

FedDefender: Client-Side Attack-Tolerant Federated Learning

Sungwon Park*
KAIST
Daejeon, South Korea
psw0416@kaist.ac.kr

Sungwon Han*
KAIST
Daejeon, South Korea
lion4151@kaist.ac.kr

Fangzhao Wu
Microsoft Research Asia
Beijing, China
wufangzhao@gmail.com

Sundong Kim
GIST
Gwangju, South Korea
sundong@gist.ac.kr

Bin Zhu
Microsoft Research Asia
Beijing, China
binzhu@microsoft.com

Xing Xie
Microsoft Research Asia
Beijing, China
xingx@microsoft.com

Meeyoung Cha
IBS & KAIST
Daejeon, South Korea
mcha@ibs.re.kr

ABSTRACT

Federated learning enables learning from decentralized data sources without compromising privacy, which makes it a crucial technique. However, it is vulnerable to model poisoning attacks, where malicious clients interfere with the training process. Previous defense mechanisms have focused on the server-side by using careful model aggregation, but this may not be effective when the data is not identically distributed or when attackers can access the information of benign clients. In this paper, we propose a new defense mechanism that focuses on the client-side, called FedDefender, to help benign clients train robust local models and avoid the adverse impact of malicious model updates from attackers, even when a server-side defense cannot identify or remove adversaries. Our method consists of two main components: (1) attack-tolerant local meta update and (2) attack-tolerant global knowledge distillation. These components are used to find noise-resilient model parameters while accurately extracting knowledge from a potentially corrupted global model. Our client-side defense strategy has a flexible structure and can work in conjunction with any existing server-side strategies. Evaluations of real-world scenarios across multiple datasets show that the proposed method enhances the robustness of federated learning against model poisoning attacks.

CCS CONCEPTS

• **Security and privacy** → **Software and application security; Intrusion/anomaly detection and malware mitigation**; • **Computing methodologies** → **Artificial intelligence**.

KEYWORDS

Federated Learning, Client-Side Defense, Model Poisoning Attack, Knowledge Distillation, Meta Learning

*Equal contribution to this work

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '23, August 6–10, 2023, Long Beach, CA, USA.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0103-0/23/08...\$15.00
<https://doi.org/10.1145/3580305.3599346>

ACM Reference Format:

Sungwon Park, Sungwon Han, Fangzhao Wu, Sundong Kim, Bin Zhu, Xing Xie, and Meeyoung Cha. 2023. FedDefender: Client-Side Attack-Tolerant Federated Learning. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23)*, August 6–10, 2023, Long Beach, CA, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3580305.3599346>

1 INTRODUCTION

Federated learning has become a popular model training method to guarantee the minimum level of data privacy [5]. In each training round of federated learning, clients optimize their local models and send updates to the central server to aggregate them to produce a global model of the entire data distribution. Thus federated learning enables clients to jointly train a global model without directly sharing their private training data, making it a privacy-friendly solution [41]. Federated learning has been rapidly adopted by various applications that require data privacy [26, 30, 33, 43].

Despite its advantages, federated learning is vulnerable to attacks due to its decentralized nature. Any client can easily participate in the training process, introducing the possibility of maliciousness [3, 27]. For example, federated learning systems assume that all participants are benign and that their data can help improve the performance of resulting models. This leaves room for *model poisoning attacks*, wherein malicious users deceive the system by pretending to be benign and sending “poisoned” updates to the central server. This type of attack can disrupt parameter optimization [34] and adversely impact the model’s performance [13], undermining the integrity of the federated learning system.

Existing studies focus on using robust aggregation on the server side as a defense against model poisoning attacks. A central server can be trained to preserve updates from likely benign clients while discarding updates from likely corrupted clients or adversaries. Statistics like the trimmed mean or median can be used for outlier-resistant aggregation instead of simply averaging all updates. Fu *et al.*, extended this method by introducing a concept of confidence computed from residuals of repeated median estimator [11]. Other detection algorithms like Norm Bound, Multi-Krum, and FoolsGold assign lower weights to outlier updates and reduce their adverse effect [4, 12, 40]. However, these methods are limited in their ability to detect adversaries in the real world, where local datasets are no longer independent and identically distributed (i.e., non-IID). Non-IID makes benign local updates diverse and indistinguishable from corrupted ones [12]. Furthermore, extant defenses have been

breached by newer attacks that introduce elaborate local updates [3, 9, 35], leading server-side strategies alone to be obsolete.

This paper takes a step further by introducing a client-side defense strategy named FedDefender, which can run alongside existing server-side defense strategies to enhance resilience against poisoning attacks in federated learning. This strategy modifies the local training process of benign clients by obtaining robust local models even when the server-side defense fails to filter out corrupt updates from malicious clients.

FedDefender consists of two unique training components: (1) *attack-tolerant local meta update* that finds local parameters that are less susceptible to noisy training by malicious clients and (2) *attack-tolerant global knowledge distillation* that extracts the correct information from a noisy global model.

In our attack-tolerant local meta update component, FedDefender generates a synthetic batch of corrupted data from the local data, including randomly flipped labels. The local model is then perturbed with this noisy batch of samples through one iteration update before the main objective is optimized using the clean batch of samples. This meta update is analogous to a vaccination process, where synthetic noise is introduced to prevent the model from collapsing in the face of attacks, and to find the model parameters that are more tolerant to noise. To generate more realistic noise during training, FedDefender discovers k -nearest neighbors from the local model’s embeddings to replace the original label.

Meanwhile, the main objective of the attack-tolerant global knowledge distillation is to convey only useful knowledge from the possibly corrupted global model on the central server. Conventional knowledge distillation cannot be used in federated learning attack scenarios because attackers’ malicious influence on the corrupted global model can disrupt local training. Instead, we leverage the fact that the last layer of a model is more prone to overfitting from noisy updates than the intermediate layer [21, 38]. FedDefender can learn more reliable information from the global model without overfitting by distilling knowledge only to the intermediate layer via an auxiliary network. We also add a self-knowledge distillation objective to calibrate the prediction so that it works better for distillation.

Experiments show that FedDefender is highly effective in various scenarios when applied in conjunction with existing server-side defense strategies. Comparison with several recent baselines such as Norm bound, Multi-Krum, and ResidualBase shows that combining our method with these server-side strategies consistently results in better performance. For instance, in the label poisoning attack scenario, the classification accuracy increased by 18%, 20%, 17%, and 17% for the CIFAR-10, CIFAR-100, TinyImageNet, and FEMNIST datasets, respectively. We further show that FedDefender is able to effectively defend against advanced poisoning attacks like LIE [3], STAT-OPT [9], and DYN-OPT [35].

The major contributions and findings of our work include:

- Proposing a unique client-side defense strategy, FedDefender, that trains robust local models to thwart model poisoning attacks in federated learning.
- Designing an attack-tolerant local meta update that helps discover noise-tolerant parameters for local models by utilizing a synthetically corrupted training set.

- Introducing an attack-tolerant global knowledge distillation technique that efficiently aligns the local model’s knowledge to the global data distribution while reducing the adverse effects of false information in the possibly-corrupted global model.
- Showcasing that FedDefender can easily be applied in combination with any server-side defense strategies to enhance accuracy by 17-20% under poisoning attacks across various datasets.

The code and implementation details are available at the following URL: <https://github.com/deu30303/FedDefender/>.

2 RELATED WORK

2.1 Poisoning Attacks in Federated Learning

In recent years, there is a growing concern about security issues in machine learning, leading to the emergence of various model poisoning attacks [8, 36]. Federated learning is particularly susceptible to these attacks because malicious users can easily access intermediate processes of the training and send poisoning updates to the central server [27, 44]. Attackers can attack the training without access to the entire data distribution, as the data is decentralized and cannot be shared among clients [9, 37].

Model poisoning attacks can be divided into two categories: targeted attacks and untargeted attacks [13, 34]. In targeted attacks, malicious clients inject a backdoor into the central server so that any instance with the backdoor trigger will be classified as “targeted” without degrading the overall performance of the model. In untargeted attacks, on the other hand, attackers aim to degrade model performance indiscriminately across all classes. A simple and widely used untargeted attack is the label-flipping attack [45], in which malicious clients corrupt local models by randomly flipping labels of training samples from original classes to other classes.

Recently, some studies have used benign clients’ partial information to improve the stealthiness of untargeted attack performance. For example, Baruch *et al.* infer the standard deviation and the intensity factor based on benign client’s updates and inject perturbed model updates accordingly [3]. Fang *et al.* inject malicious updates by adding constant opposite direction noise estimated from benign clients’ updates mean [9]. Shejwalkar *et al.* propose a model poisoning attack by calculating a dynamic malicious update opposite to the direction of benign model updates [35].

This paper focuses on defending against *untargeted attacks* in federated learning using the client-side defense.

