# FLTracer: Accurate Poisoning Attack Provenance in Federated Learning

Xinyu Zhang*, Qingyu Liu*, Zhongjie Ba, *Member, IEEE*, Yuan Hong, *Senior Member, IEEE*, Tianhang Zheng, Feng Lin, *Senior Member, IEEE*, Li Lu, *Member, IEEE*, and Kui Ren, *Fellow, IEEE*

*Abstract*—**Federated Learning (FL) is a promising distributed learning approach that enables multiple clients to collaboratively train a shared global model. However, recent studies show that FL is vulnerable to various poisoning attacks, which can degrade the performance of global models or introduce backdoors into them. In this paper, we first conduct a comprehensive study on prior FL attacks and detection methods. The results show that all existing detection methods are only effective against limited and specific attacks. Most detection methods suffer from high false positives, which lead to significant performance degradation, especially in not independent and identically distributed (non-IID) settings. To address these issues, we propose FLTracer, the first FL attack provenance framework to accurately detect various attacks and trace the attack time, objective, type, and poisoned location of updates. Different from existing methodologies that rely solely on cross-client anomaly detection, we propose a Kalman filter-based cross-round detection to identify adversaries by seeking the behavior changes before and after the attack. Thus, this makes it resilient to data heterogeneity and is effective even in non-IID settings. To further improve the accuracy of our detection method, we employ four novel features and capture their anomalies with the joint decisions. Extensive evaluations show that FLTracer achieves an average true positive rate of over 96.88% at an average false positive rate of less than 2.67%, significantly outperforming SOTA detection methods.** [1]

## I. INTRODUCTION

FEDERATED learning (FL) [1], a machine learning technique that allows multiple parties to collaborate on data analysis without revealing their private data, has been deployed at scale in the last decade [2], [3]. Under FL, a central server (e.g., cloud server) distributes a global model to multiple clients (e.g., smartphones and IoT devices). These clients individually train local models based on the global model with their own data. Then, the central server collects local model updates, aggregates them into a new global model, and starts the next round. However, recent studies show that malicious clients can apply a variety of poisoning attacks to the global model by disrupting their local models: *Untargeted attacks* that attempt to slow down the learning process, decrease the overall accuracy of the task, or even make the global model unusable [4], [5], [6]. *Backdoor attacks* (aka. targeted attacks) that aim to inject the backdoor into a model so that it incorrectly predicts Trojan data samples (i.e., samples containing attacker-chosen triggers) as the attacker-chosen class, without affecting the predictions for clean data samples [7], [8], [9].

To better understand the security threat posed to FL by poisoning attacks, we conduct a comprehensive empirical

study of existing attacks. Our study assesses the effectiveness, stability, and robustness of 14 untargeted and 16 backdoor attacks using 9 metrics (see § VI-B). Our findings indicate that the majority of untargeted attacks lead to a significant decrease in the accuracy or stability of the global model. Furthermore, most backdoor attacks can rapidly inject backdoors into the global model within approximately 30 training rounds, which persists for a duration exceeding 200 rounds.

However, existing detection methods suffer from limitations in accurately detecting attacks, especially in non-IID settings, for two primary reasons. First, adversaries can employ various stealthy attacks, such as untargeted adaptive attacks [10] that modify small portions of update parameters, and backdoor attacks [8], [11] that employ covert strategies involving data and code modifications. These stealthy attacks present difficulties for current detection methods [12], [5] that rely solely on identifying anomalous updates, resulting in a low true positive rate (TPR). The introduction of malicious updates can compromise the accuracy of the global model or inject backdoors. Second, in non-IID settings, the presence of highly heterogeneous feature and label distributions causes significant variation in benign client updates [13]. Consequently, existing cross-client detection methods [14] are prone to misclassifying and removing these highly heterogeneous updates from training, resulting in a high false positive rate (FPR). The lack of unique data from misclassified benign clients further diminishes the accuracy of the global model.

In this paper, we introduce *attack provenance* in FL, which accurately detects adversaries and traces the attack time, objective, type, and poisoned location of model updates. To our best knowledge, we propose the first FL attack provenance framework, FLTracer, as illustrated in Fig. 1. Regarding detection, FLTracer aims to minimize the introduction of poisoned data from a wide range of attacks (achieving high TPR) while maximizing benign data utilization (for low FPR), particularly in non-IID settings. Specifically, FLTracer extracts features from the model updates of each client and applies cross-client and cross-round detection to identify them. Based on the detection results, FLTracer aggregates the benign updates and provides attack tracing for malicious updates.

FLTracer addresses the aforementioned two limitations. Firstly, we introduce a novel update decoupling technique to generate four new sensitive features that represent the impact of various attacks from both value and function perspectives. Detecting these fine-grained features with *joint decisions* is more accurate than detecting updates directly, especially for stealthy attacks. Secondly, to further reduce the FPR of
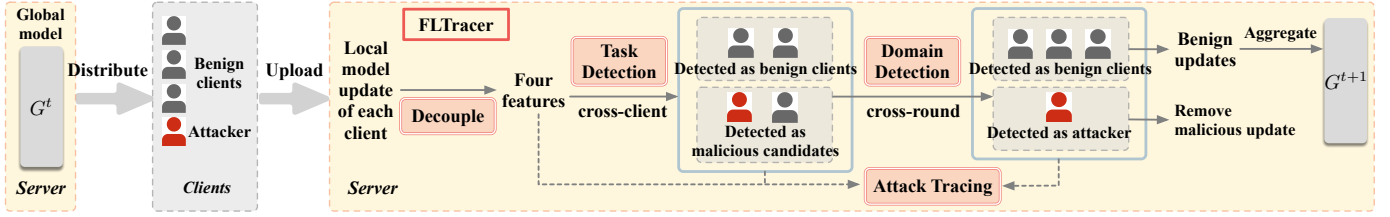
---

Figure 1: Overview of Attack Provenance (FLTracer), which accurately detects adversaries and traces the attack.

backdoor attacks in non-IID settings, we propose task and domain detection to capture *intrinsic differences* between adversary and benign clients. The detections are based on the following observations: (1) the adversary's training task is highly inconsistent with benign clients due to backdoor attacks involving learning both original and backdoor tasks, and (2) the adversary's feature domain (such as training data and code) changes drastically before and after the attack, in contrast to benign clients. Therefore, we propose task similarity (TSim) to measure the training task differences across clients in the same round, as well as domain distance (DDist) and a new *Kalman filter-based* domain detection to monitor and capture domain changes of the same client across rounds.

Our detection is resilient to data heterogeneity and effective in non-IID settings, as demonstrated by the following analysis. In real-world FL scenarios, such as autonomous driving and multiple edge devices, the data distribution collected by vehicles, security cameras, and other IoT devices tends to remain relatively constant or change slowly w.r.t. time, weather conditions, light, etc [15]. Different from prior works [16], [17], which consider benign clients' data distribution to remain constant, our framework also accommodates asymptotic distributions. For benign clients with fixed data distribution, we use the newly proposed Kalman filter to estimate the domain state of each client. Adversaries can be easily identified by their abrupt changes in domain state when launching attacks. For benign clients with gradually changing data, the parameters of the Kalman filter are dynamically updated to adapt to the evolving stream state of each client. Similarly, adversaries can be distinguished from benign clients by their abrupt state changes rather than gradual transitions. In cases where certain benign clients' data undergoes abrupt changes, we sacrifice those clients to ensure that no poisoned data is introduced.

Furthermore, our framework can trace attack attributions based on our detailed analysis of multiple features. Once an adversary is detected, the attack tracing can be provided: i) *the timing of the attack*, ii) *the objective and type of the attack*, and iii) *the specific poisoned location of model updates*.

Therefore, our contributions are summarized as follows:

- We conduct a comprehensive study of current FL poisoning attacks, including 14 untargeted and 16 backdoor attacks, along with 4 detection methods, in IID and non-IID settings. We evaluate these attacks using 9 metrics and demonstrate that existing attacks pose severe threats to FL.
- To the best of our knowledge, we propose the first FL attack provenance framework called FLTracer. It not only accurately detects adversaries with minimum global model performance drop, but also provides attack tracing. FLTracer utilizes 4 new sensitive features to analyze the different

crucial impacts of attacks on model updates, enabling the accurate detection of a wide range of attacks. To avoid the effects of data heterogeneity, FLTracer further leverages TSim and DDist to capture intrinsic anomalies of adversaries across clients and rounds, respectively.
- We conduct extensive experiments to evaluate FLTracer on 6 types of datasets and 7 types of model architectures. The results show high accuracy in IID and non-IID settings. In the non-IID settings, FLTracer achieves an average TPR of over $96.88\%$ and an average FPR of less than $2.67\%$, significantly outperforming state-of-the-art (SOTA) detection methods. It is worth noting that our detection is instantaneous, with detection validity limited to the current round.

## II. BACKGROUND

### A. Federated Learning

FL generates a shared model by iteratively training the global model on the central server, with the cooperation of $N$ distributed clients. In each round, the central server aggregates a global model from several local models trained individually by the clients using their local data. Specifically, at round $t$, the central server first broadcasts the current global model parameters $G^t$ to a randomly selected subset of $n$ clients. Then, each selected client initializes its local model parameters $L_i^t = G^t$ and updates $L_i^t$ with its local dataset. In this process, the selected client updates the local model individually by the stochastic gradient descent optimization for $E$ local epochs with a local learning rate. After obtaining $L_i^{t+1}$, each client calculates and sends the local model update $\theta_i^{t+1} = L_i^{t+1} - G^t$ to the server. Finally, the central server collects all updates and aggregates the new global model as $G^{t+1} = G^t + \eta \sum_{i=1}^{n} \theta_i^{t+1}$, where $\eta = \frac{1}{n}$ is the global learning rate. The new global model will be iteratively trained until achieving the predefined high accuracy on the test dataset or number of rounds.

### B. Poisoning Attacks on Federated Learning

**Untargeted attacks:** intend to slow down the training and eventually decrease the accuracy of the global model, including three main types. *Add noise attack* [12], [18] affects the training of the global model by altering the local update. It transforms the update $\theta_i$ to $\theta_i^* = \theta_i + \delta$, where $\delta$ is a random noise. *Sign-flipping attack* [19] is another local update modification attack. It flips the signs of parameters in an update, i.e., $\theta_i^* = -\theta_i$. It has a greater impact on the training than add noise attack. *Dirty label attack* [20] in centralized learning is a targeted attack, where the adversary misleads the training samples from one source class to another target class. However, under FL, it behaves like an untargeted attack, since the aggregation of the central server relieves this misdirection.

**Backdoor attacks (aka. targeted attacks):** aim to inject a backdoor into a model, such that the model predicts Trojan data as the attacker-chosen class, while making unaffected predictions on clean data. The process of the backdoor attack consists of a Trojan data synthesis phase and a backdoor injection phase. For the synthesis of Trojan data, the adversary embeds the trigger into the clean data and mislabels these data as the target category. There are two main types of backdoor triggers. A *patch trigger* is a small pattern consisting of a few pixels, such as a $3 \times 3$ checkerboard [21] or a reverse-engineered trigger [22]. A *perturbation trigger* is a perturbation of the same size as the training image, such as a random noise perturbation [23], a sinusoidal signal perturbation [24], or a reflected image [25].

For the backdoor injection phase, the adversary uses synthetic Trojan data to train backdoor models, including three injection strategies. *BadNets attack* [21] employs a classic injection strategy by training a local model with Trojan data containing backdoors and then poisoning the global model. *Model replacement attack* (MRA) [7] is a one-shot backdoor attack that amplifies the malicious updates by using a scaling factor before uploading to the server. *Distributed backdoor attack* (DBA) [11] is an FL-specific injection strategy that splits the backdoor trigger into multiple independent triggers and embeds the different triggers in multiple client datasets.

**Adaptive attacks:** are proposed to counter existing detection and defense methods. Experimental evidence has demonstrated that existing methods exhibit inadequate detection capabilities against adaptive attacks (see § VI-B). **Adaptive untargeted attacks** include Manipulating the Byzantine (*MB attack*) [5], *Fang attack* [4], and *LIE attack* [26], to counter Byzantine-robust aggregations. They require the cooperation of multiple malicious clients for effective execution and formulate an optimization problem to construct a malicious update that not only mimics the benign updates, but also significantly diminishes the model's performance. For **Adaptive backdoor attacks**, *Blind backdoor attack* [8] utilizes multi-task learning and the multiple gradient descent algorithm to balance both the original and backdoor tasks, achieving high performance on both tasks while maintaining stealthiness. We propose the *Feature Replacement Attack* (FRA) to manipulate only the weights of feature extraction layers, which are observed to be the critical updates for backdoor attacks. This approach aligns with the findings of the *Subnet replacement attack* [10], which effectively injects covert backdoors into models by modifying a limited set of model parameters.

*C. Kalman Filter*

The Kalman filter [27] is a recursive solution that utilizes the state-space model of a linear system to optimally estimate the system's state by integrating past observations with a dynamic model of the system's behavior. It is suitable for applications involving time-varying dynamical systems. At its core, the Kalman filter uses a state-space model that represents the underlying system to be estimated. The *state vector*, denoted as $x^t$, at *time step* $t$ can be estimated and predicted from the state vector $x^{t-1}$ at the previous moment, and the control input $u^t$, expressed as $x^t = Ax^{t-1} + Bu^t + w$. Here, $w$ is the *process noise*. The *state transition pattern* $A$ and the *control input matrix* $B$ are system parameters that describe the system's behavior, physical laws, and the relationships between the state variables $x$ and the *control inputs* $u$. It is worth noting that in time-varying dynamical systems, $A$ and $B$ may change over time. In such cases, the parameters can be estimated or updated based on the available historical measurement data and the estimation process. Then, the Kalman filter can be used to predict the next state vector.

## III. PROBLEM FORMULATION

*A. System Model*

The FL system commonly involves two entities: an aggregator deployed on the central server and multiple clients. The aggregator distributes the global model and aggregates updates. The clients provide local updates trained on their datasets. To prevent updates of malicious clients from joining the global model, we develop FLTracer, deployed on the central server, as depicted in Fig. 2. FLTracer has no restriction on datasets, i.e., the training data can be IID or non-IID, regardless of whether it is among benign or malicious clients. We assume that the distribution of a client's data across training rounds remains unchanged or gradually changes over time, which relaxes the assumption of unchanged data distribution in [17], [16]. Specifically, we discuss a practical scenario where clients' real-world driving dataset changed over time in § VI-D.
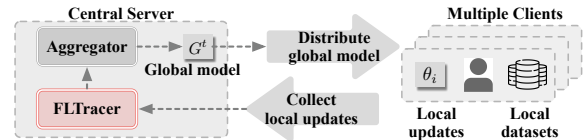


Figure 2: The system model

*B. Threat Model*

Following the threat model in the previous works [12], [7], [14], [28], [5], we consider a strong adversary who can compromise up to $50\%$ of the clients. However, the adversary cannot compromise the central server. The compromised clients will attempt to poison the global model, i.e., degrade its performance and inject backdoors. Table I summarizes the adversary's objectives and attack types. The details of the attacks are described in § II-B. Specifically, we consider five types of untargeted attacks for the adversary who wants to degrade the performance of the global model, and five types of backdoor attacks for backdoor injection. Both untargeted and backdoor attacks include adaptive attacks (w.r.t. existing defenses). For the central server, we assume it is honest and fully trusted. Thus, the inference attacks and corresponding countermeasures are not considered.

