) [4](#), [5](#), [12](#)
5. Chen, X., Liu, C., Li, B., Lu, K., Song, D.: Targeted backdoor attacks on deep learning systems using data poisoning. arXiv preprint arXiv:1712.05526 (2017) [3](#)
6. Chen, Y., Su, L., Xu, J.: Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. Proceedings of the ACM on Measurement and Analysis of Computing Systems **1**(2), 1–25 (2017) [4](#), [5](#)
7. Fung, C., Yoon, C.J., Beschastnikh, I.: Mitigating sybils in federated learning poisoning. arXiv preprint arXiv:1808.04866 (2018) [4](#), [7](#), [12](#)
8. Goldblum, M., Tsipras, D., Xie, C., Chen, X., Schwarzschild, A., Song, D., Madry, A., Li, B., Goldstein, T.: Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses. arXiv preprint arXiv:2012.10544 (2020) [2](#)
9. Gu, T., Dolan-Gavitt, B., Garg, S.: Badnets: Identifying vulnerabilities in the machine learning model supply chain. arXiv preprint arXiv:1708.06733 (2017) [3](#), [4](#), [5](#), [6](#)
10. Guerraoui, R., Rouault, S., et al.: The hidden vulnerability of distributed learning in byzantium. In: International Conference on Machine Learning. pp. 3521–3530. PMLR (2018) [4](#), [5](#)
11. Huang, L., Joseph, A.D., Nelson, B., Rubinstein, B.I., Tygar, J.D.: Adversarial machine learning. In: Proceedings of the 4th ACM workshop on Security and artificial intelligence. pp. 43–58 (2011) [3](#)
12. Ji, Y., Zhang, X., Ji, S., Luo, X., Wang, T.: Model-reuse attacks on deep learning systems. In: Proceedings of the 2018 ACM SIGSAC conference on computer and communications security. pp. 349–363 (2018) [3](#)
13. Kairouz, P., McMahan, H.B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A.N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al.: Advances and open problems in federated learning. arXiv preprint arXiv:1912.04977 (2019) [2](#)
14. Knaan, Y.: Under the hood of the pixel 2: how ai is supercharging hardware. Google AI (2017) [1](#)
15. Konečný, J., McMahan, H.B., Yu, F.X., Richtárik, P., Suresh, A.T., Bacon, D.: Federated learning: Strategies for improving communication efficiency. arXiv preprint arXiv:1610.05492 (2016) [1](#), [4](#)
16. Li, Q., Wen, Z., Wu, Z., Hu, S., Wang, N., Li, Y., Liu, X., He, B.: A survey on federated learning systems: vision, hype and reality for data privacy and protection. arXiv preprint arXiv:1907.09693 (2019) [2](#)
17. Li, S., Cheng, Y., Wang, W., Liu, Y., Chen, T.: Learning to detect malicious clients for robust federated learning. arXiv preprint arXiv:2002.00211 (2020) [4](#), [12](#)

18. Li, T., Sahu, A.K., Talwalkar, A., Smith, V.: Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine* **37**(3), 50–60 (2020) [2](#)
19. Li, X., Huang, K., Yang, W., Wang, S., Zhang, Z.: On the convergence of fedavg on non-iid data. arXiv preprint arXiv:1907.02189 (2019) [5](#)
20. Liu, Y., Ma, S., Aafer, Y., Lee, W.C., Zhai, J., Wang, W., Zhang, X.: Trojaning attack on neural networks. *NDSS* (2018) [2, 3](#)
21. McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communication-efficient learning of deep networks from decentralized data. In: *Artificial intelligence and statistics*. pp. 1273–1282. PMLR (2017) [1, 4](#)
22. Rubinstein, B.I., Nelson, B., Huang, L., Joseph, A.D., Lau, S.h., Rao, S., Taft, N., Tygar, J.D.: Antidote: understanding and defending against poisoning of anomaly detectors. In: *Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement*. pp. 1–14 (2009) [3](#)
23. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*. pp. 618–626. IEEE Computer Society (2017). <https://doi.org/10.1109/ICCV.2017.74>, <https://doi.org/10.1109/ICCV.2017.74> [7](#)
24. Steinhardt, J., Koh, P.W., Liang, P.: Certified defenses for data poisoning attacks. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. pp. 3520–3532 (2017) [3](#)
25. Sun, Z., Kairouz, P., Suresh, A.T., McMahan, H.B.: Can you really backdoor federated learning? arXiv preprint arXiv:1911.07963 (2019) [3](#)
26. Tolpegin, V., Truex, S., Gursoy, M.E., Liu, L.: Data poisoning attacks against federated learning systems. In: Chen, L., Li, N., Liang, K., Schneider, S.A. (eds.) *Computer Security - ESORICS 2020 - 25th European Symposium on Research in Computer Security, ESORICS 2020, Guildford, UK, September 14-18, 2020, Proceedings, Part I. Lecture Notes in Computer Science*, vol. 12308, pp. 480–501. Springer (2020). https://doi.org/10.1007/978-3-030-58951-6_24, https://doi.org/10.1007/978-3-030-58951-6_24 [3](#)
27. Wang, H., Sreenivasan, K., Rajput, S., Vishwakarma, H., Agarwal, S., Sohn, J.y., Lee, K., Papailiopoulos, D.: Attack of the tails: Yes, you really can backdoor federated learning. arXiv preprint arXiv:2007.05084 (2020) [3](#)
28. Wu, C., Yang, X., Zhu, S., Mitra, P.: Mitigating backdoor attacks in federated learning. arXiv preprint arXiv:2011.01767 (2020) [4, 12](#)
29. Xie, C., Huang, K., Chen, P.Y., Li, B.: Dba: Distributed backdoor attacks against federated learning. In: *International Conference on Learning Representations* (2019) [3](#)
30. Yin, D., Chen, Y., Kannan, R., Bartlett, P.: Byzantine-robust distributed learning: Towards optimal statistical rates. In: *International Conference on Machine Learning*. pp. 5650–5659. PMLR (2018) [4, 5, 12](#)
31. Zhang, J., Chen, J., Wu, D., Chen, B., Yu, S.: Poisoning attack in federated learning using generative adversarial nets. In: *2019 18th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/13th IEEE International Conference On Big Data Science And Engineering (Trust-Com/BigDataSE)*. pp. 374–380. IEEE (2019) [2](#)
32. Zou, M., Shi, Y., Wang, C., Li, F., Song, W., Wang, Y.: Potrojan: powerful neural-level trojan designs in deep learning models. arXiv preprint arXiv:1802.03043 (2018) [3](#)