2.2 Server-Side Defenses Against Poisoning Attacks in Federated Learning

In federated learning, dimension-wise averaging is a commonly used and effective method for aggregating local updates [7, 28]. However, this naive averaging method is vulnerable to model poisoning attacks, which may succeed even with just a single malicious model update [1, 4]. To address this threat, two server-side defense strategies have been proposed: coordinate-wise aggregation and detection-wise aggregation. Both aim to filter out malicious updates from adversary clients during the aggregation and ensure the integrity of the federated learning process.

Coordinate-wise Aggregation. Coordinate-wise aggregation introduces outlier-resistant operations instead of averaging. For example, *Trimmed Mean* eliminates the largest and smallest values in

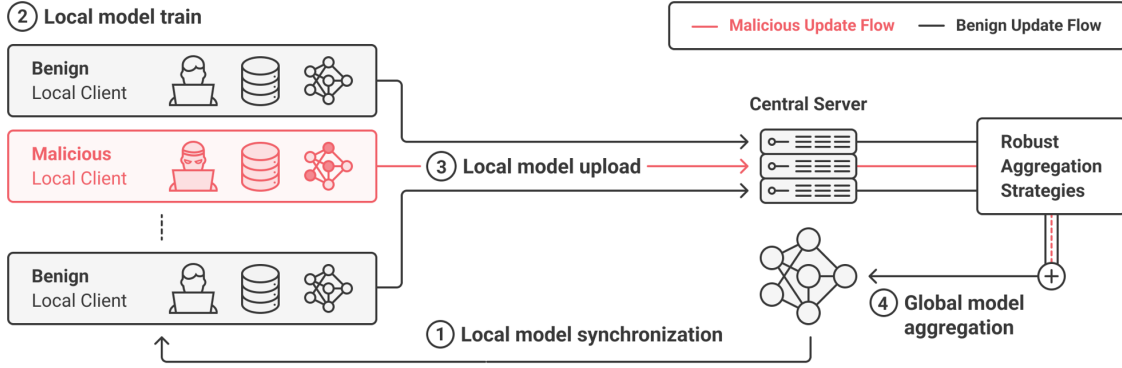


Figure 1: Illustration of federated learning under poisoning attacks. Attackers attempt to send malicious updates to the central server to corrupt the aggregated model. Unlike existing works that primarily focus on robust aggregation on the server side (stage ④), FedDefender focuses on training robust local models by benign clients (stage ②) to protect against malicious updates from adversaries.

each dimension, before computing the mean [47]. *Median* is another dimension-wise aggregation algorithm that computes the median of the updates in place of averaging [46]. However, both approaches are ineffective when the data distribution is non-IID as they may overlook underrepresented updates. To tackle this limitation, *ResidualBase* proposes a residual-based aggregation method [11], which calculates residuals of each parameter in the local model using a repeated median estimator. These residuals are then used to determine parameter confidence and detect false updates.

Detection-wise Aggregation. To mitigate the adverse effect of malicious clients’ updates, detection-wise aggregation adjusts learning rates based on the abnormality score of each local update. *Norm Bound* disregards clients with the norm of local updates being above a certain threshold by exploiting the observation that malicious clients often produce updates with a larger variance and norm than benign clients [40]. *Krum* introduces a Byzantine-resilient algorithm assuming that malicious updates would be placed far from benign updates in the Euclidean space [4]. More specifically, it selects a single local client update that is the most similar to its $n - m - 2$ neighboring updates to produce the global model, where m is the expected number of malicious local clients. *Krum* can be extended to *Multi-Krum* by iteratively running the algorithm to select multiple local updates, which shows better robustness than the original *Krum*. *FoolsGold* identifies grouped actions of targeted attacks based on the similarity score among local updates [12]. Unlike benign clients, malicious clients in a targeted attack scenario share the common loss objective and tend to have a more similar update pattern than benign clients. In this context, *FoolsGold* adjusts the learning rates of local models in proportion to the degree of diversity in each local update.

While the client-side defense has been relatively under-investigated, FL-WBC proposed a client-side defense strategy specifically designed to mitigate backdoor attacks [39]. In contrast, our proposed method, FedDefender, stands as a pioneering client-side defense strategy against untargeted attacks in federated learning. By complementing existing server-side defense strategies, it significantly

enhances the overall resilience against model poisoning attacks when used in conjunction with them.

3 PROBLEM FORMULATION

We will outline model poisoning attacks in federated learning and defenses against them in this section. Figure 1 highlights the basic flow of the training process in federated learning when malicious attackers are present.

3.1 Federated Learning

Overview. Federated learning aims to optimize a global model, denoted as F_θ , by utilizing data distributed across N local clients. We denote the loss objective and local dataset in the k -th client by \mathcal{L}_k and \mathcal{D}_k , respectively. The primary objective for training θ can be expressed as:

$$\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \frac{\sum_{k=1}^N |\mathcal{D}_k| \cdot \mathcal{L}_k(\theta, \mathcal{D}_k)}{\sum_{k=1}^N |\mathcal{D}_k|}. \quad (\text{Eq. 1})$$

In this paper, we use FedAvg [28], a widely used aggregation method as the default setting. FedAvg comprises four stages in each communication round:

- ① **Local model synchronization:** At the beginning of each round t , a random subset of local clients is selected, and these clients synchronize their model parameters by downloading the current global model parameter θ^t .
- ② **Local model train:** Each selected client k trains the downloaded model with her private data \mathcal{D}_k for a few epochs and computes the local model update $\Delta\theta_k^t = \theta_k^t - \theta^t$.
- ③ **Local model upload:** The local model update $\Delta\theta_k^t$ is sent to the central server.
- ④ **Global model aggregation:** The global model is updated by averaging all local updates received from the clients at the central server, as shown in Eq. 2:

$$\theta^{t+1} = \theta^t + \sum_{k=1}^N \frac{|\mathcal{D}_k|}{|\mathcal{D}|} \Delta\theta_k^t, \quad (\text{Eq. 2})$$

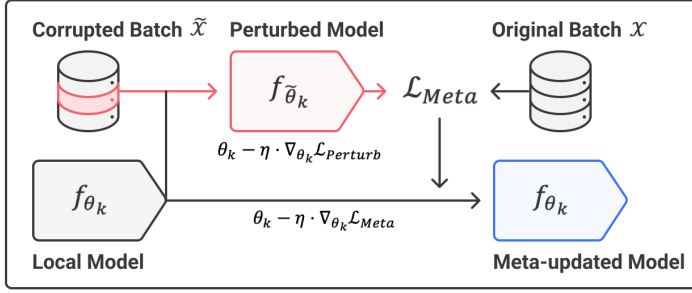
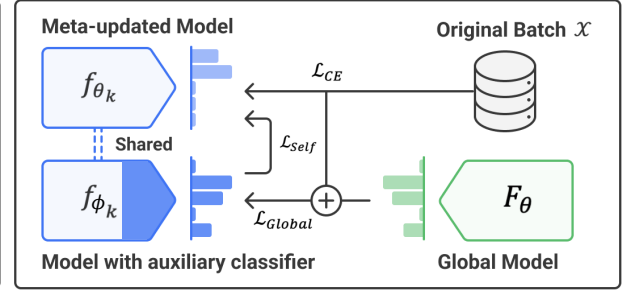
Step.1 Attack-Tolerant Local Meta Update (Sec. 4.1)**Step.2 Attack-Tolerant Global Knowledge Distillation** (Sec. 4.2)

Figure 2: The overall design of FedDefender framework. FedDefender comprises two steps: (1) attack-tolerant local meta-update, which finds local model parameters that are less prone to overfitting to noise, and (2) attack-tolerant global knowledge distillation, which aims to convey correct global knowledge from a potentially contaminated global model, regularizing the local model to mitigate data bias.

where $|\mathcal{D}| = \sum_{k=1}^N |\mathcal{D}_k|$. This process is repeated for each communication round.

Threat model. We hypothesize that all benign clients are honest and follow the proposed protocol, whereas malicious clients do not. Adversaries aim to disrupt the convergence and deteriorate the performance of the global model (i.e., an untargeted attack) by poisoning the local models of malicious clients. We also assume that the central server is trustworthy and performs the required task faithfully. Regarding the attacker’s capability, we assume that adversaries can bypass the verification process and infiltrate the federated system as participating clients (e.g., during optimization). We consider two scenarios for the threat model: (1) attackers have no access to benign clients’ information, and (2) attackers have partial access to benign client information, such as the standard deviation or average updates from benign clients.

3.2 Federated Learning When Facing Poisoning Attacks

Poisoning attacks rely on sending corrupted local updates to the central server, ultimately contaminating the trained global model. Several malicious objectives can be used for such attacks. For instance, attackers may adopt the same loss objective as legitimate clients for optimization, but use deliberately corrupted training set $\tilde{\mathcal{D}}_k$ to convey false information (i.e., $\mathcal{L}_k(\theta, \tilde{\mathcal{D}}_k)$). Additionally, attackers can also introduce an adversarial objective $\tilde{\mathcal{L}}_k$ that disrupts the overall convergence (i.e., $\tilde{\mathcal{L}}_k(\theta, \mathcal{D}_k)$).