Regarding the adversary's capabilities, we consider a strong adversary who has full control over the compromised clients and knows the detailed design of our detection mechanism and the aggregation algorithm. Specifically, the adversary has five capabilities. ① Update modification (modif.): It can modify updates sent to the central server. ② Dataset modif.: It controls the local dataset of the compromised clients. Both the data and associated labels can be modified. ③ Code modif.: The

Table I: Specific attacks we considered with their corresponding objective, type, and capability.

| Objective | Attack type - capability | | Objective | Attack type - capability |
|---|---|---|---|---|
| Degrading performance (Untargeted attacks) | Add noise - ①④⑤ | | Inject the backdoor (Backdoor attacks) | BadNets attack (BN) - ②④⑤ |
| | Sign-flipping - ①④⑤ | | | Model replacement(MRA) -①②④⑤ |
| | Dirty label - ②④⑤ | | | Distributed backdoor(DBA) - ②④⑤ |
| | MB attack* - ①③④⑤ | | | Blind backdoor(Blind)*- ②③④⑤ |
| | Fang attack* - ①③④⑤ | | | Feature replacement(FRA)*-①②④⑤ |

"*" denotes adaptive untargeted attacks or adaptive backdoor attacks.
♯ with the capabilities of ④ and ⑤, all attacks will be more effective.

adversary is able to modify the local learning process by modifying the loss function. ④ Hyperparameter modif.: It can modify the hyperparameters to optimize the attacks, such as the number of epochs and the local learning rate. ⑤ Timing control: The adversary can decide the optimal time to launch attacks. Table I summarizes the capabilities of each attack.

## C. Design Goal

We aim to design a framework that can accurately detect a variety of attacks and provide attack tracing in both IID and non-IID settings, achieving the following goals:

**Accuracy:** The framework should guarantee the accuracy of the global model against SOTA poisoning attacks (including adaptive) in both IID and non-IID settings. Specifically, the detection should achieve a high TPR and a low FPR, so that the vast majority of benign clients can join the training.

**Backdoor-proof:** The framework should prevent backdoor injection into the global model even if the adversary performs a strong and stealthy backdoor attack.

**Provenance:** The framework should provide the attack tracing, including the attack time, objective, type, and specific poisoned location of the model update.

## IV. KEY OBSERVATIONS FOR ATTACK PROVENANCE

In this section, we present three key observations of FLTracer, as depicted in Fig. 3. To achieve attack provenance, we propose four new features (see § IV-A) to analyze update anomalies at different critical locations of the update. To further improve the detection accuracy of backdoor attacks, we introduce TSim (see § IV-B) and DDist (see § IV-C) to capture intrinsic differences between malicious and benign clients from task and domain perspectives, respectively.
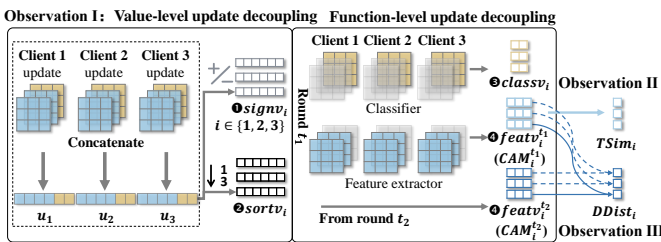


Figure 3: Overview of the key observations for FLTracer

## A. Poisoned Portions Analysis

**Observation I:** Previous methods [12], [5] typically distinguish between malicious and benign updates by detecting anomalies in the concatenation of the update values of the client $i$ ($u_i$). However, our results in § VI-C show that these methods are only effective for certain specific attacks for two primary reasons: First, the magnitude of the update values is

larger in shallow layers compared to deep layers, leading to the limited reflection of deep layer anomalies in $u_i$. Second, adaptive attacks such as FRA selectively modify critical updates, causing the unmodified normal values in $u_i$ to mask the anomalies in the modified portions. Therefore, we decouple updates into four sensitive features to detect anomalies, as decoupled training outperforms joint training [29].

**Update decoupling at the value and function level can analyze the specific poisoned portions of model updates by various attacks.** As described in § III-B, adversaries have the ability to execute a variety of attacks, either by directly manipulating updates or by using stealthy strategies such as modifying the dataset or code during training.

For attacks that directly modify updates, we decouple the concatenation of updates ($u_i$) into symbolic and numeric values. ❶ Sign vector (signv) is the concatenation of the signs of each element in $u_i$, reflecting the adversary's alterations to the update symbols. ❷ Sort vector (sortv) is the same size as $u_i$, and each element in $\text{sortv}_i$ indicates the order of the corresponding element ($u_i$) within $u_{\{i \in [n]\}}$. Specifically, for the 1st element, we arrange the 1st element of $n$ clients in ascending order. If the sorting result ranges from client $s_1$ to client $s_n$, then the 1st element for each client $i$ is $\text{sortv}_i^1 = s_i$. The sortv captures changes in the order of all elements for a client. This property ensures a higher level of robustness, as each element contributes equally to the anomaly of an update, thereby balancing anomalies in both shallow and deep layers.

For attacks that modify the dataset and code, the changed training process results in an alteration in the model's functionality. Thus, we decouple the updates regarding the classifier and the feature extractor. ❸ Classifier vector (classv) is the concatenation of weights and biases of the classification layer, capturing the adversary's modification of the decision boundary. ❹ Feature extractor vector (featv) extracts the fine-grained anomaly of the feature extractor, indicating the adversary's modifications to the feature space. For $k$-layer neural networks, such as Transformer [30], Principal Component Analysis (PCA) [31] is employed to compute anomaly scores ($\rho_i^\kappa \in \mathbb{R}$) for each encoder layer (or block) within the feature extractor. All anomaly scores of the encoder layer are concatenated as $\text{featv}_i = \{\rho_i^1, \cdots, \rho_i^k\} \in \mathbb{R}^k$.

Especially for the $k$-layer CNN, the detection of backdoors is challenging (see § VI-B). Considering that the convolution kernel serves as a processing unit [32], we introduce a novel feature extractor, called the Convolution Anomaly Matrix (CAM). This matrix contains the anomaly vectors ($A_i^\kappa$) from all convolutional layers. The $\text{CAM}_i$ for client $i$ is computed as

$$\text{CAM}_i = \{A_i^1, A_i^2, \cdots, A_i^k\}, \quad (1)$$

$$A_i^\kappa = w_i^\kappa \circledast w_{\{i \in [n]\}}^\kappa. \quad (2)$$

Here, $w_i^\kappa \in \mathbb{R}^{C_i^\kappa \times C_2^\kappa \times L_1^\kappa \times L_2^\kappa}$ denotes the convolution kernels in layer $\kappa$, where $C_1^\kappa, C_2^\kappa, L_1^\kappa \times L_2^\kappa$ represents the numbers of output channels, input channels, and kernel size. The operation $\circledast$ calculates the anomaly vector $A_i^\kappa \in \mathbb{R}^{C_1^\kappa \times C_2^\kappa}$ for client $i$ among all participants (see Appendix A). *In this paper, we describe the detection approach for the most complex features (i.e., $\text{CAM}_i$ and $A_i$). For the relatively simpler features $\text{featv}_i$ and $\rho_i$, our detection can use them instead of $\text{CAM}_i$ and $A_i$.*
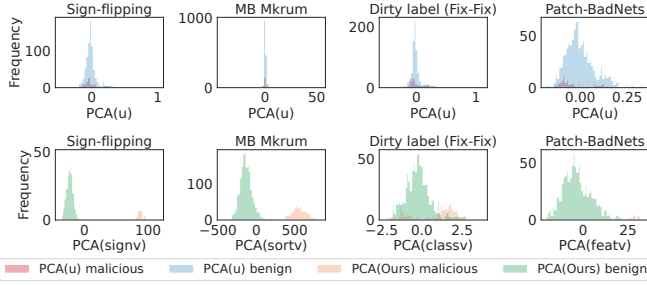
Figure 4: The distribution of dimensionality reduction results for the four features of malicious and benign clients. The top and bottom rows show the results for $u$ and our features.

To illustrate the advantages of our features over the feature in the previous work [5] (i.e., $u$), we perform four attacks on CIFAR10 in the non-IID settings. Fig. 4 shows the distribution of the PCA results for these four features. Our features discriminate the malicious and benign clients better than $u$.

### B. Task Consistency Analysis

**Observation II:** In FL, the central server collaborates with $N$ clients to perform the shared task, resulting in a consistent pattern within the convolution kernels of their local model. However, in a backdoor attack, this pattern deviates from the normal fixed pattern because the models are trained for both backdoor and regular tasks. To distinguish between backdoor and normal models, we leverage the similarity of convolution kernels across clients. For two convolution feature vectors $A_1$ and $A_2$, we define the *Task Similarity* (TSim $\in \mathbb{R}$) as

$$\text{TSim}_{A_1,A_2} = \frac{(A_1 - \overline{A_1}) \cdot (A_2 - \overline{A_2})}{\|A_1 - \overline{A_1}\| \ \|A_2 - \overline{A_2}\|}, \tag{3}$$

where $\overline{A_i} \in \mathbb{R}$ is the average of $A_i$. Here we use an adjusted cosine similarity metric, which mitigates the effect of different client data domains [33]. In essence, TSim focuses solely on assessing the correlation between convolution kernels within a convolutional layer, regardless of the overall value of $A_i$.

**A successful backdoor attack involves compromising the deep layers' weights in the feature extractor.** To analyze this, we specifically examine the variations between the layers in backdoor models and normal models. We extract the `CAMs` for each participant and plot the TSim in each layer (see Fig. 5). Notably, the deep layers of the backdoor models exhibit more anomalies, suggesting that the backdoor features tend to be embedded in the deeper layers. This phenomenon can be attributed to the fact that the model's shallow layers mainly capture texture information, whereas the deeper layers primarily contain structural details and possess a larger receptive field [32]. It is worth noting that some normal models with highly heterogeneous data exhibit similar anomalies to the adversary (see Fig. 5 (b) and (d)).

Therefore, in § V-B, we present a Task Detection approach to assess the consistency of tasks among participants within a round. Models that incorporate additional tasks, such as backdoor models or models trained with highly diverse data, are identified as potential malicious candidates. Furthermore, we introduce Observation III to differentiate between backdoor models and models trained using highly heterogeneous data.
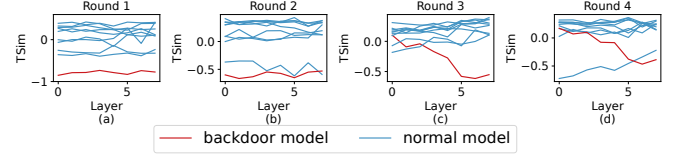


Figure 5: TSim of different layers in the backdoor and normal models each round. Lower TSim indicates more anomalous.

### C. Domain Consistency Analysis

**Observation III:** In a consistent data distribution, each client maintains fixed domain features. The fluctuation in the domain features of each processing unit (i.e., convolution kernel) between two consecutive rounds represents the domain variation of a client. We define the *Domain Distance* (DDist $\in \mathbb{R}$) between two convolution feature vectors of client $i$ in different rounds $t_1$ and $t_2$ as:

$$\text{DDist}_{A_i^{t_1},A_i^{t_2}} = \frac{1}{len(A_i^{t_1})} \cdot \sum_{j=1}^{len(A_i^{t_1})} (A_{i,j}^{t_1} - A_{i,j}^{t_2}), \tag{4}$$

where $j$ denotes each element in $A_i^{t_1}$ and $A_i^{t_2}$. The average anomaly value $(1/len(A_i^{t_1}))$ is used to balance anomalies for feature maps of different sizes.

**In non-IID settings, we accurately identify malicious clients through domain distance analysis. This approach rests on the observation that benign clients' domains remain consistent across rounds, whereas malicious ones exhibit inconsistency.** A common strategy for backdoor adversaries is to train a model normally and then launch a backdoor attack after the model converges [7]. The introduction of a backdoor trigger in the training data results in a sudden and drastic shift in the data domain. Conversely, the data domain of benign clients stays constant or undergoes gradual fluctuations over time. This divergence in behavior allows for anomaly detection by scrutinizing the domain consistency of a given client across rounds. To show the differences between malicious and benign clients, we plot their DDist across rounds (see Fig. 6). The DDist of malicious clients remains consistent prior to the attack, but exhibits drastic variation after the attack, clearly diverging from that of the benign clients.

Therefore, we propose a Kalman filter-based Domain Detection approach for detecting domain shift in clients across rounds in § V-C. Clients are identified as malicious if they exhibit anomalies in both TSim and DDist. Additionally, the consistency of DDist can help distinguish potential malicious candidates identified by Task detection (due to their highly heterogeneous data) from actual malicious clients.
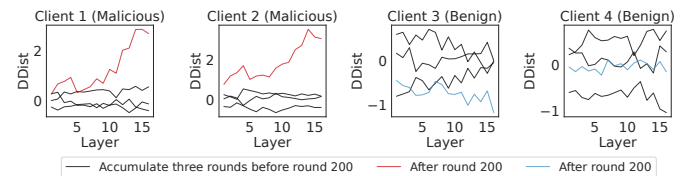


Figure 6: DDist of different layers for malicious and benign clients before and after the attack (launched at round 200).