We denote the subset of malicious clients as \mathcal{S}_m and the subset of benign clients as \mathcal{S}_b from the entire set of N clients \mathcal{S} (i.e., $\mathcal{S} = \mathcal{S}_m \cup \mathcal{S}_b$). To defend against poisoning attacks, the optimal objective for training the global model weight θ can be defined as:

$$\min_{\theta} \mathcal{L}^*(\theta) = \min_{\theta} \frac{\sum_{k \in \mathcal{S}_b} |\mathcal{D}_k| \cdot \mathcal{L}_k(\theta, \mathcal{D}_k)}{\sum_{k \in \mathcal{S}_b} |\mathcal{D}_k|}. \quad (\text{Eq. 3})$$

One approach to achieving this objective is to implement a server-side robust aggregation strategy for the global model aggregation (i.e., stage ④), which aims to retain the updates from benign clients while discarding all false updates from malicious clients during the global aggregation phase. The objective of robust aggregation can

be expressed as follows:

$$\theta^{t+1} = \theta^t + \frac{\sum_{k=1}^N \mathbb{1}_{\{k \in \mathcal{S}_b\}} \cdot \Delta \theta_k^t}{|\mathcal{S}_b|}, \quad (\text{Eq. 4})$$

where $\mathbb{1}$ is an indicator function. $|\mathcal{D}_k|$ is removed in order to mitigate the attack scenario where attackers try to manipulate the size of the local training data.

In contrast to existing defenses that focus on robust aggregation strategies, FedDefender aims to directly solve Eq. 3 by redesigning the local model train process (stage ②) and modifying the local updates of legitimate clients $\Delta \theta_k^t$. The detail is described in the next section.

4 DESCRIPTION OF OUR METHOD

FedDefender comprises the following two steps,

- **Step 1. Attack-Tolerant Local Meta Update (Section 4.1):** We train the benign local model f_{θ_k} in a robust manner via meta-learning. The goal is to discover the model’s parameters that produce accurate predictions even after it has been perturbed by noisy information. To achieve this, we first generate a noisy synthetic label, then apply one gradient update to perturb the local network parameters. Afterward, the gradient for the perturbed network to predict the correct outputs is computed and utilized to optimize the original local model.
- **Step 2. Attack-Tolerant Global Knowledge Distillation (Section 4.2):** If the global aggregation defense cannot block updates from malicious clients, the global model F_{θ} is no longer trustworthy. Given the meta-updated local model from the previous step, we apply global knowledge distillation to an auxiliary classifier using intermediate feature maps to neutralize the adverse impact on the contaminated global model F_{θ} . Self-knowledge distillation is applied between the auxiliary classifier and the original classifier to further incorporate the global knowledge into the deeper layers of the local model.

Figure 2 present the overall pipeline of our proposed defense. Our technique can mitigate model poisoning attacks in federated learning. We will now provide a detailed description of each component of FedDefender next.

4.1 Attack-Tolerant Local Meta Update

If the local model parameters θ_k are only optimized with a traditional supervised loss term (e.g., cross-entropy), it is susceptible to overfitting and being influenced by noise generated by malicious clients. Our key idea is to learn noise-tolerant parameters θ_k in a way that "vaccinates" the local model f_{θ_k} against model poisoning with synthetic noise, drawing inspiration from recent meta-learning works [10, 23, 32]. The proposed local meta-update replicates the training context with a model poisoning attack and makes the network less sensitive to noisy perturbations.

Local model poisoning with synthetic noise. Let us denote a mini-batch from local dataset \mathcal{D}_k as $\mathcal{X} = \{(\mathbf{x}_i, y_i)\}_{i=1}^B$, where \mathbf{x} is an input instance and y is the corresponding one-hot label. We want to generate synthetic batches $\tilde{\mathcal{X}} = \{(\mathbf{x}_i, \tilde{y}_i)\}_{i=1}^B$ with label noise to simulate the poisoning attack and perturb the local model.

Excessive perturbation can overly deform the model's decision boundary, leading to degraded performance. To mitigate this risk of severe deformation, we create more realistic noisy labels that resemble the distribution of \mathbf{y} by transferring labels from similar samples. For each $(\mathbf{x}, y) \in \mathcal{X}$, we calculate its feature representation $h_{\mathbf{x}}$ from the model's backbone network. Then, a random instance from top- k nearest neighbors in the representation space is randomly selected to replace the original label:

$$\tilde{\mathcal{X}} = \{(\mathbf{x}, \tilde{y}) | (\mathbf{x}, y) \in \mathcal{X} \text{ and } \tilde{y} = \text{Sample}_y(\mathcal{N}_k(\mathbf{x}, \theta_k))\}, \quad (\text{Eq. 5})$$

where $\mathcal{N}_k(\mathbf{x}, \theta)$ indicates the set of k -nearest neighbors of \mathbf{x} from the representation space made by f_{θ} . $\text{Sample}_y(\cdot)$ is the random selector function to extract the synthetic label. Since the nearest neighbors $\mathcal{N}_k(\mathbf{x}, \theta)$ are likely to share the same label as \mathbf{x} , this can generate noise within an acceptable range (i.e., 5-20% error rate).

Given the synthetic batch samples $\tilde{\mathcal{X}}$, we perturb the local model parameter θ_k using one gradient descent step:

$$\mathcal{L}_{\text{Perturb}} = \frac{1}{|\tilde{\mathcal{X}}|} \sum_{\mathbf{x}, \tilde{y} \in \tilde{\mathcal{X}}} H(\tilde{y}, f_{\theta_k}(\mathbf{x})) \quad (\text{Eq. 6})$$

$$\tilde{\theta}_k \leftarrow \theta_k - \eta \nabla_{\theta_k} \mathcal{L}_{\text{Perturb}}, \quad (\text{Eq. 7})$$

where $H(p, q)$ is the conventional cross entropy between p and q , and η is a learning rate.

Local model correction with meta update. To discover the model parameter that is less susceptible to noisy training by malicious clients, we propose a meta update to guide the model training. Given the perturbed model $\tilde{\theta}_k$, we optimize a classification loss $\mathcal{L}_{\text{Meta}}$ (Eq. 8) for one gradient descent step to encourage the perturbed local model $f_{\tilde{\theta}_k}$ to give correct predictions from \mathcal{X} . Note that the optimization process is applied to the parameters of the original local model θ_k (Eq. 9), although the loss calculation is based on the perturbed model's parameters $\tilde{\theta}_k$ (Eq. 8). This prevents the local model from being contaminated by synthetic noise during the training.

$$\mathcal{L}_{\text{Meta}} = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}, y \in \mathcal{X}} H(y, f_{\tilde{\theta}_k}(\mathbf{x})) \quad (\text{Eq. 8})$$

$$\theta_k \leftarrow \theta_k - \eta \nabla_{\theta_k} \mathcal{L}_{\text{Meta}} \quad (\text{Eq. 9})$$

4.2 Attack-Tolerant Global Knowledge Distillation

Given the meta-updated network f_{θ_k} , our next step is to enhance the performance of the refined local model θ_k through global knowledge distillation. In the case of non-IID settings, learning relying on local data leads to representation bias, as local data distribution differs from the global distribution. Regularization techniques, including global knowledge distillation, can be employed to control the local updates, thus preventing the local drift fallacy [15, 42, 49].

In the scenario of a model poisoning attack, however, the credibility of the global model F_{θ} can be compromised if the global defense fails to block malicious updates from hostile clients. This can lead to suboptimal results when the F_{θ} is perturbed. We propose attack-tolerant global knowledge distillation to mitigate the adverse effects of a potentially corrupted global model F_{θ} .

Refined knowledge distillation with an auxiliary network.

We start with the fact that a deeper layer in deep neural networks is much easier to overfit to noise (i.e., memorization) due to the inherent nature of gradient descent-based optimization [2, 21, 38]. In this context, FedDefender transfers global knowledge to a shallow intermediate part of the local model to reduce the adverse effect of false information. We attach an auxiliary classifier on top of intermediate layers to produce a new model parameter ϕ_k , and perform knowledge distillation on it (i.e., θ_k and ϕ_k share the shallow section of the entire network). Given an input data point \mathbf{x} from mini-batch \mathcal{X} , we regard the output probability from the global model F_{θ} as a target to train the local model parameter with auxiliary classifier ϕ_k . The global knowledge distillation loss is defined as follows:

$$F_{\theta}(\mathbf{x}, \tau) = \text{Sharpen}(F_{\theta}(\mathbf{x}), \tau) \quad (\text{Eq. 10})$$

$$\mathcal{L}_{\text{Global}} = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}, y \in \mathcal{X}} H(F_{\theta}(\mathbf{x}, \tau), f_{\phi_k}(\mathbf{x})), \quad (\text{Eq. 11})$$

where $\text{Sharpen}(\cdot)$ is the function that adjusts the confidence with sharpening temperature τ [17].

Since the global model's prediction is still not reliable in the attack scenario, we further introduce a simple strategy to refine target probabilities for robust knowledge distillation. Given a data point (\mathbf{x}, y) from \mathcal{X} , we calculate the scale coefficient α , which represents the cosine similarity between original label \mathbf{y} and the global model prediction $F_{\theta}(\mathbf{x})$. FedDefender then generates a refined label $\hat{\mathbf{y}}$ by applying a linear sum to calibrate the global information as follows:

$$\hat{\mathbf{y}} = (1 - \alpha) \cdot \mathbf{y} + \alpha \cdot F_{\theta}(\mathbf{x}, \tau). \quad (\text{Eq. 12})$$

If the global model's prediction is well aligned with the ground-truth labels, the scale coefficient becomes large and the final label weighs more on the global model's knowledge. Meanwhile, if the prediction is far different from the truth, we neglect the global model's information. We construct a modified batch $\hat{\mathcal{X}}$ with the refined labels $\hat{\mathbf{y}}$, and transfer the global knowledge to ϕ_k with this batch (Eq. 13).

$$\mathcal{L}_{\text{Global}} = \frac{1}{|\hat{\mathcal{X}}|} \sum_{\mathbf{x}, \hat{\mathbf{y}} \in \hat{\mathcal{X}}} H(\hat{\mathbf{y}}, f_{\phi_k}(\mathbf{x})) \quad (\text{Eq. 13})$$

Auxiliary self-knowledge distillation. We propose an additional design to improve the deeper layers of the local model, whose

weights are excluded from global knowledge distillation, in learning calibrated global knowledge. Motivated by the notion of self-knowledge distillation [14, 29, 48], we extract knowledge from the local model using an auxiliary classifier ϕ_k and transfer it back to the original local model θ_k . By doing so, we can distill the calibrated global knowledge learned from the refined label (Eq. 12) to the latter part of the model, effectively regularizing the entire model. Moreover, the output from the auxiliary head can be considered as a different view (i.e., augmentation) of the same instance. Maximizing the agreement between these two views can further enhance the model training. The self-knowledge distillation loss between the auxiliary classifier f_{ϕ_k} and the original model f_{θ_k} is defined as follows.

$$\mathcal{L}_{Self} = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} KL(f_{\theta_k}(\mathbf{x}, \tau) \| f_{\phi_k}(\mathbf{x}, \tau)), \quad (\text{Eq. 14})$$

where $KL(p||q)$ is a Kullback–Leibler (KL) divergence between two probabilities p and q . Finally, the final loss for our global knowledge distillation is given in Eq. 15.

$$\mathcal{L}_{KD} = \mathcal{L}_{Global} + \mathcal{L}_{Self} \quad (\text{Eq. 15})$$

This global knowledge distillation loss is optimized in conjunction with original cross-entropy loss to train the FedDefender (Eq. 18).

$$\mathcal{L}_{CE} = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} H(\mathbf{y}, f_{\theta_k}(\mathbf{x})) \quad (\text{Eq. 16})$$

$$\mathcal{L}_{total} = \mathcal{L}_{CE} + \mathcal{L}_{KD} \quad (\text{Eq. 17})$$

$$\theta_k \leftarrow \theta_k - \eta \nabla_{\theta_k} \mathcal{L}_{total} \quad (\text{Eq. 18})$$

5 EXPERIMENTS

We evaluate the robustness of FedDefender under various attacks on multiple datasets and compare it with server-side global aggregation baselines. We also examine how the model components and hyperparameters affect the overall performance.

5.1 Experimental Setup

Data settings. We use four image classification benchmark datasets in our experiments. The first two are CIFAR-10 and CIFAR-100, each containing 60,000 images of 32x32 pixels [20]. CIFAR-10 comprises 10 classes, such as airplanes, cats, and dogs, while CIFAR-100 comprises 100 classes. The third dataset is TinyImageNet, which comprises 100,000 images of 64x64 pixels and 200 classes. The fourth is the Federated Extended MNIST (FEMNIST) dataset, which contains 805,263 images of handwriting digits/characters of the size of 28x28 pixels [6].

The Dirichlet distribution is used to simulate the non-IID characteristics of CIFAR-10, CIFAR-100, and TinyImageNet datasets on federated learning. We denote the Dirichlet distribution by $Dir(N, \beta)$, where N is the total number of local clients and β represents the concentration parameter controlling the degree of non-IIDness of decentralized local data distributions. The probability $p_{k,j}$ is sampled from $Dir(N, \beta)$ and assigns the proportion of class j in k -th client’s dataset. With this non-IID distribution strategy with a low β value, local clients will have a disparate class distribution from one another. The default values for N and β are set to 20 and 0.5,

Table 1: Performance improvement with FedDefender on classification accuracy in Scenario-1 over four datasets. Our model brings non-trivial improvement for all server-side baseline algorithms for both last and best accuracy.

Method	CIFAR-10		CIFAR-100		TinyImageNet		FEMNIST	
	Last	Best	Last	Best	Last	Best	Last	Best
No Defense	68.80	71.96	42.97	43.90	30.37	38.98	18.88	23.81
+ FedDefender	78.17	79.96	51.76	51.92	35.59	39.68	22.11	24.48
Median	62.02	64.76	36.33	37.54	26.73	32.53	12.00	20.03
+ FedDefender	72.96	76.28	39.89	41.29	26.92	32.88	15.75	21.24
Trimmed Mean	75.52	76.03	47.93	47.95	36.76	38.24	19.37	23.17
+ FedDefender	81.06	81.96	52.90	53.41	38.85	39.37	21.80	23.95
Norm Bound	67.36	70.07	42.51	45.12	30.78	35.87	18.50	23.22
+ FedDefender	74.05	75.84	48.39	49.09	31.69	36.17	20.55	24.36
Multi-Krum	73.09	75.03	47.75	47.83	37.26	38.54	20.55	23.30
+ FedDefender	81.87	82.77	53.15	53.35	38.98	39.48	22.43	24.36
ResidualBase	73.61	75.10	44.80	45.13	35.05	38.60	19.44	23.86
+ FedDefender	79.28	80.83	50.62	50.98	36.22	39.24	22.41	24.27

respectively. In the case of FEMNIST, on the other hand, the writers of digits/characters in the data are randomly distributed to N clients. N is set to 20 for FEMNIST as well.

Implementation details. The number of communication rounds is set to 100, with 1 epoch per round for all federated learning experiments. Half of the clients (i.e., $10 = N/2$ clients) are randomly selected in each round for communication to make the federated setting more realistic. ResNet18 [16] is used as the default backbone network following the literature in federated learning [15, 24]. The learning rate (η), weight decay, and momentum for the SGD optimizer are set to 0.01, 1e-5, and 0.9, respectively. The batch size is set to 64. For knowledge distillation objectives, the temperature τ is set to 2. Random crop, color jitter, and random horizontal flip are used for data augmentations in local model training.

Model poisoning attack scenario. We divide a total of N clients into benign and malicious clients with a ratio of 8:2 (i.e., attacker ratio = 20%) as the default setting. For example, given N to be 20, we assign four clients to play the adversarial role from a pool of 20 clients. Because we randomly select 10 out of 20 clients in every communication round, the number of malicious clients in each round may vary. We consider two representative scenarios of an untargeted poisoning attack on the federated learning framework as follows:

Scenario-1) No access to benign clients’ information: This scenario assumes that malicious clients cannot obtain any information about benign clients. We consider the label flipping attack as the model poisoning attack in this case, as it does not require prior knowledge of the training data or the benign clients’ update information. For instance, in CIFAR-10, a class ‘dog’ image may have a random label of ‘ship’ or ‘horse’. This false model update is sent to the central server after training the local network with randomly flipped noisy labels.

Scenario-2) Partial access to benign client information: This scenario represents a more advanced attack in which malicious clients can access local updates from benign clients. Updates from

Table 2: Performance improvement with FedDefender on classification accuracy in Scenario-2 over four datasets. Our model brings non-trivial improvement for all server-side baseline algorithms for both last and best accuracy.