### D. Attack Tracing Analysis

We present an analysis of attack tracing using newly introduced features and observations. Firstly, FLTracer is capable
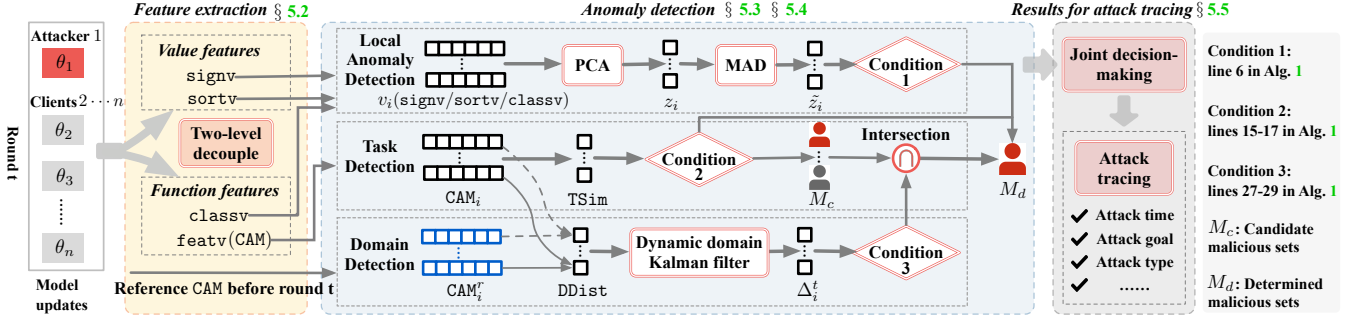
Figure 7: FLTracer (Attack Provenance) Framework. Upon receiving the client updates $\theta_{\{i\in[n]\}}$, FLTracer uses a two-level decoupling technique to generate four features representing different attacks (see § V-A). Then, it uses Local Anomaly Detection to identify anomalies in `signv`, `sortv`, and `classv`. For the complex feature, `featv` (CAM), FLTracer monitors task and domain anomalies concurrently by Task and Domain Detection (see § V-B and § V-C). In the results phase, FLTracer employs joint decision-making to accurately identify adversaries and provides attack tracing to support the detection process (see § V-D).

of identifying the adversary and pinpointing the time of attack when any feature detects an anomaly. Secondly, FLTracer enables the analysis of attack objectives and types in § III-B by detecting anomalies within the features (see Table II). The Add noise, MB attack, and Fang attack all alter the update. To capture their anomalies concerning the order of each element in the update, we employ the `sortv`. For the Sign-flipping attack, which preserves the absolute values of the update while flipping only the symbols, `signv` proves effective in its detection. For the Dirty label attack, which manipulates dataset labels without affecting the training data itself, we employ `classv`. This attack primarily modifies decision boundaries (i.e., classifier), with minimal impact on the feature space (i.e., feature extractor). For most backdoor attacks (e.g., BadNets, DBA, Blind, and FRA attacks) that inject backdoors into the original feature space, we use `featv` (CAM) to analyze the feature extractor of the backdoor model. Specifically, for the MRA attack, although it injects a backdoor by modifying the training data, it also directly scales up the updates of the backdoor model. Thus, detection using `sortv` is sufficient.

Table II: Attack Tracing

| Attack Type | Crucial Impact | Feature |
|---|---|---|
| Add noise, MB attack, Fang attack | Alter the values of the update, thus change their order among the participating clients | `sortv` |
| Sign-flipping | Only flip the symbols of the update | `signv` |
| Dirty label | Mainly modify the classifier | `classv` |
| BadNets, DBA, Blind, FRA | Mainly modify the feature extractor, leads to the anomaly of TSim and DDist | `featv` (CAM) |
| MRA | Scale up the values of the update | `sortv` |

## V. SYSTEM DESIGN

In this section, we present the framework of FLTracer, as depicted in Fig. 7. We define key notations in Table III.

### A. Value and Function Level of Decoupling

Based on Observation I, four newly proposed features represent the specific poisoned portions of different attacks on updates from the perspectives of value characteristics and model functions. For further information on feature extraction (labeled ❶ to ❹), please refer to § IV-A.

Table III: Key notations

| Term | Description |
|---|---|
| $N, m$ | Total number of clients, the total number of malicious clients |
| $n$ | Participating clients (participants) each round |
| $t, i$ | Round index, Client index |
| $G^t$ | Global model at round $t$ |
| $\theta_i^t$ | Local model update of client $i$ at round $t$ |
| $u_i$ | The concatenation of update values of client $i$ |
| $k$ | Total number of the layers |
| $\kappa$ | Layer index |
| $w_i^\kappa$ | Values of convolution kernels of client $i$ in layer $\kappa$ |
| $A_i^\kappa$ | Anomaly vector of the client $i$ in layer $\kappa$ |
| $\alpha_i^\kappa \in \mathbb{R}$ | Task similarity of client $i$ in layer $\kappa$ |
| $\beta_i^{\kappa,t} \in \mathbb{R}$ | Domain distance of client $i$ in layer $\kappa$ at round $t$ |
| $\text{CAM}_p = \{A_p^\kappa\}$ | Fixed feature pattern calculated by participants each round |
| $\text{CAM}_i^r = \{A_i^{\kappa,r}\}$ | Reference feature pattern of client $i$ |
| $\delta_i^\kappa$ | Difference between predicted and actual domain distance |

### B. Anomaly Detection

We propose Local anomaly detection to capture the anomalies of `signv`, `sortv`, and `classv`. For backdoor attacks, we analyze the feature extractor from the consistency of the task and domain perspectives. We propose Task and Domain detection to detect anomalies in CAM. Our method is applicable to other neural networks without convolution layers by replacing CAM and $A$ with `featv` and $\rho$ in the detection.

**Local Anomaly Detection:** As described in the LOCAL-ANOMALYDETECT function of Alg. 1, PCA is utilized to decrease the feature dimensionality (denoted as $v_i$) and acquire $z_i$. Then the median absolute deviation (MAD) algorithm [34] is employed to compute the anomaly scores of the $z_i$, which is robust in the presence of malicious clients. The anomaly scores ($\tilde{z}_i$) are represented by the following Eq.(5). Assuming a normal distribution, a value of $b = 1.4826$ is applied, and any client whose $\tilde{z}_i > \lambda$ (where $\lambda = 2.5$) is considered to have a 99.38% probability of being a malicious client [34].

$$\tilde{z}_i = f_{\text{MAD}}(z_i) = \frac{|z_i - \text{med}_q z_q|}{\xi}, \xi = b\,\text{med}_i\{|z_i - \text{med}_q z_q|\} \quad (5)$$

**Task Detection to Select Malicious Candidates:** This method identifies discrepancies in task similarity among $n$ participants in a round, as described in the TASKDETECT function of Alg. 1. Specifically, we require $\text{CAM}_{\{i\in[n]\}}$ to create a fixed feature pattern $\text{CAM}_p = \{A_p^1, A_p^2, \cdots, A_p^k\}$ (which represents the median of $n$ CAMs) for this round. Then we compute the

set of task similarities ($T_i$) between $\text{CAM}_i$ and pattern $\text{CAM}_p$ per layer, as expressed below:

$$T_i = \{\alpha_i^1, \alpha_i^2, \cdots, \alpha_i^k\}, \tag{6}$$

$$\alpha_i^\kappa = \text{TSim}_{A_p^\kappa, A_i^\kappa} = \frac{(A_p^\kappa - \overline{A_p^\kappa}) \cdot (A_i^\kappa - \overline{A_i^\kappa})}{\|A_p^\kappa - \overline{A_p^\kappa}\| \, \|A_i^\kappa - \overline{A_i^\kappa}\|} \tag{7}$$

where $\alpha_i^\kappa \in [-1, 1]$ is adopted from Eq.(3). A client is considered a malicious candidate if their $T_i$ value is abnormal. From § IV-B, we summarize three conditions that indicate abnormal $T_i$. First, $\min(T_i) \leq \tau, \tau \sim -1$ (see Fig. 5(a)), which suggests that certain layers in the $\text{CAM}_i$ of the malicious client differ significantly from benign ones. Second, $\max(T_i) \leq 0$ (see Fig. 5(b)), which indicates a negative correlation in the similarity between $\text{CAM}_i$ and the pattern $\text{CAM}_p$ for each layer. Third, the slope of the line connecting from $\alpha_i^1$ to $\alpha_i^k$ (in $T_i$) is less than 0 (see Fig. 5(c) and (d)). This implies that the similarity ($\alpha$) of the malicious client decreases as the layer goes deeper, in opposition to the trend of the benign one.

**Domain Detection to Determine Malicious Clients:** This method identifies malicious clients among malicious candidates by monitoring a client's change in domain distance across rounds based on § IV-C. We accumulate and construct a reference feature pattern $\text{CAM}_i^r = \{A_i^{1,r}, A_i^{2,r}, \cdots, A_i^{k,r}\}$ for each client. $\text{CAM}_i^r$ is the average of multiple $\text{CAM}$s for client $i$ in a few rounds. Then, we calculate the set of domain distances ($D_i^t$) between $\text{CAM}_i$ and $\text{CAM}_i^r$ by layer at round $t$, as:

$$D_i^t = \{\beta_i^{1,t}, \beta_i^{2,t}, \cdots, \beta_i^{k,t}\}, \tag{8}$$

$$\beta_i^{\kappa,t} = \text{DDist}_{A_i^{\kappa,t}, A_i^{\kappa,r}} = \frac{1}{len(A_i^{\kappa,r})} \cdot \sum_{j=1}^{len(A_i^{\kappa,r})} (A_{i,j}^{\kappa,t} - A_{i,j}^{\kappa,r}) \tag{9}$$

where $\beta_i^{\kappa,t}$ is adopted from Eq.(4). Specifically, we present a Kalman filter-based method to monitor changes in domain distance between rounds to determine malicious clients.

### C. Kalman Filter-based Domain Detection

We present a Kalman filter-based method to model the domain state of each client and detect malicious clients by comparing the predicted domain and actual domain states. If the difference exceeds certain conditions (explained later), we classify the client as malicious. Refer to the Kalman filter in § II-C, which models a system at time step $t$ with the equation $x^t = Ax^{t-1} + Bu^t + w$. Here, the core is to define each component of the system, including the state vector $x$, state transition pattern $A$, control input matrix $B$, control input $u$, and noise $w$. In addition, to ensure consistency with FL training, it is necessary to satisfy the following equations for the local model update $\theta_i^t$ and global model update $G^t$:

$$\theta_i^t = \mathbf{F_i^1} \cdot \theta_i^{t-1} + \mathbf{F_i^2} \cdot G^{t-1} + \mathbf{H}_i \circ \Omega_i + \mathbf{W}_i \tag{10}$$

$$G^t = \qquad\qquad G^{t-1} + \qquad \Omega_i + \mathbf{P}_i \cdot \Theta_i^t \tag{11}$$

where $\Omega_i$ represents the impact of client $i$'s datasets on updates and $\Theta_i^t = \sum_{l=1}^n \theta_l^t - \theta_i^t$ denotes the sum of updates from other clients. Client-specific parameters include $\mathbf{F_i^1}, \mathbf{F_i^2}, \mathbf{H}_i, \mathbf{W}_i$, and $\mathbf{P}_i$. Eq.(10) indicates that each client $i$'s update is related to their update at the previous round ($\theta_i^{t-1}$), their dataset ($\Omega_i$), and the starting training model ($G^{t-1}$). Eq.(11) indicates that

---

**Algorithm 1** Anomaly Detection Functions

1: **function** LOCALANOMALYDETECT($v_{\{i \in [n]\}}$)
2:     Initialize $M_d = \emptyset$
3:     $z_{\{i \in [n]\}} \leftarrow \text{PCA}(v_{\{i \in [n]\}}, \text{componenets} = 1)$
4:     $\tilde{z}_i \leftarrow f_{\text{MAD}}(z_i)$ using Eq.(5), $i \in [n]$
5:     **for** $i \in [n]$ **do**
6:         **if** $\tilde{z}_i > \lambda$ **then** $M_d \leftarrow M_d \cup \{i\}$
7:     **return** $M_d$     ▷ Determined malicious sets
8: **function** TASKDETECT($\text{CAM}_{\{i \in [n]\}}, k$)
9:     Initialize $M_d, M_c \leftarrow \emptyset$
10:     **for** layer $\kappa \in [k]$ **do**
11:         $A_p^\kappa \leftarrow \text{med}_i A_i^\kappa$
12:         $\alpha_i^\kappa \leftarrow \text{TSim}(A_p^\kappa, A_i^\kappa)$ using Eq.(7), $i \in [n]$
13:     $T_i \leftarrow$ Generated using $\alpha_i^\kappa$ and Eq.(6), $i \in [n]$
14:     **for** $i \in [n]$ **do**
15:         **if** $\min(T_i) \leq \tau$ **then** $M_d \leftarrow M_d \cup \{i\}$
16:         **else if** $\max(T_i) \leq 0$ **then** $M_c \leftarrow M_c \cup \{i\}$
17:         **else if** $\text{slope}(T_i) \leq 0$ **then** $M_c \leftarrow M_c \cup \{i\}$
18:     **return** $M_d, M_c$ ▷ Determined and candidate malicious sets
19: **function** DOMAINDETECT($\text{CAM}_{\{i \in [n]\}}, \text{CAM}_{\{i \in [n]\}}^r, k, M_c, M_d$)
20:     **for** layer $\kappa \in [k]$ **do**
21:         $\beta_i^\kappa \leftarrow \text{DDist}(A_i^{\kappa,r}, A_i^{\kappa,t})$ using Eq.(9), $i \in [n]$
22:     $D_i^t \leftarrow$ Generated using $\beta_i^\kappa$ and Eq.(8), $i \in [n]$
23:     $\hat{D}_i^t \leftarrow$ Estimated using Eq.(13), $i \in [n]$
24:     $\Delta_i^t = \{\delta_i^{\kappa,t}\} \leftarrow$ Calculated using $D_i^t, \hat{D}_i^t$ and Eq.(14), $i \in [n]$
25:     **for** $i \in [n]$ **do**
26:         **if** $\text{slope}(\Delta_i^t) > 0$ **then**
27:             **if** $f_{\text{MAD}}(\text{slope}(\Delta_i^t)) > \lambda$ **then** $M_c \leftarrow M_c \cap \{i\}$
28:         **else if** $f_{\text{MAD}}(\delta_i^{\kappa,t}) > \lambda$ **then** $M_c \leftarrow M_c \cap \{i\}$
29:         **else if** $\text{mean}(\Delta_i^t) > 0$ **then** $M_c \leftarrow M_c \cap \{i\}$
30:     **return** $M_d \leftarrow M_d \cup M_c$   ▷ Determined malicious sets

---

the global model is related to the previous round's global model ($G^{t-1}$), the training data of client $i$ ($\Omega_i$), and the updates of other participants, except client $i$ ($\Theta_i^t$).

**Dynamic Domain State for FL:** To model the domain state, we first calculate the domain distance set of $\theta_i^t$, $G^t$, and $\Theta_i^t$ and obtain $D_i^t$, $D(G)_i^t$, and $D(\Theta)_i^t$ using Eq.(8). We define the state vector as $X_i^t = [D_i^t, D(G)_i^t] \in \mathbb{R}^{2k \times 1}$ and the control input as $Z_i^t = D(\Theta)_i^t \in \mathbb{R}^{k \times 1}$. According to Eq.(10) and Eq.(11), the dynamic domain state of the FL for each client $i$ is modeled as follows, where $\circ$ denotes the Hadamard product:

$$X_i^t = \begin{bmatrix} \mathbf{F_i^1} & \mathbf{F_i^2} \\ \mathbf{0}_{k \times k} & \mathbf{1}_{k \times k} \end{bmatrix} \cdot X_i^{t-1} + \begin{bmatrix} \mathbf{H}_i \\ \mathbf{1}_{k \times 1} \end{bmatrix} \circ \Omega_i + \begin{bmatrix} \mathbf{0}_{k \times 1} \\ \mathbf{P}_i \end{bmatrix} \begin{bmatrix} \mathbf{0}_{k \times 1} \\ Z_i^t \end{bmatrix} + \begin{bmatrix} \mathbf{W}_i \\ \mathbf{0}_{k \times 1} \end{bmatrix} \tag{12}$$

Here, $\mathbf{F_i^1}, \mathbf{F_i^2} \in \mathbb{R}^{k \times k}$ are the state transition patterns for dynamics of the system, $\Omega_i \in \mathbb{R}^{k \times 1}$ is the dataset anomaly matrix for identifying each client's dataset, $\mathbf{H}_i \in \mathbb{R}^{k \times 1}$ and $\mathbf{P}_i \in \mathbb{R}^{k \times k}$ are the control input matrices to model the effect of the control input $Z_i^t$ on the state, and $\mathbf{W}_i$ is the process noise reflecting the uncertainty in the system dynamics.