CIFAR-10	LIE		STAT-OPT		DYN-OPT		CIFAR-100	LIE		STAT-OPT		DYN-OPT	
	Last	Best	Last	Best	Last	Best		Last	Best	Last	Best	Last	Best
No Defense	73.40	76.68	73.74	73.99	62.37	63.13	No Defense	48.99	49.43	23.40	24.16	49.33	49.70
+ FedDefender	79.69	82.78	79.77	80.46	68.78	69.53	+ FedDefender	52.88	53.46	29.49	29.49	53.86	54.21
Median	55.71	64.80	53.37	58.17	74.71	76.65	Median	46.40	47.12	16.95	17.02	46.81	46.91
+ FedDefender	60.58	72.97	60.88	66.01	81.86	83.07	+ FedDefender	50.46	52.25	23.21	23.62	51.11	51.49
Trimmed Mean	39.19	76.29	64.46	68.34	72.80	75.18	Trimmed Mean	47.45	47.77	25.24	25.27	47.80	47.87
+ FedDefender	42.86	79.37	82.53	83.05	81.29	82.18	+ FedDefender	51.48	51.91	29.80	29.89	52.64	52.86
Norm Bound	12.30	27.90	64.15	64.70	72.54	75.09	Norm Bound	47.42	47.77	23.61	23.66	47.83	48.01
+ FedDefender	14.00	38.34	82.45	82.85	81.47	82.42	+ FedDefender	52.15	52.76	30.45	30.60	53.09	53.19
Multi-Krum	41.46	75.80	75.56	76.10	71.49	75.25	Multi-Krum	48.23	48.66	25.34	25.62	48.23	48.41
+ FedDefender	46.30	79.05	82.33	82.80	80.94	82.43	+ FedDefender	51.79	52.57	31.42	31.52	52.96	53.23
ResidualBase	74.53	77.04	70.24	70.39	74.40	76.93	ResidualBase	48.68	48.88	26.46	26.57	48.96	49.03
+ FedDefender	83.19	83.65	77.81	78.73	83.48	83.90	+ FedDefender	53.11	53.55	31.41	31.41	54.39	54.70

TinyImageNet	LIE		STAT-OPT		DYN-OPT		FEMNIST	LIE		STAT-OPT		DYN-OPT	
	Last	Best	Last	Best	Last	Best		Last	Best	Last	Best	Last	Best
No Defense	35.07	37.92	26.82	26.82	30.66	30.66	No Defense	20.07	22.30	19.67	22.43	19.41	22.63
+ FedDefender	37.61	39.27	29.62	29.62	32.50	32.52	+ FedDefender	21.87	24.03	21.68	24.13	21.42	24.29
Median	26.79	36.67	18.27	18.39	20.12	20.33	Median	17.81	23.15	10.56	14.60	19.14	22.52
+ FedDefender	33.42	36.90	21.68	23.63	23.65	23.89	+ FedDefender	20.58	23.39	15.02	19.16	21.01	24.16
Trimmed Mean	33.15	35.87	26.52	26.54	26.40	26.47	Trimmed Mean	19.01	22.21	19.48	22.83	18.19	21.71
+ FedDefender	35.70	37.28	27.26	27.29	31.42	31.48	+ FedDefender	21.67	23.73	21.99	24.24	20.70	23.77
Norm Bound	32.43	35.45	23.61	23.70	27.17	27.19	Norm Bound	20.50	22.70	19.14	21.97	19.55	22.25
+ FedDefender	35.59	37.42	26.66	26.66	31.45	31.49	+ FedDefender	21.64	23.79	21.59	24.15	20.89	23.88
Multi-Krum	33.38	36.32	26.39	26.40	26.31	26.31	Multi-Krum	20.25	22.43	19.88	23.16	19.56	22.82
+ FedDefender	35.22	37.04	27.83	27.95	31.34	31.44	+ FedDefender	21.93	23.44	22.06	24.16	20.96	23.87
ResidualBase	33.61	38.21	28.10	28.10	27.93	27.96	ResidualBase	20.21	22.92	19.07	21.53	19.28	23.08
+ FedDefender	36.73	38.73	28.49	28.51	31.85	31.88	+ FedDefender	22.66	24.01	22.01	23.97	21.49	24.15

malicious clients are computed using the statistics of benign updates. We consider three variations:

- (1) LIE algorithm infers both the standard deviation and intensity factors from the updates of benign clients. Then, the perturbed noise is generated by multiplying the two to deceive the server-side defense [3].
- (2) STAT-OPT injects a static inverse unit vector into the central server, computed as the opposite direction of the noise from the mean of the benign clients' updates [9].
- (3) DYN-OPT injects dynamically perturbed noise, where the maximum distance from any other updates is bounded by the maximum distance between any two benign updates [35].

5.2 Defense Performance Evaluation

Baselines. We have implemented six server-side defense strategies as baselines for untargeted attacks: (1) No Defense refers to the typical Federated Averaging (FedAvg) algorithm without any robust aggregation strategy. (2) Median is an aggregation algorithm that computes the median of each dimension of the updates rather than the average. (3) Trimmed Mean is an aggregation algorithm that computes the mean of local updates by removing the largest and smallest values. (4) Norm Bound is an algorithm that removes

updates if the norm of the local update is above a certain threshold. (5) Multi-Krum is an algorithm that iteratively selects local updates using the Krum method, which selects a single honest client by calculating the Euclidean distance between the client's update and the updates of its neighbors. (6) ResidualBase is an algorithm that computes parameter confidence using the residuals of each parameter from the repeated mean.

Evaluation. All models are evaluated using the same attacks and experimental settings (e.g., client pool, ratio of attacker numbers, attack strategy, number of local epochs, optimizer, and communication rounds) to ensure a fair comparison. We report the top-1 classification accuracy on the test set, including the last and best accuracy. The reported accuracy (labeled Last and Best) is calculated by averaging the accuracy of the last and best five rounds, respectively.

Result. Table 1 compares the performance of defense algorithms against the label flipping attack described in Scenario 1. We can see that FedDefender consistently improves the performance when applied to existing server-side defense algorithms. Despite unfiltered malicious updates from failures in server-side defense or the absence of a defense strategy (i.e., No Defense), incorporating FedDefender can lower the risk of performance degradation.

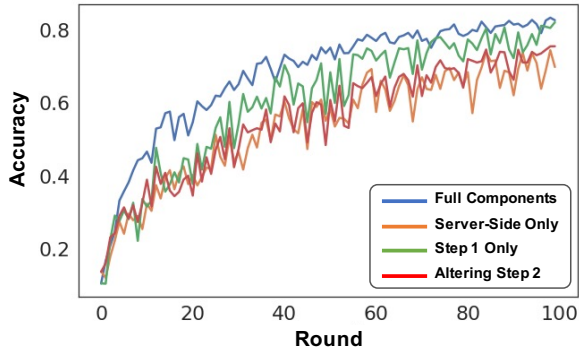


Figure 3: Performance comparison of ablations across communication rounds on CIFAR-10. Removing or altering any component results in a performance drop.

Table 2 reports the results against advanced attacks in Scenario 2, where malicious clients have access to information from benign clients. FedDefender consistently outperforms the baselines across all cases. Server-side defenses alone can lead to removing updates of benign users during training, which can cause the performance to drop. For example, existing server-side defense strategies perform poorly in LIE cases, even compared to No Defense. However, adding FedDefender can alleviate this issue and improve performance, highlighting the effectiveness of our approach.

5.3 Component Analyses

Ablation study. We have conducted an ablation study by removing and altering the key components. We compare the following variations:

- Full Components: Include all components.
- Step 1 Only: Only using the local meta update component (Section 4.1), with Step 2 removed.
- Altering Step 2: Altering knowledge distillation (Section 4.2) with a conventional alternative, which naively distills global knowledge without calibration.
- Sever-Side Only: A model that uses only the server-side defense method without any attack-tolerant local updates.

Figure 3 confirms that the full model provides the best performance among the ablations. Due to space limitations, we report results when our method is used along with Multi-Krum. Ablations using Step 1 only shows a more robust performance than using the server-side defense strategy alone. Altering Step 2 to a conventional knowledge distillation, on the other hand, degrades performance compared with entirely removing Step 2. This is likely due to corrupted information in the global model, F_θ , and demonstrates the need for a carefully designed attack-tolerant knowledge distillation.

Comparison with alternative global regularization approaches. The proposed global knowledge distillation (Section 4.2) may be replaced with existing global regularization techniques that don’t consider malicious attacks. We have tested the three alternatives that replace the global knowledge distillation with existing methods FedProx [25], Scaffold [19], or Moon [24] on top of the local meta update module. These existing methods also regularize the local

Table 3: Performance comparison of the proposed knowledge distillation technique with existing global regularization approaches. The proposed knowledge distillation technique achieves the best accuracy.

Global Knowledge Distillation	No Defense		Multi-Krum	
	Last	Best	Last	Best
Full Components	78.17	79.96	81.87	82.78
Local Meta Update + Scaffold	74.81	78.52	80.32	80.57
Local Meta Update + FedProx	74.67	78.01	80.83	82.01
Local Meta Update + Moon	73.13	75.12	78.05	78.99

model and address the local drift fallacy caused by the disagreement between local and global data distribution.