To fit these parameters for each client, we accumulate three rounds of updates and global models for all clients and compute the corresponding $\{X_i^t\}$ and $\{Z_i^t\}$. We then fit the parameters $\mathbf{F_i^1}, \mathbf{F_i^2}, \Omega_i, \mathbf{H}_i, \mathbf{P}_i$, and $\mathbf{W}_i$ for each client $i$ using a one-layer fully connected neural network.

In the case where the client's data distribution changes slowly over time, we need to update the dataset parameter $\Omega_i$ each round and leave the other parameters unchanged.

**Domain State Prediction and Detection:** After obtaining the above parameters, given Eq.(10), we can predict the domain state of client $i$ at round $t$ with the input $D_i^{t-1}$ and $G^{t-1}$:

$$\hat{D}_i^t = \mathbf{F_i^1} \cdot D_i^{t-1} + \mathbf{F_i^2} \cdot G^{t-1} + \mathbf{H}_i \circ \Omega_i + \mathbf{W}_i. \quad (13)$$

The difference between the predicted domain state ($\hat{D}_i^t$) and the actual domain state ($D_i^t$) represents the change in the domain of client $i$, which can be expressed as

$$\Delta_i^t = D_i^t - \hat{D}_i^t = \{\delta_i^{1,t}, \cdots, \delta_i^{k,t}\} \quad (14)$$

where $\delta_i^{\kappa,t} = \beta_i^{\kappa,t} - \hat{\beta}_i^{\kappa,t}$ represents the domain state anomaly of client $i$ in the layer $\kappa$. We summarize four conditions based on § IV-C to select abnormal $\Delta_i^t$, as described in the DOMAIN-DETECT function of Alg. 1. First, the slope of the line connecting from $\delta_i^{1,t}$ to $\delta_i^{k,t}$ is greater than 0 (see Fig. 6), which means that the malicious client's domain changes become larger as the layer gets deeper. This condition is consistent with our Observation II and similar to the condition in line 17 of Alg. 1. Second, the MAD value (using Eq.(5)) for the slope is abnormal, and $\lambda$ is equal to that in line 6 of Alg. 1. Third, the MAD value for the domain change of any layer is abnormal. Fourth, the average of $\delta_i^{\kappa,t}$ is greater than 0, meaning that the actual domain state for the malicious client is always greater than the predicted state. This is because the malicious client's training data changes drastically compared to before.

### D. Results for Attack Tracing

After receiving the detection results of four features, the joint decision module identifies clients as malicious if any feature detects anomalies, and determines the remaining clients as benign. Then, *FLTracer determines the launch time of the attack and utilizes the detection details of each feature to trace the attack's objective, type, and poisoned portions of updates.*

The type and objective of attacks are discerned as follows: anomalies in `signv`, `sortv`, `classv`, and `featv` (CAM) correspond to sign-flipping, adaptive untargeted, dirty label, and backdoor attacks, respectively. The anomalies in the first three vectors indicate a reduction in the performance of the global model (untargeted), while in `featv`, they indicate the injection of a backdoor into the model (backdoor). We evaluate the percentage of directly traced attack objectives and types, with FLTracer achieving a success rate of $93.75\%$ and $80\%$ for tracing untargeted and backdoor attacks, as well as $84.38\%$, $78.00\%$, $95.00\%$, and $80.00\%$ for tracing sign-flipping, adaptive untargeted, dirty label, and backdoor attacks, respectively. Next, we analyze the poisoned portions of updates in the attacks. Any anomalies in `signv`, `sortv`, `classv`, and `featv` indicate that the opponent altered the symbols, sorting, classifier, and feature extractor, respectively. Specifically, in `featv`, certain layers exhibit a large Tsim, which implies the possible embedding of a backdoor within those layers. In situations where there are multiple adversaries in FL, utilizing attack tracing can identify colluding clients launching similar attacks.

## VI. EVALUATION

### A. Experimental Setup

**Datasets and Data distribution:** We evaluate FLTracer on multiple datasets for different learning tasks, including three

Table IV: Specific attacks in attack assessment.

| Untargeted attack type / Specific attack | | Targeted attack type / Specific attack | |
|---|---|---|---|
| Add noise [18] | Add noise | BadNets [21] | Patch-BN, Noise[23]-BN |
| Sign-flipping[19] | Sign-flipping | Model[7] replacement | Patch-MRA, Noise-MRA |
| Dirty label [20] | Fix-Fix, Fix-Rnd, Rnd-Fix, Rnd-Rnd | Distributed backdoor[11] | Patch-DBA, DBA-MRA |
| MB attack* [5] | MB MKrum, MB Bulyan, MB Median, MB Trmean, MB Max, MB Sum | Blind [8] backdoor* | Patch-Blind, Noise-Blind |
| Fang attack* [4] | Fang Mkrum, Fang Trmean | Feature replacement* | FRA combine with the above eight attacks |

"*" denotes untargeted adaptive or adaptive backdoor attacks.

image classification datasets (MNIST [35], EMNIST [36], and CIFAR10 [37]), a traffic sign recognition dataset (GT-SRB [38]), a human activity recognition dataset (HAR [39]), and a real-world driving dataset (BDD100K [40]). HAR is a 6-class class-imbalanced signal dataset representing human activity collected from smartphones of 30 real-world users. BDD100K is a large-scale driving video dataset of 100K real-world videos collected from diverse locations. We use its 100K image dataset, each of which is sampled at the 10th second from the video. We perform the 10-class road object classification task on these class-imbalanced data, where the distribution of the data changes slowly w.r.t weather conditions, scene types, and different times of the day.

We generate both IID and non-IID data distributions for each dataset, as listed in Table X (column 6). To simulate real non-IID settings, we use a Dirichlet distribution [41] with the degree of non-IID ($d$) set to $0.5$, following the setup in [6], [7], [28], [11]. We utilize the code in [7], [8] to distribute both the labels and data of the dataset, as depicted in Fig. 14. Note that a smaller value of $d$ indicates more heterogeneous data, and we compare the impact of non-IID degrees with one class per client, $d = 0.5$, $d = 5$ and $d \to \infty$. Since HAR is derived from real-world users, we do not need to distribute the data, and consider each user a client. The BDD100K dataset is partitioned into three subsets (i.e., daytime, dusk, and nighttime) and images in each part are further distributed using a $d = 0.5$. We emulate client training based on real-world changes in data distribution, starting with daytime data, then adding dusk data, and finally adding nighttime data.

**Model architectures:** We conduct experiments on 7 model architectures, as listed in Table X (column 7), including fully connected-based, convolution-based, and transformer-based, covering most of the popular model structures.

**Learning parameters:** We set up the standard training process using the AdamW/SGD optimizer for ViT/other model structures for local training, and the global learning rate is $\eta = 1/n$. Table X (columns 8-13) shows parameters. Following the setup in [5], [7], we launch untargeted attacks at the beginning of training and backdoor attacks in the 200th round when the model has converged. See Appendix B for more details.

**Evaluation Metrics:** First, we assess 14 untargeted attacks of 5 types and 16 backdoor attacks of 5 types, as listed in Table IV. Table V presents 9 metrics to access the performance of the normal task and backdoor task against attacks from three aspects: effectiveness, stability, and robustness. For the normal task, we measure the overall *accuracy* and *model stability* of the global model. The effectiveness of the attack is measured by *clean accuracy drop*. In addition, a set of category

Table V: Evaluation metrics

| Metrics | Definition |
|---|---|
| (E) Accuracy | # correct predictions / # total predictions |
| (S) Model stability | Standard deviation of accuracy |
| (E, D) Clean accuracy drop | # clean vanilla accuracy (w/o attack and detection) - # accuracy after the attack (and detection) |
| (E) Category accuracy$_c$ | # correct predictions of class $c$ / # total predictions of class $c$ |
| (E) Worst category accuracy | Minimum category accuracy$_c$ in all categories |
| (E) Best category accuracy | Maximum category accuracy$_c$ in all categories |
| (S) Category accuracy stability | Standard deviation of all category accuracies |
| (E, D) Attack success rate | # correct poisoned predictions / # total predictions |
| (R) Backdoor injection rounds | Number of rounds to reach 90% attack success rate |
| (R) Backdoor removal rounds | Number of rounds to remove backdoors (attack success rate=10%) when training with clean data |
| (D) True positive rate (TPR) | # correctly identified malicious clients / # total malicious clients |
| (D) False positive rate (FPR) | # benign clients incorrectly identified as malicious / # total benign updates |
| (D) Area under the curve | the area under the Receiver Operating Characteristic (ROC) curve |

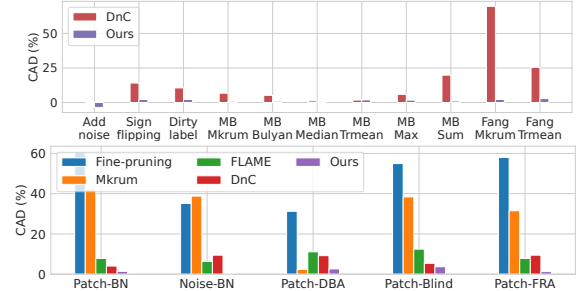(E)-Effectiveness; (S)-Stability; (R)-Robustness; (D)-Detection performance



Figure 8: Clean accuracy drop (CAD) of different methods on CIFAR10 (ResNet18). Top: against untargeted attacks. Bottom: against backdoor attacks. FLAME is unable to prevent backdoor injections against Patch-FRA even if only two clients remain for aggregation (attack success rate=99.97%).

metrics is introduced to further evaluate the effectiveness and stability of untargeted attacks. These metrics include *category accuracy*, *worst category accuracy*, *best category accuracy*, and *category accuracy stability*. For the backdoor task, the *attack success rate* is used to assess the effectiveness of backdoor attacks. Additionally, the *backdoor injection rounds* and *backdoor removal rounds* are utilized to evaluate the robustness of backdoor attacks.

Second, to evaluate the detection performance, we select representative attacks (11 untargeted and 5 backdoor attacks). We assess detection effectiveness based on three aspects: detection results, performance impact, and backdoor-proofing of the global model. We use TPR and FPR for detection results. We use *clean accuracy drop* to measure the accuracy decrease of the global model after attack and detection when compared to the clean model. For the highly imbalanced BDD100K dataset (where the number of images in the largest class is $100\times$ that of the smallest class), we use *area under the curve* (AUC) to show model accuracy. The *attack success rate* is used to assess the effectiveness of backdoor detection.

**Baselines:** We compare FLTracer with four SOTA methods, including all three categories in § VII: Byzantine-robust aggregation, update anomaly detection, and model sanitization. 1) *MKrum* [12] detect untargeted and backdoor attacks. Its initial strategy selects $n-2m-2$ clients for aggregation, but we reduce its high FPR by selecting $n-m-2$ clients. 2) *DnC* [5] detect untargeted and backdoor attacks. We adjust the tolerance rate $c$ to detect the maximum number of malicious clients. 3) *FLAME* [14] is designed to detect backdoor attacks, and we adopt its parameters for detection. 4) *Fine-pruning* [42] is a defense method used to mitigate the effectiveness of backdoor attacks. We choose $10\%$ clean samples from the test dataset randomly, then prune $50\%-70\%$ neurons with a $5\%$ change interval. We examine the effectiveness of the Differential Privacy (DP)-based method to counter backdoor attacks in Appendix D. The results show that the adversary can alter their attack strategy to successfully inject a backdoor.

Results for non-IID settings are reported here, while results for IID settings can be found in Appendix G. The runtime of FLTracer is analyzed in Appendix E. Our method's TPR and FPR are instantaneous, allowing erroneously classified clients to participate in subsequent training rounds.

## B. Assessment of Previous Attacks and Detection

We assess the effectiveness, stability, and robustness of untargeted and backdoor attacks and detection methods with 9 metrics. Here, we present important remarks, and for detailed analysis and results, please refer to Appendix C.

**Untargeted attacks**, like Add noise, Sign-flipping, and Fix-Fix attacks, significantly impact the model's stability. **Untargeted adaptive attacks**, like MB and Fang attacks, have a significant impact on the model's accuracy and category accuracy. Out of 16 **backdoor attacks** (including adaptive attacks), backdoors can quickly be injected in approximately 30 rounds and persist for over 200 rounds. Injecting backdoors becomes faster and more successful as the global model converges. **Current detection methods** only achieve high accuracy in detecting a small fraction of specific attacks, but at the cost of a high FPR in non-IID settings.

## C. Performance of FLTracer

Based on the assessment results, we select effective and representative attacks for the detection evaluation, covering five types of untargeted and five types of backdoor attacks in Table IV. FLTracer achieves the design goals in § III-C.