Table 3 compares our method with these alternatives over two different global aggregation strategies: simple averaging (i.e., No Defense) and Multi-Krum. FedDefender outperforms in both the last and best accuracy across both defense cases. When the global model is heavily perturbed (i.e., No Defense), the existing global knowledge distillation method experiences a significant drop in accuracy compared with when the global model is less perturbed (i.e., Multi-Krum). We also find that our proposed global knowledge distillation leads to smaller accuracy fluctuation across training rounds compared with other methods (see Figure 6 in the Appendix). These results demonstrate the critical role that our attack-tolerant global knowledge distillation plays.

Qualitative analyses. We investigate how well FedDefender discovers the attack-tolerant parameters for benign local models. We are inspired by the experiment to visualize the loss landscape in [22] and follow it as follows. We add random direction perturbations to the model parameter (for example, for two random directions X and Y in Fig. 4a) and examine how the accuracy changes. Results in Fig. 4a indicate that our local meta update (step 1) produces a higher and smoother accuracy surface compared with when this step is missing. Resistance to random parameter perturbations suggests that, when utilizing FedDefender, each benign client can find a solution with flat minima in the loss curve within the parameter space.

Next, we examine how well the correct global knowledge is conveyed to the local model by FedDefender under poisoning attacks. Based on the experiment, we may exactly compute a hypothetical global model that is immune to attacks (i.e., a clean global model) by removing all adversarial clients in the global aggregation phase. Mathematically, the corrupted and clean global models in round t are defined as follows:

$$\theta_{\text{corrupted}}^{t+1} = \theta^t + \frac{\sum_{k=1}^N \mathbb{1}_{\{k \in S_b \cup S_m\}} \cdot \Delta \theta_k^t}{|S_b| + |S_m|} \quad (\text{Eq. 19})$$

$$\theta_{\text{clean}}^{t+1} = \theta^t + \frac{\sum_{k=1}^N \mathbb{1}_{\{k \in S_b\}} \cdot \Delta \theta_k^t}{|S_b|}, \quad (\text{Eq. 20})$$

where S_b and S_m represent the sets of benign and malicious clients, respectively. We compute the similarity between the local models with two global models: one is corrupted and the other is clean. Average cosine similarity between model predictions over the same dataset is used to evaluate the model similarity.

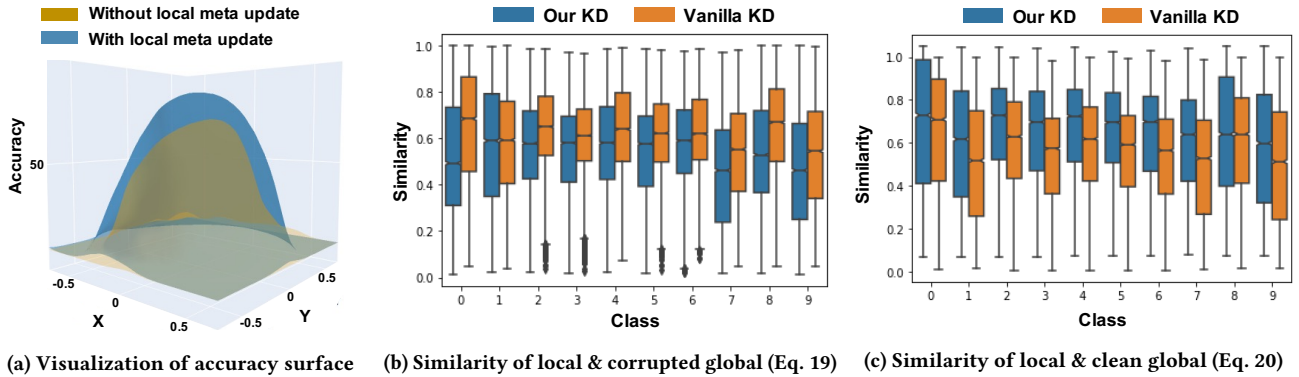


Figure 4: Qualitative analyses on FedDefender. (a) The accuracy surface of the global model with or without a local meta update after adding small perturbations to the model parameter. X and Y refer to two random directions for the perturbation. (b-c) The box plots the similarity between the local and global models (either corrupted or clean) with/without our proposed knowledge distillation.

Table 4: Performance comparison over different experimental settings with varying hyper-parameters. The results demonstrate that FedDefender consistently enhances the performance when it is applied on top of the existing aggregation strategy.

Client number (N)	Multi-Krum		Multi-Krum+Ours	
	Last	Best	Last	Best
10	82.53	83.21	84.66	84.94
15	81.26	82.75	86.11	86.89
20	73.09	75.03	81.87	82.77
25	65.18	73.81	81.67	82.34

(a) Effect of the number of clients

Attacker ratio (p_a)	Multi-Krum		Multi-Krum+Ours	
	Last	Best	Last	Best
10	74.11	76.72	80.75	82.74
15	73.74	76.56	81.73	82.96
20	73.09	75.03	81.87	82.77
25	72.52	75.85	77.99	81.97

(b) Effect of the percentage of attackers

Non-IIDness (β)	Multi-Krum		Multi-Krum+Ours	
	Last	Best	Last	Best
1.00	79.49	80.74	84.62	84.98
0.75	79.65	80.56	83.70	84.24
0.50	73.09	75.03	81.87	82.77
0.25	66.66	68.71	71.68	74.05

(c) Effect of the level of non-IIDness

Figures 4b and 4c show the relationship between local and global models for each class in CIFAR-10. The results show that local models in FedDefender have a substantially higher similarity to the clean global model (0.64) compared with the corrupted global model (0.53). In contrast, a vanilla knowledge distillation does not

guarantee the same ability to filter out contaminated knowledge from the global model (clean: 0.55 vs corrupted: 0.61). This suggests that our method can effectively distill more correct and calibrated knowledge from the corrupted global model.

Robustness test. To test robustness, we consider different experimental settings and vary the key hyperparameters. These include (a) the number of participating clients N , (b) the percentage of attackers that infiltrate the system p_a , and (c) the level of non-IIDness in the distributed local dataset, which is controlled by the beta parameter β in Dirichlet distribution. A lower beta parameter leads to higher non-IIDness. Table 4 shows the results for using Multi-Krum as a baseline aggregation strategy on CIFAR-10.

The complexity of data training increases as we increase the number of clients, the percentage of attackers, and the level of non-IIDness. Irrespective of these changes, FedDefender consistently demonstrates significant performance improvements.

6 CONCLUSION

This paper proposed FedDefender, a client-side approach to improve existing server-side defense strategies against model poisoning attacks. We proposed two attack-tolerant training components for benign clients: (1) attack-tolerant local meta update and (2) attack-tolerant global knowledge distillation. With these components, our method helps mitigate model poisoning attacks and produces more trustworthy results, even when server-side defenses fail to filter out malicious updates. As a result, our method has achieved a meaningful robustness improvement against various model poisoning attacks when used in conjunction with existing server-side defense strategies. We hope that our technique can serve as a foundation for further research in client-side robust federated learning.

ACKNOWLEDGEMENT

This research was supported by the Institute for Basic Science (IBS-R029-C2), Microsoft Research Asia, the Potential Individuals Global Training Program (2022-00155958), and the IITP grant (RS-2023-00216011) by the Ministry of Science and ICT in Korea.