**Accuracy against untargeted attacks:** Table VI displays the TPR and FPR of detection methods against untargeted attacks. The results show that, on average, our FLTracer achieves a significantly higher TPR ($63.42\%$ higher than MKrum and $21.05\%$ higher than DnC) with the lowest FPR ($29.54\%$ less than MKrum and $19.77\%$ less than DnC). In 10 training rounds, FLTracer uses data from seven fewer malicious clients and 26 more benign clients than DnC, improving global model accuracy. Fig. 8 (top) illustrates the global model's clean accuracy drop under untargeted attacks, indicating that FLTracer outperforms DnC with an average improvement of $13.5\%$, which is consistent with our analysis. The clean accuracy drop for FLTracer, even with a TPR as low as $75\%$ under the Dirty label attack on CIFAR10 (ResNet18), is only $2.2\%$. This is due to the Dirty label attacl's relatively ineffective (see Fig. 17). In addition, we observe that FLTracer outperforms the clean vanilla model with a clean accuracy drop of $-3.59\%$. We will explain the reason behind this improvement in no-attack settings. Moreover, FLTracer shows greater stability in comparison to MKrum and DnC, accurately detecting most

9

Table VI: Comparing TPR (%) and FPR (%) of MKrum, DnC, and Ours against untargeted attacks. The **bold** and blue mark the best and second performances, respectively.

| Dataset (Model) | Attack | MKrum [12] | | DnC [5] | | FLTracer (Ours) | |
|---|---|---|---|---|---|---|---|
| | | TPR | FPR | TPR | FPR | TPR | FPR |
| EMNIST (Simple-Net) | Add Noise | 88.50 | 27.88 | **100.0** | **0.00** | **100.0** | **0.00** |
| | Sign-flipping | 62.50 | 34.38 | 87.00 | 28.63 | **90.50** | **1.88** |
| | Dirty label (Fix-Fix) | 68.00 | 33.00 | 86.00 | 23.13 | **99.50** | **3.25** |
| | MB Mkrum | 25.00 | 31.25 | 92.00 | 7.00 | **99.00** | **0.25** |
| | MB Bulyan | 25.00 | 31.25 | **100.0** | **0.00** | **100.0** | **0.00** |
| | MB Median | 25.00 | 31.25 | **100.0** | **0.00** | **100.0** | **0.00** |
| | MB Trmean | 25.00 | 31.25 | **100.0** | **0.00** | 99.50 | 0.13 |
| | MB Max | 25.00 | 31.25 | **100.0** | **0.00** | **100.0** | 0.13 |
| | MB Sum | 24.75 | 31.31 | 93.00 | 1.75 | **98.00** | **0.13** |
| | Fang Mkrum | 23.50 | 31.63 | 73.00 | 21.13 | **89.00** | **2.63** |
| | Fang Trmean | 25.00 | 31.25 | **100.0** | **0.13** | 88.00 | 4.25 |
| CIFAR10 (AlexNet) | Add Noise | 83.50 | 29.13 | **100.0** | **0.00** | **100.0** | 0.13 |
| | Sign-flipping | 53.00 | 36.75 | 72.00 | 33.25 | **100.0** | **0.00** |
| | Dirty label (Fix-Fix) | 42.50 | 39.38 | 70.25 | 32.63 | 67.50 | 7.81 |
| | MB Mkrum | 11.50 | 34.63 | 57.25 | 45.31 | **92.00** | **4.50** |
| | MB Bulyan | 12.50 | 34.38 | 70.25 | 34.94 | **96.00** | **4.81** |
| | MB Median | 25.00 | 31.25 | **100.0** | **0.00** | **100.0** | 0.38 |
| | MB Trmean | 25.00 | 31.25 | **100.0** | **0.00** | **100.0** | 0.13 |
| | MB Max | 24.25 | 31.44 | 82.25 | 30.38 | **95.00** | **7.00** |
| | MB Sum | 5.250 | 36.19 | 29.00 | 45.31 | **94.00** | **5.13** |
| | Fang Mkrum | 1.000 | 37.25 | 2.50 | 50.88 | **100.0** | **5.38** |
| | Fang Trmean | 2.000 | 37.00 | 84.41 | 30.96 | **98.00** | **3.13** |
| CIFAR10 (ResNet18) | Add Noise | 71.50 | 32.13 | 93.50 | 10.00 | **100.0** | **0.63** |
| | Sign-flipping | 61.00 | 34.75 | 65.50 | 39.75 | **82.50** | **6.00** |
| | Dirty label (Fix-Fix) | **82.00** | **14.25** | 52.30 | 45.28 | 75.00 | 6.56 |
| | MB Mkrum | 15.75 | 33.56 | 55.75 | 43.94 | **100.0** | **5.81** |
| | MB Bulyan | 17.25 | 33.19 | 39.75 | 40.69 | **94.00** | **7.13** |
| | MB Median | 24.75 | 31.31 | **100.0** | 0.06 | **100.0** | **2.38** |
| | MB Trmean | 25.00 | 31.25 | **100.0** | **0.00** | 100.0 | 4.19 |
| | MB Max | 14.00 | 34.00 | 31.50 | 41.50 | **100.0** | **6.31** |
| | MB Sum | 4.250 | 36.44 | 31.75 | 51.81 | **99.00** | **2.00** |
| | Fang Mkrum | 16.00 | 33.50 | 48.44 | 39.30 | **100.0** | **3.00** |
| | Fang Trmean | 15.75 | 33.56 | 30.83 | 52.08 | **86.50** | **2.44** |
| **Average** | | 31.82 | 32.49 | 74.19 | 22.72 | **95.24** | **2.95** |

Table VII: Comparing TPR (%) and FPR (%) of MKrum, FLAME, DnC, and Ours against backdoor attacks. "BN", "DBA", and "FRA" denote the BadNets, distributed backdoor, and feature replacement attack, respectively. The **bold** and blue mark the best and second performances, respectively.

| Dataset (Model) | Attack | MKrum [12] | | FLAME [14] | | DnC [5] | | FLTracer (Ours) | |
|---|---|---|---|---|---|---|---|---|---|
| | | TPR | FPR | TPR | FPR | TPR | FPR | TPR | FPR |
| CIFAR10 (ResNet18) | Patch-BN | 91.00 | 23.22 | 67.00 | 31.11 | **100.0** | **0.22** | **100.0** | 2.33 |
| | Noise-BN | 92.00 | 23.11 | 54.00 | 30.56 | 80.00 | 9.11 | **100.0** | **3.89** |
| | Patch-DBA | 80.50 | 29.88 | 54.50 | 24.22 | 94.00 | 11.13 | **100.0** | **2.13** |
| | Patch-Blind | 73.00 | 25.22 | 40.00 | 31.33 | 50.00 | 16.67 | **99.00** | **6.78** |
| | Patch-FRA | 75.00 | 25.00 | 35.00 | 29.67 | 55.00 | 28.22 | **100.0** | **4.78** |
| CIFAR10 (VGG16) | Patch-BN | 80.00 | 24.44 | 58.00 | 31.78 | 98.00 | 11.33 | **99.00** | **3.89** |
| | Noise-BN | 81.00 | 24.33 | 60.00 | 32.89 | 79.00 | 19.00 | **98.00** | **2.56** |
| | Patch-DBA | 80.50 | 29.88 | 50.00 | 28.11 | 95.50 | 8.38 | **100.0** | **0.11** |
| | Patch-Blind | 68.00 | 25.78 | 47.00 | 33.44 | 72.00 | 21.33 | **93.00** | **5.56** |
| | Patch-FRA | 75.00 | 25.00 | 51.00 | 33.67 | 52.00 | 22.00 | **99.00** | **3.44** |
| GTSRB (ResNet34) | Patch-BN | 80.00 | 24.44 | 44.00 | 29.33 | 93.00 | 24.56 | **100.0** | **0.00** |
| | Noise-BN | 81.00 | 24.33 | 49.00 | 27.22 | 92.00 | 24.67 | **100.0** | **0.00** |
| | Patch-DBA | 69.00 | 32.75 | 58.50 | 28.13 | 91.00 | 33.75 | **99.50** | **0.13** |
| | Patch-Blind | 75.00 | 25.00 | 42.00 | 29.22 | 99.00 | 21.56 | **100.0** | **0.00** |
| | Patch-FRA | 74.00 | 25.11 | 35.00 | 30.11 | 71.00 | 31.56 | **100.0** | **0.22** |
| HAR (DNN) | Patch-BN | 99.00 | 25.25 | 28.50 | 31.75 | 85.50 | 30.25 | **99.50** | **2.13** |
| | Patch-DBA | 97.50 | 25.63 | 53.50 | 27.00 | 83.50 | 34.75 | **99.00** | **3.38** |
| | Patch-Blind | 94.50 | 26.38 | 56.00 | 26.38 | 92.50 | 33.25 | **98.50** | **4.25** |
| | Patch-FRA | 98.50 | 25.38 | 91.00 | 20.38 | 81.50 | 32.38 | **100.0** | **4.63** |
| BDD100K (ViT) | Patch-BN | 90.00 | 23.33 | 35.00 | 30.56 | **100.0** | **0.00** | **100.0** | 0.22 |
| | Noise-BN | 83.00 | 24.11 | 39.00 | 30.67 | **100.0** | **0.00** | **100.0** | 0.11 |
| | Patch-DBA | 88.00 | 23.56 | 34.00 | 26.25 | 96.50 | 25.88 | **99.50** | **3.75** |
| | Patch-Blind | 87.00 | 23.67 | 31.00 | 29.89 | **100.0** | **0.00** | 97.00 | 0.33 |
| | Patch-FRA | 84.00 | 24.00 | 31.00 | 30.56 | **100.0** | **0.00** | 98.00 | 0.33 |
| **Average** | | 83.19 | 25.37 | 47.67 | 29.34 | 85.88 | 18.33 | **99.13** | **2.29** |

Results against MRA are not reported because the malicious updates scale up to 10×, so extreme anomalies are easily captured for almost all detections.
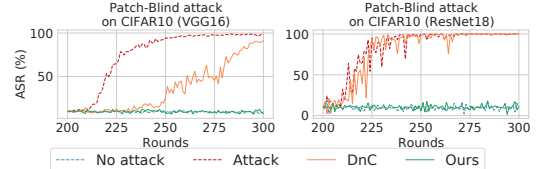


Figure 9: Attack success rate (ASR) of the different attacks.

settings across two datasets, three models, and 11 attacks. However, DnC can only achieve precise detection on simple datasets/models, like EMNIST/SimpleNet. It exhibits a high FPR on CIFAR10. This is due to the heightened sensitivity of our sortv, signv, and classv in comparison to the updates concatenation in Mkrum and DnC. Thus FLTracer can handle stealthy modifications even on complex models.

**Accuracy against backdoor attacks:** Table VII displays the TPR and FPR of detections against backdoor attacks. Our FLTracer outperforms three SOTA methods with the highest TPR and lowest FPR on CIFAR10, GTSRB, and HAR datasets, demonstrating its capability to accurately distinguish between malicious and benign clients even in complex datasets/models. Specifically, we achieve an average TPR of 99.13% at an extremely low FPR of 2.29%, which is 23.08% lower than MKrum, 27.05% lower than FLAME, and 16.04% lower than DnC. Per round, FLTracer detects the most malicious clients while using helpful data from two more benign clients out of 10. We report the lowest clean accuracy drop for each method that can prevent the backdoor injection (see Fig. 8 (bottom)). The results reveal that FLTracer's clean accuracy drop averages only 1.8%, which is 46.1%, 28.8%, 7.3%, and 5.7% lower than Fine-pruning, MKrum, FLAME, and DnC, respectively. This indicates that our FLTracer can prevent backdoor injection while maintaining global model accuracy.

**Accuracy under no-attack** are analyzed in Appendix F.

**Backdoor-proof:** Recall that backdoor accuracy can achieve 90% with an average of 36 backdoor injection rounds (see Fig. 21). This means at least 70 malicious clients must be detected (TPR > 70%) to reduce the attack success rate below 90%. Table VII shows that FLTracer can successfully address 25 backdoor injections in 25 settings while marginally impairing the global model accuracy. In contrast, other methods, such as FLAME, are unable to prevent Patch-FRA backdoor injections unless every update is excluded from the aggregation. The attack success rate of various attacks is shown in Fig. 9. FLTracer can successfully prevent backdoor injection even at the lowest TPR, particularly Patch-Blind attacks on CIFAR10 (VGG16), whereas DnC cannot.

### D. Impact of Different Data Distribution

**Different degrees of non-IID:** Table VIII displays the TPR and FPR of detections under varying degrees of non-IID. It indicates that the detection performance is worse as $d$ decreases due to a more heterogeneous data distribution and greater obscurity of the adversary. FLTracer surpasses the other three methods in attaining the highest TPR and FPR across

Table VIII: TPR (%) and FPR (%) of detection against Patch-BN attack on CIFAR10 (ResNet18) under different degrees of non-IID. "one class" represents each client has only one class. The **bold** and blue mark the best and second performances.

| Data distribution | MKrum [12] | | FLAME [14] | | DnC [5] | | FLTracer (Ours) | |
|---|---|---|---|---|---|---|---|---|
| | TPR | FPR | TPR | FPR | TPR | FPR | TPR | FPR |
| one class | 77.00 | 24.78 | 20.00 | 32.44 | 89.00 | 27.78 | **94.00** | **7.67** |
| $d = 0.5$ | 91.00 | 23.22 | 67.00 | 31.11 | 100.0 | 2.33 | 100.0 | 0.22 |
| $d = 5$ | 90.00 | 23.33 | 68.00 | 26.22 | 100.0 | 0.00 | 100.0 | 0.00 |
| $d \rightarrow +\infty$ (IID) | 88.00 | 23.56 | 96.00 | 26.56 | 100.0 | 0.00 | 100.0 | 0.00 |

all degrees of $d$. In particular, under the one-class setting, FLTracer achieves a TPR that is $5\%$ higher and an FPR that is $20\%$ lower compared to the second best method.

**Training data over Time:** While clients' data changes gradually, we detect and update the dataset parameter $\Omega$ of clients identified as benign per round, see § V-C. During rounds 0-200, every client normally trains using only daytime data. In rounds 200-300, each client adds the dusk data, while adversaries launch a backdoor attack. In rounds 300-400, each client adds nighttime data, while adversaries continue their backdoor attack. Fig. 10 shows the successful injection of the backdoor in a real-world scenario where the dataset changes gradually. Our detection effectively thwarts the injection of the backdoor and achieves an AUC value comparable to no-attack.
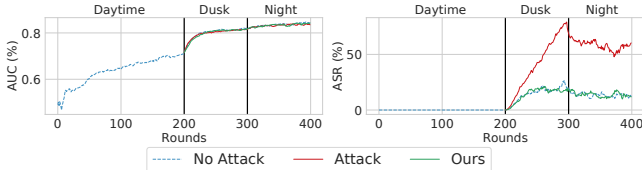


Figure 10: Area under the curve (AUC) of the main task and attack success rate (ASR) on BDD100K. Data distribution of clients change over time (daytime, dusk, and night). We launch Patch-BN attack after round 200.

### E. Impact of Selecting Different $\tau$

In FLTracer, we utilize a threshold $\tau \sim -1$ in TASKDETECT function (line 15 of Alg. 1) to identify malicious clients. Fig. 11 presents the TPR and FPR of FLTracer at different $\tau$. The results indicate that smaller $\tau$ corresponds to lower TPR and lower FPR, resulting in fewer malicious clients. However, the impact of $\tau$ on TPR and FPR is not significant, with a maximum change of $\leq 2\%$ for both TPR and FPR. Hence, the choice of $\tau$ can be adjusted based on the desired trade-off between accuracy and backdoor-proof (larger $\tau$ for higher accuracy and a smaller $\tau$ for enhanced backdoor-proof).
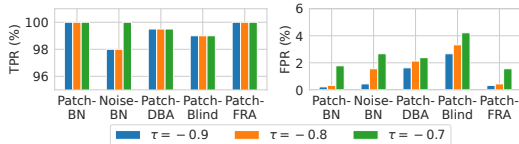


Figure 11: The impact of different $\tau$ on CIFAR10 (ResNet18).