REFERENCES

- [1] Sana Awan, Bo Luo, and Fengjun Li. 2021. Contra: Defending against poisoning attacks in federated learning. In *Proceedings of ESORICS*. Springer, 455–475.
- [2] Robert Baldock, Hartmut Maennel, and Behnam Neyshabur. 2021. Deep learning through the lens of example difficulty. In *Advances in NeurIPS*, Vol. 34. 10876–10889.
- [3] Gilad Baruch, Moran Baruch, and Yoav Goldberg. 2019. A little is enough: Circumventing defenses for distributed learning. In *Advances in NeurIPS*, Vol. 32.
- [4] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. 2017. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Advances in NeurIPS*, Vol. 30.
- [5] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, et al. 2019. Towards federated learning at scale: System design. In *Proceedings of MLSys*, Vol. 1. 374–388.
- [6] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. 2018. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097* (2018).
- [7] Jianmin Chen, Xinghao Pan, Rajat Monga, Samy Bengio, and Rafal Jozefowicz. 2016. Revisiting distributed synchronous SGD. *arXiv preprint arXiv:1604.00981* (2016).
- [8] Yao Chen, Yijie Gui, Hong Lin, Wensheng Gan, and Yongdong Wu. 2022. Federated Learning Attacks and Defenses: A Survey. *arXiv preprint arXiv:2211.14952* (2022).
- [9] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. 2020. Local model poisoning attacks to {Byzantine-Robust} federated learning. In *Proceedings of USENIX Security*. 1605–1622.
- [10] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of ICML*. 1126–1135.
- [11] Shuhao Fu, Chulin Xie, Bo Li, and Qifeng Chen. 2019. Attack-resistant federated learning with residual-based reweighting. *arXiv preprint arXiv:1912.11464* (2019).
- [12] Clement Fung, Chris JM Yoon, and Ivan Beschastnikh. 2020. The Limitations of Federated Learning in Sybil Settings. In *Proceedings of RAID*.
- [13] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733* (2017).
- [14] Sungwon Han, Sungwon Park, Sungkyu Park, Sundong Kim, and Meeyoung Cha. 2020. Mitigating embedding and class assignment mismatch in unsupervised image classification. In *Proceedings of ECCV*. Springer, 768–784.
- [15] Sungwon Han, Sungwon Park, Fangzhao Wu, Sundong Kim, Chuhuan Wu, Xing Xie, and Meeyoung Cha. 2022. FedX: Unsupervised Federated Learning with Cross Knowledge Distillation. In *Proceedings of ECCV*. 691–707.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of CVPR*. 770–778.
- [17] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2, 7 (2015).
- [18] Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. [n. d.]. Byzantine-Robust Learning on Heterogeneous Datasets via Bucketing. In *Proceedings of ICLR*.
- [19] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. 2020. Scaffold: Stochastic controlled averaging for federated learning. In *Proceedings of ICML*. 5132–5143.
- [20] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [21] Souvik Kundu, Qirui Sun, Yao Fu, Massoud Pedram, and Peter Beerel. 2021. Analyzing the confidentiality of undistillable teachers in knowledge distillation. In *Advances in NeurIPS*, Vol. 34. 9181–9192.
- [22] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2018. Visualizing the Loss Landscape of Neural Nets. In *Advances in NeurIPS*.
- [23] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. 2019. Learning to learn from noisy labeled data. In *Proceedings of CVPR*. 5051–5059.
- [24] Qinbin Li, Bingsheng He, and Dawn Song. 2021. Model-contrastive federated learning. In *Proceedings of CVPR*. 10713–10722.
- [25] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2020. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems* 2 (2020), 429–450.
- [26] Bing Luo, Xiang Li, Shiqiang Wang, Jianwei Huang, and Leandros Tassioulas. 2021. Cost-effective federated learning in mobile edge networks. *IEEE Journal on Selected Areas in Communications* 39, 12 (2021), 3606–3621.
- [27] Lingjuan Lyu, Han Yu, and Qiang Yang. 2020. Threats to federated learning: A survey. *arXiv preprint arXiv:2003.02133* (2020).
- [28] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS*. 1273–1282.
- [29] Sungwon Park, Sungwon Han, Sundong Kim, Danu Kim, Sungkyu Park, Seunghoon Hong, and Meeyoung Cha. 2021. Improving unsupervised image clustering with robust learning. In *Proceedings of CVPR*. 12278–12287.
- [30] Sungwon Park, Sundong Kim, and Meeyoung Cha. 2022. Knowledge sharing via domain adaptation in customs fraud detection. In *Proceedings of AAAI*, Vol. 36. 12062–12070.
- [31] Krishna Pillutla, Sham M Kakade, and Zaid Harchaoui. 2022. Robust aggregation for federated learning. *IEEE Transactions on Signal Processing* 70 (2022), 1142–1154.
- [32] Sachin Ravi and Hugo Larochelle. 2017. Optimization as a model for few-shot learning. In *Proceedings of ICLR*.
- [33] Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni, Spyridon Bakas, Mathieu N Galtier, Bennett A Landman, Klaus Maier-Hein, et al. 2020. The future of digital health with federated learning. *NPJ digital medicine* 3, 1 (2020), 1–7.
- [34] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciuc, Christoph Studer, Tudor Dumitras, and Tom Goldstein. 2018. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *Advances in NeurIPS*, Vol. 31.
- [35] Virat Shejwalkar and Amir Houmansadr. 2021. Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning. In *Proceedings of NDSS Symposium*.
- [36] Virat Shejwalkar, Amir Houmansadr, Peter Kairouz, and Daniel Ramage. 2022. Back to the drawing board: A critical evaluation of poisoning attacks on production federated learning. In *Proceedings of the IEEE Symposium on Security and Privacy*. 1354–1371.
- [37] Hira Shahzadi Sikandar, Huda Waheed, Sibgha Tahir, Saif UR Malik, and Waqas Rafique. 2023. A Detailed Survey on Federated Learning Attacks and Defenses. *Electronics* 12, 2 (2023), 260.
- [38] Cory Stephenson, Suchismita Padhy, Abhinav Ganesh, Yue Hui, Hanlin Tang, and Sueyeon Chung. 2021. On the geometry of generalization and memorization in deep neural networks. In *Proceedings of ICLR*.
- [39] Jingwei Sun, Ang Li, Louis DiValentin, Amin Hassanzadeh, Yiran Chen, and Hai Li. 2021. Fl-wbc: Enhancing robustness against model poisoning attacks in federated learning from a client perspective. *Advances in NeurIPS* 34 (2021), 12613–12624.
- [40] Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. 2019. Can you really backdoor federated learning? *arXiv preprint arXiv:1911.07963* (2019).
- [41] Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. 2021. A field guide to federated optimization. *arXiv preprint arXiv:2107.06917* (2021).
- [42] Chuhan Wu, Fangzhao Wu, Lingjuan Lyu, Yongfeng Huang, and Xing Xie. 2022. Communication-efficient federated learning via knowledge distillation. *Nature communications* 13, 1 (2022), 1–8.
- [43] Chuhan Wu, Fangzhao Wu, Lingjuan Lyu, Yongfeng Huang, and Xing Xie. 2022. FedCTR: Federated Native Ad CTR Prediction with Cross Platform User Behavior Data. *ACM Transactions on Intelligent Systems and Technology* (2022).
- [44] Chuhan Wu, Fangzhao Wu, Tao Qi, Yongfeng Huang, and Xing Xie. 2022. FedAttack: Effective and covert poisoning attack on federated recommendation via hard sampling. In *Proceedings of ACM SIGKDD*.
- [45] Han Xiao, Huang Xiao, and Claudia Eckert. 2012. Adversarial label flips attack on support vector machines. In *Proceedings of ECAI*. IOS Press, 870–875.
- [46] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. 2018. Generalized byzantine-tolerant sgd. *arXiv preprint arXiv:1802.10116* (2018).
- [47] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. 2018. Byzantine-robust distributed learning: Towards optimal statistical rates. In *Proceedings of ICML*. 5650–5659.
- [48] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. 2019. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of ICCV*. 3713–3722.
- [49] Weiming Zhuang, Yonggang Wen, Xuesen Zhang, Xin Gan, Daiying Yin, Dongzhan Zhou, Shuai Zhang, and Shuai Yi. 2020. Performance optimization of federated person re-identification via benchmark analysis. In *Proceedings of ACM MM*. 955–963.

A APPENDIX

A.1 Training Details

The overall training process is presented in Algorithm 1, and the FedDefender is open-sourced at <https://github.com/deu30303/FedDefender>.

The default backbone network for our experiments is ResNet18 [16], following prior work in federated learning [15, 24]. For the attack-tolerant local meta update (Step 1), we use a first-order approximation to speed up computation. The number of nearest neighbors, k , used to perturb the label is set to 10. In attack-tolerant global knowledge distillation (Step 2), we attach an auxiliary classifier to the output of the second residual block layer’s feature map.

We followed the details from original works for implementing server-side defense baselines. When deciding hyperparameters of Multi-Krum and Norm Bounding algorithms, we assumed that the central server already knows the upper bound of attacker numbers. The two hyperparameters in ResidualBase algorithm, confidence interval and clipping threshold, are set to 2.0 and 0.05 respectively.

In our experiment, Scenario-2 assumes a central-server aggregation agnostic attack, where malicious users have access to benign clients’ update information. We evaluate using three attack algorithms: LIE, STAT-OPT, and DYN-OPT. For LIE, we set the intensity factor to 0.3. For DYN-OPT, we set the initial intensity factor and threshold to 10 and $1e-5$, respectively. The optimal intensity factor value is then determined by repeatedly finding the optimal value between the smallest and largest values, until the intensity factor change falls under the threshold.

A.2 Further Results on Defense Performance

Comparison with other possible defense strategy. We have conducted additional comparison experiments with more recent baselines, including [18, 31, 39] under a label flipping attack scenario (Scenario 1). FL-WBC is a client-side defense, while the others are server-side defenses. Based on the experimental results in Table 5, we can confirm that our proposed model still offers a performance improvement for the local model when integrated with other server-side defenses. It is important to note that FL-WBC was originally designed exclusively for backdoor attacks, resulting in smaller performance gains compared with our method in untargeted attack scenarios.

Detection recall of Multi-Krum. Figure 5 shows the detection recall kernel density plot of the Multi-Krum algorithm, plotted

against the level of non-IID data distribution. As non-IID data distribution becomes more extreme (i.e., lower values of β), updates from benign clients become increasingly diverse. This makes it harder for server-side defenses to identify malicious client updates. In this light, we propose a client-side defense strategy, FedDefender, to achieve additional robustness and account for inherent model poisoning attacks during training.