### F. Impact of Percentage of Malicious Client

Fig. 12 compares FLTracer and DnC's accuracy and attack success rate for varying percentages of malicious clients. We select two untargeted and two backdoor attacks. For untargeted
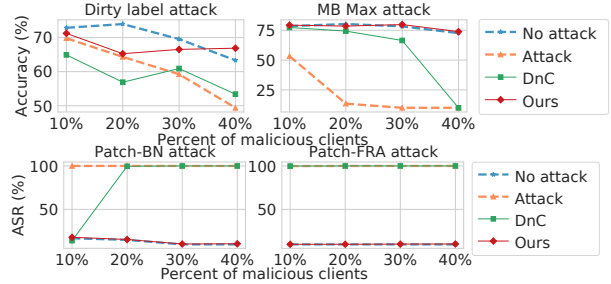


Figure 12: Impact of the percentage of malicious clients on accuracy and attack success rate (ASR) in CIFAR10 (ResNet18). Top-untargeted attacks. Bottom-targeted attacks.

attacks, DnC experiences a significant decrease in accuracy. However, this has minimal effect on FLTracer. Specifically, even with the worst detection result in the Dirty label attack, FLTracer experiences a smaller decline in accuracy than DnC ($2.45\%$ vs. $10.9\%$). In the MB Max attack, FLTracer can effectively prevent the accuracy drop, which sharply decreases under DnC, particularly with $40\%$ of malicious clients (accuracy: $10\%$). For backdoor attacks, backdoors are injected more rapidly as the number of malicious clients increases. In the Patch-BN attack, FLTracer is robust with varying numbers of malicious clients, whereas DnC is ineffective with more than $10\%$ of malicious clients. FLTracer also has a more stable detection effect against the Patch-FRA attack.

### G. Ablation Study

Table IX shows the impact of removing different components of FLTracer against untargeted attacks (Sign-flipping and Fang Mkrum) and backdoor attacks (Patch-BN and Patch-Blind). The results indicate that all three detections are crucial in enhancing effectiveness, specifically in achieving high TPR and low FPR. We observe that Local Anomaly Detection contributes the most to the high TPR in untargeted attack detection. Task Detection and Domain Detection effectively increase TPR and decrease FPR in backdoor attack detection.

Table IX: TPR/ FPR of FLTracer for removing certain detections on CIFAR10 (ResNet18). "J" denotes the joint decision-making with all detections.

| | Attacks | TPR(%) | FPR(%) | TPR-J(%) | FPR-J(%) |
|---|---|---|---|---|---|
| Remove Local Anomaly Detection | Sign-flipping | 14.00 | 3.25 | 82.50 | 6.00 |
| | Fang Mkrum | 34.50 | 2.81 | 100.0 | 3.00 |
| | Patch-BN | 100.0 | 2.33 | 100.0 | 2.33 |
| | Patch-Blind | 99.00 | 6.78 | 99.00 | 6.78 |
| Remove Task Detection | Sign-flipping | 82.50 | 6.00 | 82.50 | 6.00 |
| | Fang Mkrum | 87.00 | 1.19 | 100.0 | 3.00 |
| | Patch-BN | 97.00 | 11.00 | 100.0 | 2.33 |
| | Patch-Blind | 93.00 | 13.90 | 99.00 | 6.78 |
| Remove Domain Detection | Sign-flipping | 82.50 | 6.00 | 82.50 | 6.00 |
| | Fang Mkrum | 87.00 | 1.19 | 100.0 | 3.00 |
| | Patch-BN | 100.0 | 0.78 | 100.0 | 2.33 |
| | Patch-Blind | 100.0 | 14.70 | 99.00 | 6.78 |

### H. Advanced Adaptive Attacks

We assume that the adversary knows our FLTracer and designs advanced attacks to evade individual detection strategies. We assess their efficacy on CIFAR10 (ResNet18). Notably, combining two evasion strategies weakens the attack strength, as evading multiple detections entails adjusting fewer parameters (otherwise, it can be easily detected).

**Evade local anomaly detection:** We assume that the adversary knows all benign client updates to evade anomaly detection by `signv`, `sortv`, and `classv`. The adversary first trains normally to obtain benign updates $(u^0)$ and computes the original `signv`$^0$, `sortv`$^0$, and `classv`$^0$ based on all other benign updates. Then, the adversary modifies the update to the maximum extent possible while ensuring these three features remain unchanged to counter detection. The attack against our FLTracer only cause a clean accuracy drop of $0.09\%$ as the limited modification of updates (see Fig. 23 (a)).

**Evade task and domain detection:** A direct advanced backdoor attack is to refine the benign model with frozen feature extractor layers. For example, the adversary replaces only the bias and classifier with the backdoor model. This attempted backdoor injection almost fails, despite many rounds of training for the bias and classifier. The attack success rate is only $9.48\%$, and the clean accuracy drops by $1.31\%$ (see Fig.23(b)).

**Evade task detection:** We assume the adversary knows all benign clients' updates. The adversary can selectively retain the updates of certain layers in a backdoor model to reduce alterations of updates. The adversary first trains a backdoor model and calculates TSim $(\alpha^\kappa)$ for each layer. If $\alpha^\kappa \leq 0$, the attack on layer $\kappa$ is revoked. Additionally, if the adjusted TSim $(\alpha_{new}^\kappa)$ shows a downward trend, the attack on all convolutional layers at the current round is revoked. Under FLTracer, after 100 rounds of training, the attack's success rate is $27.09\%$ and clean accuracy drop is $0.52\%$ (see Fig.23(c)).

**Evade domain detection:** The adversary can launch the backdoor attack at the first round of training to counter the domain detection. However, injecting the backdoor into the global model while simultaneously converging the global model is difficult when its accuracy is low. Moreover, at the beginning of training, benign clients provide more general patterns, resulting in adversary updates being more anomalous. Thus, our Task Detection is sufficient to identify the adversary accurately, and the malicious candidates by Task Detection can be directly considered adversaries, i.e., $M_d = M_c \cup M_d$. This attack achieves an attack success rate of only $13.99\%$ with a clean accuracy drop of $2.41\%$, as shown in Fig. 23 (d).

**Evading a combination** of strategies one, three, and four leads to a clean accuracy drop of $0.301\%$ for new untargeted attack; an attack success rate of $10.31\%$, and a clean accuracy drop of $2.172\%$ for new backdoor attack, fails to inject backdoors.

## VII. Related Work

**Byzantine-robust Aggregation Algorithms:** give provable convergence guarantees for each round of FL aggregation against poisoning attacks caused by Byzantine client failures [43], [44], [45]. Blanchard et al. [12] proposed Krum and Multi-Krum (MKrum), which assume IID training data across clients and have knowledge of the number of malicious clients $(m)$. Krum selects one update with the minimum Euclidean distance to its other $n - m - 2$ closet updates as benign ones, while MKrum iteratively selects $n - 2m - 2$ updates with the minimum Euclidean distance from its closet update as benign. Yin et al. [46] proposed Median, which aggregates all updates by computing the median of the updates in all the dimensions.

**Updates Anomaly Detection:** aims to detect and remove malicious updates caused by poisoning attacks in each round of FL [47], [48], [49]. Fung et al. [50] proposed FoolsGold, which assumes IID training data for malicious clients and non-IID data for benign clients. Cao et al.[28] proposed FLTrust, which assumes an extra clean training dataset on the server and assigns lower trust scores to updates that have anomaly similarity with the clean updates trained by the server. Shejwalkar et al. [5] proposed DnC, which assumes that the number of malicious clients is known, computes the anomaly of each update by the singular value decomposition, and then selects the top updates with the highest anomaly as malicious ones. Nguyen et al. [14] proposed FLAME, which selects benign updates by calculating the cosine similarity and clustering their similarities, then aggregates the clipped updates to mitigate backdoors. Zhang et al. [16] and Cao et al. [17] proposed FLDetector and FedRecover, both of which assume that all clients participate in FL training in each round, and use the stored historical information of each client to estimate and recover their model update in each round. These methods assume that the data of benign clients is constant; otherwise, the updates of benign clients are also considered adversaries and are recovered to the updates before the data was changed.

**Model Sanitization:** is developed on centralized training to detect if models have backdoors and sanitize them by retraining [51], [52], [53], [22], [54], [55], [56], [57]. Liu el at. [42] proposed fine-pruning, which mitigates the backdoor by pruning neurons with low activation for clean samples and then fine-tunes the pruned model with clean samples to recover the model accuracy. Xu et al. [58] proposed MNTD, which trains a binary meta-classifier to predict whether a given model has backdoors. The meta-classifier is trained using a set of clean and backdoor shadow models trained on the same tasks as the given model. Another type of model sanitization is based on DP. Naseri et al. [59] showed that the DP mechanism could mitigate the backdoors by bounding updates and adding perturbations. However, our study in Appendix D shows that the adversary can successfully inject backdoors in a DP-based FL system by adjusting the attack strategy.

## VIII. Conclusion

In this paper, we comprehensively study the prior FL attacks and detection methods and then propose the first attack provenance framework, FLTracer. It can accurately detect various attacks and provide attack tracing by exploiting new features of various attacks. We develop task and Kalman filter-based domain detection to reduce false positives caused by data heterogeneity. Our extensive evaluations show that FLTracer has a high TPR and low FPR against various attacks and successfully addresses backdoors while only minimally affecting the global model in both IID and non-IID settings.

## References

[1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *AISTATS*, 2017.

[2] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv*, 2018.

[3] "Designing for privacy - wwdc19 - videos - apple developer," https://developer.apple.com/videos/play/wwdc2019/708, 2019.

[4] M. Fang, X. Cao, J. Jia, and N. Gong, "Local model poisoning attacks to byzantine-robust federated learning," in *USENIX*, 2020.

[5] V. Shejwalkar and A. Houmansadr, "Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning," in *NDSS*, 2021.

[6] V. Shejwalkar, A. Houmansadr, P. Kairouz, and D. Ramage, "Back to the drawing board: A critical evaluation of poisoning attacks on production federated learning," in *SP*, 2022.

[7] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *ICANN*, 2020.

[8] E. Bagdasaryan and V. Shmatikov, "Blind backdoors in deep learning models," in *USENIX*, 2021.

[9] Z. Zhang, A. Panda, L. Song, Y. Yang, M. Mahoney, P. Mittal, R. Kannan, and J. Gonzalez, "Neurotoxin: durable backdoors in federated learning," in *ICML*, 2022.

[10] X. Qi, T. Xie, R. Pan, J. Zhu, Y. Yang, and K. Bu, "Towards practical deployment-stage backdoor attack on deep neural networks," in *CVPR*, 2022.

[11] C. Xie, K. Huang, P.-Y. Chen, and B. Li, "Dba: Distributed backdoor attacks against federated learning," in *ICLR*, 2019.

[12] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," in *NeurIPS*, 2017.

[13] H. Zhu, J. Xu, S. Liu, and Y. Jin, "Federated learning on non-iid data: A survey," *Neurocomputing*, 2021.

[14] T. D. Nguyen, P. Rieger, H. Chen, H. Yalame, H. Möllering, H. Fereidooni, S. Marchal, M. Miettinen, A. Mirhoseini, S. Zeitouni *et al.*, "Flame: Taming backdoors in federated learning," in *USENIX*, 2022.

[15] A. Malki, E.-S. Atlam, and I. Gad, "Machine learning approach of detecting anomalies and forecasting time-series of iot devices," *Alexandria Engineering Journal*, 2022.

[16] Z. Zhang, X. Cao, J. Jia, and N. Z. Gong, "Fldetector: Defending federated learning against model poisoning attacks via detecting malicious clients," in *SIGKDD*, 2022.

[17] X. Cao, J. Jia, Z. Zhang, and N. Z. Gong, "Fedrecover: Recovering from poisoning attacks in federated learning using historical information," in *SP*, 2023.

[18] Z. Wu, Q. Ling, T. Chen, and G. B. Giannakis, "Federated variance-reduced stochastic gradient descent with robustness to byzantine attacks," *IEEE Transactions on Signal Processing*, 2020.

[19] L. Li, W. Xu, T. Chen, G. B. Giannakis, and Q. Ling, "Rsa: Byzantine-robust stochastic aggregation methods for distributed learning from heterogeneous datasets," in *AAAI*, 2019.

[20] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *ICML*, 2012.

[21] T. Gu, B. Dolan-Gavitt, and S. Garg, "Badnets: Identifying vulnerabilities in the machine learning model supply chain," *arXiv*, 2017.

[22] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, "Trojaning attack on neural networks," in *NDSS*, 2018.

[23] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, "Targeted backdoor attacks on deep learning systems using data poisoning," *arXiv*, 2017.

[24] M. Barni, K. Kallas, and B. Tondi, "A new backdoor attack in cnns by training set corruption without label poisoning," in *ICIP*, 2019.

[25] Y. Liu, X. Ma, J. Bailey, and F. Lu, "Reflection backdoor: A natural backdoor attack on deep neural networks," in *ECCV*, 2020.

[26] G. Baruch, M. Baruch, and Y. Goldberg, "A little is enough: Circumventing defenses for distributed learning," *NeurIPS*, 2019.

[27] R. E. Kalman, "A new approach to linear filtering and prediction problems," 1960.

[28] X. Cao, M. Fang, J. Liu, and N. Z. Gong, "Fltrust: Byzantine-robust federated learning via trust bootstrapping," in *NDSS*, 2021.

[29] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for long-tailed recognition," in *ICLR*, 2020.

[30] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021.

[31] S. Wold, K. Esbensen, and P. Geladi, "Principal component analysis," *Chemometrics and intelligent laboratory systems*, vol. 2, no. 1-3, pp. 37–52, 1987.

[32] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *ECCV*, 2014.

[33] A. Egoyan, "Can someone differentiate between cosine, adjusted cosine, and pearson correlation similarity measuring techniques," 2017.

[34] P. J. Rousseeuw and C. Croux, "Alternatives to the median absolute deviation," *Journal of the American Statistical association*, vol. 88, no. 424, pp. 1273–1283, 1993.

[35] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, 1998.

[36] G. Cohen, S. Afshar, J. Tapson, and A. Van Schaik, "Emnist: Extending mnist to handwritten letters," in *IJCNN*. IEEE, 2017.

[37] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[38] S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel, "Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark," in *IJCNN*, no. 1288, 2013.

[39] D. Anguita, A. Ghio, L. Oneto, X. Parra Perez, and J. L. Reyes Ortiz, "A public domain dataset for human activity recognition using smartphones," in *ESANN*, 2013.

[40] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving dataset for heterogeneous multitask learning," in *CVPR*, 2020.

[41] T. Minka, "Estimating a dirichlet distribution," 2000.

[42] K. Liu, B. Dolan-Gavitt, and S. Garg, "Fine-pruning: Defending against backdooring attacks on deep neural networks," in *RAID*, 2018.

[43] Y. Chen, L. Su, and J. Xu, "Distributed statistical machine learning in adversarial settings: Byzantine gradient descent," *POMACS*, 2017.

[44] H. Hashemi, Y. Wang, C. Guo, and M. Annavaram, "Byzantine-robust and privacy-preserving framework for fedml," *arXiv*, 2021.

[45] R. Guerraoui, S. Rouault *et al.*, "The hidden vulnerability of distributed learning in byzantium," in *ICML*, 2018.

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

[47] S. Awan, B. Luo, and F. Li, "Contra: Defending against poisoning attacks in federated learning," in *ESORICS*, 2021.

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

[49] P. Rieger, T. D. Nguyen, M. Miettinen, and A.-R. Sadeghi, "Deepsight: Mitigating backdoor attacks in federated learning through deep model inspection," in *NDSS*, 2022.

[50] C. Fung, C. J. Yoon, and I. Beschastnikh, "The limitations of federated learning in sybil settings," in *RAID*, 2020.

[51] Y. Li, X. Lyu, N. Koren, L. Lyu, B. Li, and X. Ma, "Neural attention distillation: Erasing backdoor triggers from deep neural networks," in *ICLR*, 2021.

[52] Y. Liu, W.-C. Lee, G. Tao, S. Ma, Y. Aafer, and X. Zhang, "Abs: Scanning neural networks for back-doors by artificial brain stimulation," in *CCS*, 2019.

[53] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao, "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks," in *SP*, 2019.

[54] G. Shen, Y. Liu, G. Tao, S. An, Q. Xu, S. Cheng, S. Ma, and X. Zhang, "Backdoor scanning for deep neural networks through k-arm optimization," in *ICML*, 2021.

[55] L. Truong, C. Jones, B. Hutchinson, A. August, B. Praggastis, R. Jasper, N. Nichols, and A. Tuor, "Systematic evaluation of backdoor data poisoning attacks on image classifiers," in *CVPR Workshops*, 2020.

[56] S. Andreina, G. A. Marson, H. Möllering, and G. Karame, "Baffle: Backdoor detection via feedback-based federated learning," in *ICDCS*, 2021.

[57] L. Zhao, S. Hu, Q. Wang, J. Jiang, C. Shen, X. Luo, and P. Hu, "Shielding collaborative learning: Mitigating poisoning attacks through client-side detection," *TDSC*, 2020.

[58] X. Xu, Q. Wang, H. Li, N. Borisov, C. A. Gunter, and B. Li, "Detecting ai trojans using meta neural analysis," in *SP*, 2021.

[59] M. Naseri, J. Hayes, and E. De Cristofaro, "Local and central differential privacy for robustness and privacy in federated learning," in *NDSS*, 2022.

[60] R. De Maesschalck, D. Jouan-Rimbaud, and D. L. Massart, "The mahalanobis distance," *Chemometrics and intelligent laboratory systems*, vol. 50, no. 1, pp. 1–18, 2000.

[61] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *NeurIPS*, 2012.

[62] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[63] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv*, 2014.

[64] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, "Learning differentially private recurrent language models," *arXiv*, 2017.

[65] H. Wang, K. Sreenivasan, S. Rajput, H. Vishwakarma, S. Agarwal, J.-y. Sohn, K. Lee, and D. Papailiopoulos, "Attack of the tails: Yes, you really can backdoor federated learning," *NeurIPS*, 2020.

## CAM EXTRACTION DETAILS

Fig. 13 depicts the process of CAM extraction. Here three clients are present, each having a model consisting of two convolutional layers. The first and second layers convolutional layers have sizes of $3\times3\times3\times3$ and $4\times4\times3\times3$, respectively. The convolution kernel's size is $L_1\times L_2 = 3\times3$.

**Step1.** We collect the first convolution kernel of each client and compute the anomaly score (i.e., $a_1$, $a_2$, and $a_3$) using the function SCORECONVKERNEL of Alg. 2.

**Step2.** (⊛ operation) We compute the anomaly scores for the convolution kernels of the same channels in layer one using the function SCORECONVKERNEL. After obtaining all anomaly scores for client 1, we combine them into an anomaly vector $A_1$. The anomaly vectors $A_2$ and $A_3$ for client 2 and client 3 can be obtained in a similar method.

**Step3.** After computing the anomaly matrices for all layers of client 1, we merge them into a convolution matrix CAM$_1$ using the function SCORECONVMATRIX of Alg. 2. Similarly, we can derive CAM$_2$ and CAM$_3$ for client 2 and client 3.
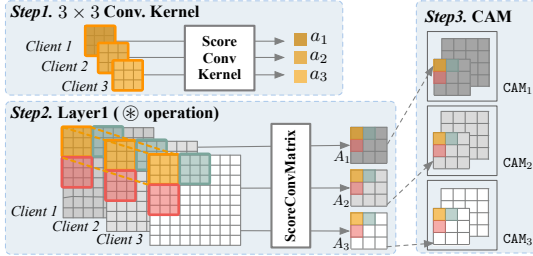


Figure 13: The CAM extraction for each client in a round.

Alg. 2 described functions used in CAM extraction. The function SCORECONVKERNEL computes anomaly scores for a set of convolution kernels ($w_{\{i\in[n]\}}$) in the same channel, where $u_i$ is a convolution kernel for client $i$. PCA and Mahalanobis Distance[60] are employed to obtain the outliers. The MAIN function computes CAM for $n$ clients. The loop in lines 7-9 computes the anomaly scores for $n$ clients in layer $\kappa$. The loop in lines 10-11 generates $A_i^\kappa$ for each client.

---

**Algorithm 2** SCORECONVMATRIX Algorithm

---

**Input:** $\theta_{\{i\in[n]\}}$: updates of $n$ clients; $k$: total number of conv. layers
**Output:** CAM$_{\{i\in[n]\}}$: convolution anomaly matrices of $n$ clients
1: **function** SCORECONVKERNEL($w_{\{i\in[n]\}}$):
2:     $w_{pca} \leftarrow$ PCA($u$, componenets = 2)
3:     $w_{md} \leftarrow$ MahalanobisDistance($u_{pca}$)
4:     $a_{\{i\in[n]\}} \leftarrow$ Normalize($u_{md}$)     **return** $a_{\{i\in[n]\}}$
5: **function** MAIN:
6:     **for** conv. layer $\kappa \in [k]$ **do**     ▷ the operation ⊛
7:         **for** each channel $\iota$ **do**
8:             $w_{\{i\in[n]\}}^\iota \leftarrow$ Conv. kernel set of channel $\iota$ for $n$ clients
9:             $a_{\{i\in[n]\}}^\iota \leftarrow$ SCORECONVKERNEL($w_{\{i\in[n]\}}^\iota$)
10:         **for** $i \in [n]$ **do**
11:             $A_i^\kappa \leftarrow$ Combine all $a_i^\iota$ of client $i$ in layer $\kappa$
12:     CAM$_i \leftarrow$ Generate using $A_i^\kappa$ and Eq.(1), $i \in [n]$
13:     **return** CAM$_{\{i\in[n]\}}$

---

## EXPERIMENTAL SETUP DETAILS

**Details on datasets and model architectures:** MNIST [35] is a 10-class class-balanced digit image classification dataset.
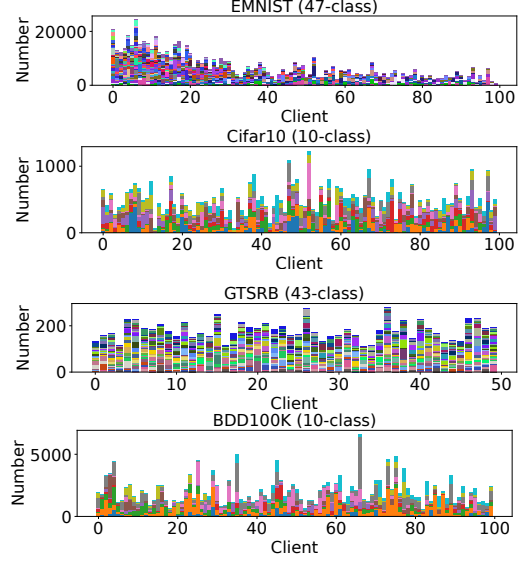


Figure 14: Data distribution of $N$ clients derived from the Dirichlet distribution with $d = 0.5$.

EMNIST [36] is a 47-class class-imbalanced digit image classification dataset. CIFAR10 [37] is a 10-class class-balanced color image classification dataset. German Traffic Sign Recognition Benchmark (GTSRB) [38] is a 43-class class-imbalanced traffic sign dataset with varying light conditions and rich backgrounds. For more details on datasets, please refer to Table X (columns 3-5).

Table X (column 7) lists the model architectures utilized for each dataset. SimpleNet comprises two convolution layers and two fully-connected layers. AlexNet [61], ResNet18 [62], VGG16 [63], and ResNet34 [62] are different architectures of convolution networks. DNN contains two fully connected layers. Vision Transformer (ViT) [30] is a Transformer based model for image recognition.

**Details on learning parameters:** Following the standard setup in [64], [7], we train local models with local learning rate $lr$, local epochs 2, and batch size 64. At each round, the $n$ selected clients train a local model to aggregate. For untargeted attacks, we start with a sustained attack from scratch of the training. We train MNIST and EMNIST with SimpleNet using $lr = 0.01$ for rounds 0-100. We train CIFAR10 with AlexNet and ResNet18 using $lr = 0.1$ for rounds 0-200. For backdoor attacks, we begin the attack when the global model is convergent, which is round 200 for CIFAR10, 100 for GTSRB (IID), and 200 for GTSRB (non-IID). We train CIFAR10 with ResNet18 and VGG16 using $lr = 0.1$ and $lr = 0.05$ for rounds 0-400. We train GTSRB with ResNet34 using $lr = 0.1$ for rounds 0-200 in the IID settings and $lr = 0.1$ for rounds 0-300 in the non-IID settings.

**Details on attacks:** For Add noise attack, we use Gaussian noise $\delta \sim \mathcal{N}(\mu, \sigma^2)$, where $\mu = 0$ and $\sigma = 0.3$. For Dirty label attacks, we increase the number of sources and target labels to enhance the attack's impact. Based on the different strategies applied to "source-target label" pairs, we evaluate four types of dirty label attacks, including the Fix-Fix attack, the Fix-Rnd attack, the Rnd-Fix attack, and the Rnd-Rnd attack. We set the proportion of poisoned samples for each malicious client

Table X: Datasets, model structures, and parameters

| | Dataset | Training input | Class | Input size | Data distribution | Model structure | $N$ | $n^*$ | $P_m^*$ | $P_p^*$ | Benign $lr$/epochs | Malicious $lr$/epochs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Un targeted | MNIST | 60000 | 10 | $28 \times 28$ | IID | SimpleNet | 100 | 10 | 20% | 0% | 0.01/2 | 0.01/2 |
| | EMNIST | 73168 | 47 | $28 \times 28$ | non-IID | SimpleNet | 100 | 10 | 20% | 0% | 0.01/2 | 0.01/2 |
| | CIFAR10 | 50000 | 10 | $32 \times 32$ | IID& non-IID | AlexNet | 100 | 10 | 20 | 0% | 0.1/2 | 0.1/2 |
| | CIFAR10 | 50000 | 10 | $32 \times 32$ | IID& non-IID | ResNet18 | 100 | 10 | 10%-40% | 0% | 0.1/2 | 0.1/2 |
| Targeted | CIFAR10 | 50000 | 10 | $32 \times 32$ | IID& non-IID | ResNet18 | 100 | 10 | 10%-40% | 3% | 0.01/2 | 0.005/4 |
| | CIFAR10 | 50000 | 10 | $32 \times 32$ | IID & non-IID | VGG100 | 100 | 10 | 10% | 3% | 0.01/2 | 0.005/4 |
| | GTSRB | 50000 | 43 | $32 \times 32$ | IID & non-IID | ResNet34 | 50 | 10 | 10% | 3% | 0.01/2 | 0.005/4 |
| | HAR | 7353 | 6 | $1 \times 561$ | non-IID | DNN | 21 | 10 | 10% | 3% | 0.01/2 | 0.005/4 |
| | BDD100K | 474706 | 10 | $64 \times 64$ | non-IID | ViT (Transformer) | 100 | 10 | 10% | 3% | 0.0001/4 | 0.00005/6 |

$^*n = 20$ for MB and Fang attacks, $P_m = 20\%$ for DBA, $P_p = 5\%$ for dirty label attack.
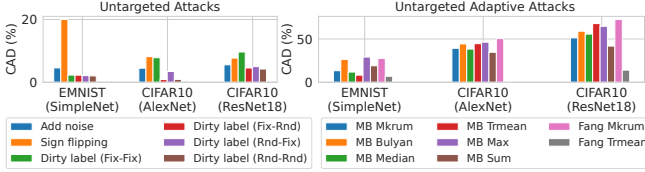In MRA, the adversary attacks for one round and scales the updates by $10\times$.



Figure 15: Clean accuracy drop (CAD) of untargeted attacks in non-IID settings. The CAD of the Sign-flipping attack on EMNIST is $81.28\%$.
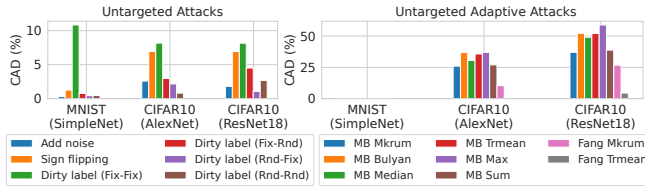


Figure 16: Clean accuracy drop (CAD) of untargeted attacks in IID settings. The CAD of untargeted adaptive attacks on MNIST is less than $0.2\%$.

to 50%. In the Fix-Fix attack, we mislabel the samples in categories 1 to 5 as category 3. In the Fix-Rnd attack, we mislabel the samples in categories 1 to 5 as random categories. In the Rnd-Fix attack, we mislabel half of the samples to category 3. In the Rnd-Rnd attack, we mislabel half of the samples into random categories. For MB and Fang attacks, we set $n = 20$ because these attacks require more than four malicious clients per round to generate malicious updates. For the backdoor attacks, we set $P_m = 10\%$ for most backdoor attacks and increase the $P_m = 20\%$ for DBA because of the distributed trigger and injection strategy. Following the setup in [7], [11], we attack at the beginning of the training for MRA and at round 200 for other backdoor attacks. In all attacks, adversaries modify their local $lr$/epochs to attack efficiently.