Comparison with other global regularization approaches across communication rounds. Our proposed architecture introduces an attack-tolerant global knowledge distillation technique to mitigate the negative effects of a malicious global model during training. As an alternative to this method, one may also use existing approaches, such as FedProx [25], Scaffold [19], or Moon [24], in conjunction with our local meta update module. Figure 6 compares the classification accuracy of our method and alternative global regularization methods across communication rounds, without the use of any server-side defense. The models with other global regularization methods exhibit significant fluctuations in accuracy when a large number of attackers participate in each round. In contrast, our method’s accuracy remains stable and consistent throughout the rounds, even when the number of participating attackers is high.

Model comparison across communication rounds. The addition of FedDefender results in significant improvement in robustness against model poisoning attacks, compared to relying on server-side defense only. Figure 7 plots the accuracy changes throughout training rounds, including six baselines: No Defense, Median, Trimmed Mean, Norm Bound, Multi-Krum, and ResidualBase.

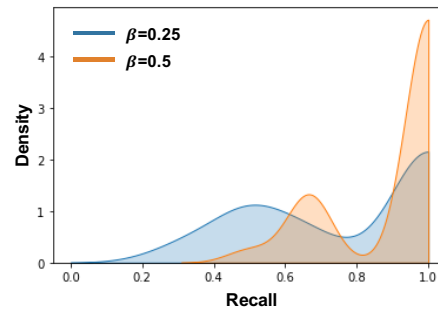


Figure 5: Detection recall plot of Multi-Krum with different levels of non-IID. The higher level of non-IID in data leads to the lower detection performance of the server-side defense strategy.

Table 5: Comparison of classification accuracy with other possible baselines in Scenario-1.

Method	CIFAR-10		CIFAR-100		TinyImageNet		FEMNIST	
	Last	Best	Last	Best	Last	Best	Last	Best
No Defense	68.80	71.96	42.97	43.90	30.37	38.98	18.88	23.81
+ FL-WBC	68.25	71.04	43.83	44.29	31.58	37.71	18.80	23.76
+ FedDefender	78.17	79.96	51.76	51.92	35.59	39.68	22.11	24.48
RFA	72.65	74.67	46.98	47.20	36.46	37.27	17.30	23.02
+ FedDefender	79.87	81.26	52.16	52.32	37.89	38.49	21.43	23.73
Bucket	70.12	72.14	48.45	48.46	37.55	38.87	19.88	22.70
+ FedDefender	77.57	79.37	53.32	54.44	39.06	40.07	21.82	24.63

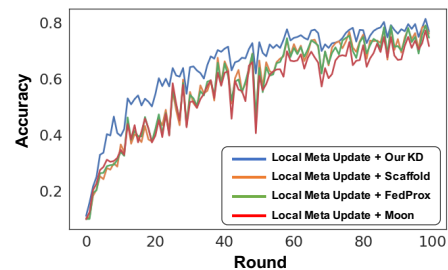


Figure 6: Performance comparison with other possible global regularization methods across communication rounds. The proposed distillation method worked the best for all rounds.

Algorithm 1 Attack tolerant local update algorithm of FedDefender.

Input: k -th Local Dataset \mathcal{D}_k , global model F_θ , k -th local model f_{θ_k} , k -th local auxiliary classifier model f_{ϕ_k} temperature τ , learning rate η
Output: One epoch updated k -th local model f_{θ_k}

```

/* In  $k$ -th local client one epoch update process */
 $f_{\theta_k} \leftarrow F_\theta$  //  $k$ -th local model synchronization by downloading the current global model
 $F_\theta.detach()$  // Global model gradient detach for knowledge distillation
for mini batch  $\mathcal{X} \in \mathcal{D}_k$  do
  /* Step 1. Attack-Tolerant Local Meta Update */
   $\tilde{\mathcal{X}} = \{(x, \tilde{y}) | (x, y) \in \mathcal{X} \text{ and } \tilde{y} = \text{Sample}_y(\mathcal{N}_k(x, \theta_k))\}$  // Perturbed mini batch  $\tilde{\mathcal{X}}$  using synthetic label  $\tilde{y}$ 
   $\mathcal{L}_{Perturb} = \frac{1}{|\tilde{\mathcal{X}}|} \sum_{x, \tilde{y} \in \tilde{\mathcal{X}}} H(\tilde{y}, f_{\theta_k}(x))$  // Local model poisoning with synthetic noises
   $\tilde{\theta}_k \leftarrow \theta_k - \eta \nabla_{\theta_k} \mathcal{L}_{Perturb}$ 
   $\mathcal{L}_{Meta} = \frac{1}{|\mathcal{X}|} \sum_{x, y \in \mathcal{X}} H(y, f_{\tilde{\theta}_k}(x))$  // Local model correction with meta update
   $\theta_k \leftarrow \theta_k - \eta \nabla_{\theta_k} \mathcal{L}_{Meta}$ 
  /* Step 2. Attack-Tolerant Global Knowledge Distillation */
   $\hat{\mathcal{X}} = \{(x, \hat{y}) | (x, y) \in \mathcal{X} \text{ and } \hat{y} = (1 - \alpha) \cdot y + \alpha \cdot F_\theta(x, \tau)\}$  // Refined mini batch  $\hat{\mathcal{X}}$  using global model  $F_\theta$ 
   $\mathcal{L}_{Global} = \frac{1}{|\hat{\mathcal{X}}|} \sum_{x, \hat{y} \in \hat{\mathcal{X}}} H(\hat{y}, f_{\phi_k}(x))$  // Calibrate global knowledge distillation
   $\mathcal{L}_{Self} = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} KL(f_{\theta_k}(x, \tau) || f_{\phi_k}(x, \tau))$  // Auxiliary Self-Knowledge Distillation
   $\mathcal{L}_{KD} = \mathcal{L}_{Global} + \mathcal{L}_{Self}$  // Calculate the entire knowledge distillation loss
   $\mathcal{L}_{CE} = \frac{1}{|\mathcal{X}|} \sum_{x, y \in \mathcal{X}} H(y, f_{\theta_k}(x))$  // Calculate the cross entropy loss using original batch  $\mathcal{X}$ 
   $\mathcal{L}_{total} = \mathcal{L}_{CE} + \mathcal{L}_{KD}$  // Calculate the total loss
   $\theta_k \leftarrow \theta_k - \eta \nabla_{\theta_k} \mathcal{L}_{total}$  // Update local model parameters via back-propagation
end

```

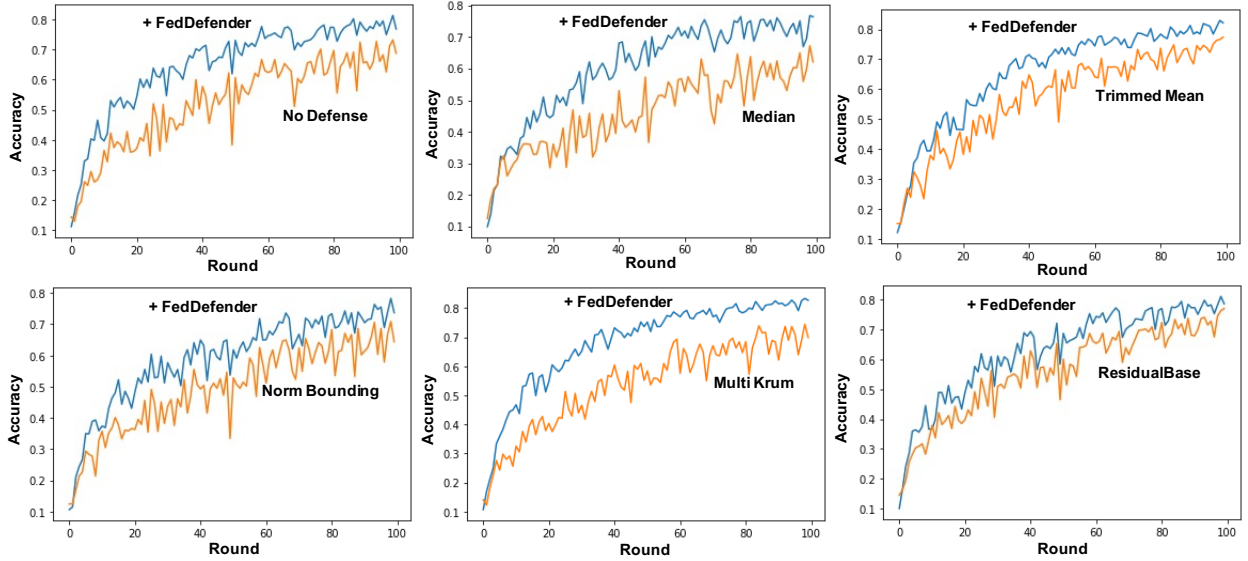


Figure 7: Performance comparison between server-side defense baselines and FedDefender-enhanced versions across rounds.