## APPENDIX C
## ASSESSMENT OF PREVIOUS ATTACKS AND DETECTION

**Untargeted attacks:** Fig. 15 shows the clean accuracy drop of untargeted and adaptive attacks in non-IID settings. Most untargeted attacks achieve a high clean accuracy drop, except for three dirty label attacks (Fix-Rnd to Rnd-Rnd). Specifically, the clean accuracy drop of MB Mkrum and Fix-Rnd on CIFAR10 (AlexNet) is $39.29\%$ and $0.85\%$. We also observe that the performance of untargeted adaptive attacks is generally stronger than that of untargeted attacks due to the ability of adaptive attacks to generate optimized perturbations. The average clean accuracy drop for adaptive and untargeted attacks is $25.8\%$ and $3.77\%$, respectively. (see Fig.16 for IID).
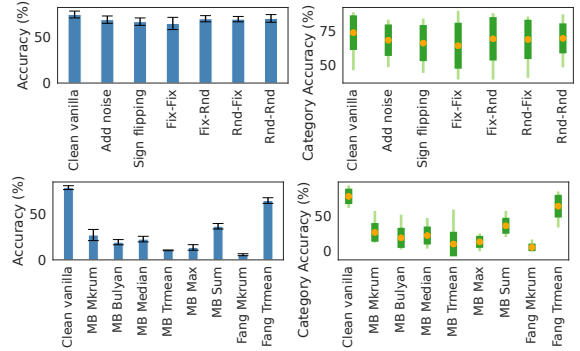


Figure 17: Accuracy and category accuracy of untargeted attacks compared to the clean vanilla on CIFAR10 (ResNet18). Top-untargeted attacks mainly affect stability. Bottom-untargeted adaptive attacks mainly affect accuracy.

Fig. 17 shows the accuracy, model stability, worst/best category accuracy, and category accuracy stability of untargeted and adaptive attacks. We observe that untargeted attacks mainly affect the stability of the model. Specifically, the model stability for the Fix-Fix attack is $4.15$, much higher than the clean vanilla model ($0.7$). In addition, untargeted adaptive attacks significantly affect the accuracy and category accuracy. The accuracy and category accuracy of the MB attack and the Fang attack are smaller than the baseline by $41.82\%$ and $41.15\%$, respectively. (see Fig.24 and 26 for other datasets).

> **Remark 1.** The 14 untargeted attacks (except for Fix-Rnd to Rnd-Rnd) lead to a significant decrease in the accuracy or stability of the global model.

**Targeted attacks:** Fig. 18 shows the attack success rate of 16 backdoor attacks in non-IID settings. Most attacks can inject backdoors with almost $100\%$ attack success rate, except for Patch-MRA and Noise-MRA on GTSRB. For different attack strategies, MRA is more difficult to inject backdoors because MRA only injects a single round. For different datasets, those with fewer categories and more training inputs are easier to inject backdoors. Specifically, the average attack success rate of $96.21\%$ for CIFAR10 is higher than that of $76.11\%$ for GTSRB. In addition, the effect of our new adaptive strategy, FRA, achieves similar results as without it. Specifically, the average attack success rate w/ and w/o FRA is $88.3\%$ and $89.51\%$, respectively. We also observe that the backdoor injection is unstable and slow at the beginning of the training, as shown in Fig. 20. Due to the low accuracy of the global model, it is difficult to converge the global model and inject backdoors into it simultaneously. (see Fig. 19 for IID).

**Remark 2.** When the global model converges, backdoors can be injected faster and more successfully, similar to the conclusion in [7].

Fig. 21 shows the backdoor injection rounds and backdoor removal rounds of backdoor attacks in non-IID settings. All attacks with and without Patch-FRA (except MRA) achieve an attack success rate of over $90\%$, with an average of only 36 and 39 backdoors injection rounds required. In almost all attacks (except MRA), the backdoor removal rounds exceed more than 200 rounds. In Patch-MRA, the adversary injects a backdoor for one round, and the global model takes more than 84 rounds to remove the backdoor. Such small backdoor injection rounds and large backdoor removal rounds imply that backdoor attacks are robust, posing a significant security risk to FL. (see Fig. 25 and Fig. 27 for other datasets).

**Remark 3.** The 16 backdoor attacks can rapidly inject backdoors and persist for much more training rounds.
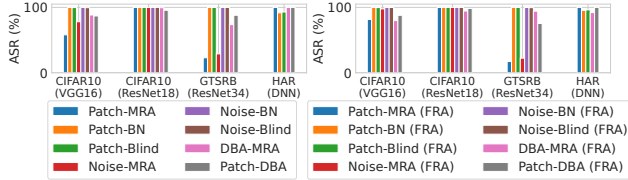


Figure 18: Attack success rate (ASR) of backdoor attacks in non-IID settings.
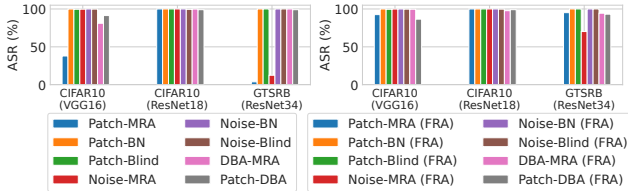


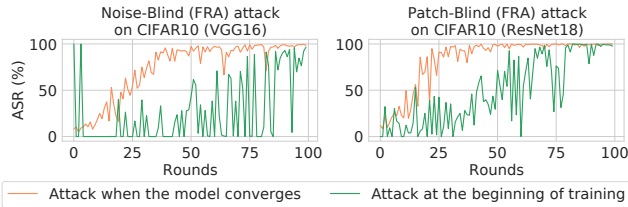Figure 19: Attack success rate (ASR) of backdoor attacks in IID settings.



Figure 20: Comparison of attack success rate between the attack when the global model converges and the attack at the beginning of training.
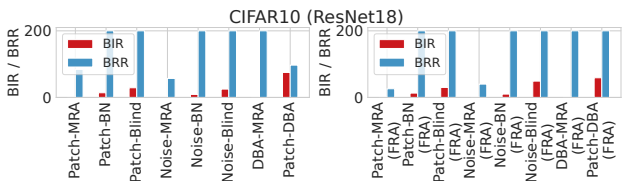


Figure 21: Backdoor injection rounds (BIR) and backdoor removal rounds (BRR) of backdoor attack in non-IID settings. A BIR of 200 indicates that more than 200 training rounds are required to remove the backdoor.

**Detection against attacks:** In non-IID settings, the effectiveness of existing methods against untargeted and backdoor attacks is unstable, as shown in Table VI (column MKrum

and DnC) and Table VII (column Mkrum, FLAME, and DnC). DnC is only effective against a few specific attacks, such as Add noise, MB Median, MB Trmean, and backdoor attacks under ViT, but is less effective against other attacks. The effectiveness of Mkrum is generally low under non-IID because Mkrum can only guarantee convergence of FL, which is not sufficient to detect malicious clients [7]. We also observe that existing methods are less effective on complex models and datasets (such as CIFAR10 and ResNet/VGG) because anomalies are more difficult to detect in complex settings.

**Remark 4.** Existing detection can only ensure high accuracy in detecting a small fraction of specific attacks, sacrificing high false positive rates (FPR) in non-IID settings.

## APPENDIX D
### DP-BASED DEFENSE AGAINST BACKDOOR ATTACKS

We evaluate the effectiveness of backdoor attacks launched under different DP strategies (see Fig. 22 (top)). We notice that adversaries can inject backdoors successfully through optimized attack strategies, such as increasing the adversary's local training epochs. Despite DP's inability to defend such optimized attacks, our method is still effective under DP. Fig. 22 (middle and bottom) show the effectiveness of our method against this attack. The results show that FLTracer can prevent backdoor injections (attack success rate$=10\%$) and the global model's accuracy can approximate the no-attack setting.
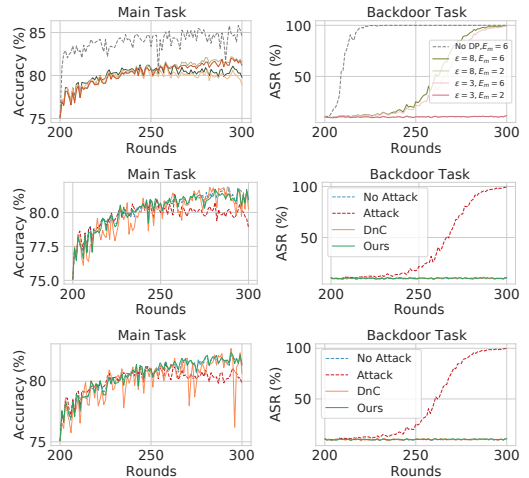


Figure 22: Accuracy and attack success rate (ASR) on CI-FAR10 (ResNet18) under DP-based method. Top-under DP ($E_m$: malicious epochs). Middle-w/ detection under DP with sensitivity $\epsilon = 3$. Bottom-w/ detection under DP with $\epsilon = 8$.

## APPENDIX E
### RUNTIME ANALYSIS

We offer a baseline runtime for attack provenance, which includes the detection time and attack tracing time of FLTracer. In non-IID settings, the average runtime of FLTracer is $8.68$ seconds per round (s/ r) for backdoor attacks on CIFAR10 (ResNet18). FLTracer has a relatively short attack tracing time overhead when compared to detection-only methods like MKRUM ($5.51$ s/ r), FLAME ($1.2$ s/ r), and DnC ($4.95$ s/ r).

Referring to the results in Table VII and Fig. 8, FLAME, despite its quick detection time, fails to prevent backdoor

injection, such as Patch-FRA. On the other hand, while both DnC and MKrum successfully prevent backdoor injection, they significantly slow down the training of the global model due to their high FPR. For instance, to prevent Patch-BadNets attacks, FLTracer only requires 37 training rounds to achieve the same level of accuracy as a global model that is protected by DnC with 100 training rounds. The global model accuracy of Mkrum is lower than that of DnC. In summary, FLTracer achieves the predefined global model accuracy in the shortest time while ensuring backdoor-proof protection.

## APPENDIX F
### PERFORMANCE OF FLTRACER UNDER NO ATTACK

Table XI lists the results of FLTracer under no attack in non-IID settings. It shows that FLTracer has little impact on the global model and increases its stability. In particular, FLTracer achieves a low $8\%$ FPR (under joint), which is close to the sum of the FPRs of the four features. This means that each feature detects malicious clients from different attack locations. The clean accuracy drop of FLTracer decreases by a total of $0.28\%$, and the model stability decreases by $0.4$. This is because without attacks, the updates of the clients with extremely heterogeneous data are significantly different from others, and aggregating these small fractions of updates decreases the global model performance. This observation matches the conclusion in [65]. (see Table XII for IID).

Table XI: Detection results under no attack in non-IID settings.

| | Baseline | signv | sortv | classv | featv | Joint |
|---|---|---|---|---|---|---|
| False positive rate | - | 1.80% | 3.30% | 3.00% | 1.70% | 8.00% |
| Accuracy | 84.38 | 84.97 | 84.99 | 84.49 | 84.37 | 84.65 |
| Clean accuracy drop | 0 | -0.59 | -0.61 | -0.11 | +0.01 | -0.28 |
| Model stability | 1.06 | 0.95 | 0.56 | 0.68 | 0.92 | 0.66 |

Table XII: Detection results under no attack in IID settings.

| | Baseline | signv | sortv | classv | featv | Joint |
|---|---|---|---|---|---|---|
| False positive rate | - | 0.20% | 3.80% | 0.60% | 1.00% | 4.50% |
| Accuracy | 89.76 | 89.86 | 89.84 | 90.04 | 89.86 | 89.85 |
| Clean accuracy drop | 0 | -0.10 | -0.08 | -0.27 | +0.10 | -0.09 |
| Model stability | 0.11 | 0.08 | 0.07 | 0.08 | 0.08 | 0.14 |

## APPENDIX G
### ADDITIONAL EXPERIMENTAL RESULTS

Table XIII: Comparing TPR(%) and FPR(%) of DnC and FLTracer (Ours) against backdoor attacks in IID settings.

| Dataset (Model) | Attack | DnC [5] | | FLTracer (Ours) | |
|---|---|---|---|---|---|
| | | TPR | FPR | TPR | FPR |
| CIFAR10 (ResNet18) | Patch-BadNets | 100.0 | 0.00 | 100.0 | 0.00 |
| | Noise-BadNets | 100.0 | 0.00 | 100.0 | 0.00 |
| | Patch-DBA | 98.50 | 4.13 | 100.0 | 0.00 |
| | Patch-Blind | 92.00 | 0.89 | 100.0 | 0.00 |
| | Patch-FRA | 100.0 | 0.00 | 100.0 | 0.11 |
| CIFAR10 (VGG16) | Patch-BadNets | 100.0 | 0.00 | 100.0 | 0.00 |
| | Noise-BadNets | 100.0 | 0.00 | 100.0 | 0.11 |
| | Patch-DBA | 100.0 | 0.00 | 100.0 | 0.00 |
| | Patch-Blind | 100.0 | 0.00 | 100.0 | 0.00 |
| | Patch-FRA | 100.0 | 0.00 | 100.0 | 0.00 |
| GTSRB (ResNet34) | Patch-BadNets | 100.0 | 0.00 | 100.0 | 0.00 |
| | Noise-BadNets | 100.0 | 0.00 | 100.0 | 0.11 |
| | Patch-DBA | 93.00 | 10.25 | 99.00 | 0.00 |
| | Patch-Blind | 100.0 | 0.11 | 100.0 | 0.00 |
| | Patch-FRA | 95.00 | 23.67 | 100.0 | 0.00 |
| **Average** | | 98.57 | 2.60 | **99.93** | **0.02** |

Table XIV: Comparing TPR(%) and FPR(%) of DnC, and FLTracer (Ours) against untargeted attacks in IID settings.

| Dataset (Model) | Attack | DnC [5] | | FLTracer (Ours) | |
|---|---|---|---|---|---|
| | | TPR | FPR | TPR | FPR |
| MNIST (SimpleNet) | Add Noise | 100.0 | 0.00 | 100.0 | 0.00 |
| | Sign-flipping | 74.00 | 10.13 | 100.0 | 0.00 |
| | Dirty label (Fix-Fix) | 98.50 | 0.38 | 98.00 | 3.63 |
| | MB Mkrum | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Bulyan | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Median | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Trmean | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Max | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Sum | 100.0 | 0.00 | 100.0 | 0.00 |
| | Fang Mkrum | 100.0 | 0.00 | 100.0 | 0.00 |
| | Fang Trmean | 100.0 | 0.00 | 100.0 | 0.00 |
| CIFAR10 (AlexNet) | Add Noise | 100.0 | 0.00 | 100.0 | 0.00 |
| | Sign-flipping | 99.50 | 22.50 | 100.0 | 0.00 |
| | Dirty label (Fix-Fix) | 100.0 | 1.13 | 100.0 | 2.94 |
| | MB Mkrum | 99.50 | 1.13 | 92.00 | 1.63 |
| | MB Bulyan | 98.50 | 1.88 | 92.25 | 2.69 |
| | MB Median | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Trmean | 100.0 | 0.00 | 96.00 | 0.50 |
| | MB Max | 99.75 | 1.50 | 92.00 | 2.25 |
| | MB Sum | 97.50 | 4.25 | 94.50 | 4.50 |
| | Fang Mkrum | 34.25 | 27.69 | 100.0 | 0.00 |
| | Fang Trmean | 98.25 | 2.25 | 100.0 | 0.00 |
| CIFAR10 (ResNet18) | Add Noise | 99.00 | 0.25 | 99.50 | 0.38 |
| | Sign-flipping | 68.00 | 23.75 | 100.0 | 0.13 |
| | Dirty label (Fix-Fix) | 99.49 | 9.38 | 100.0 | 1.56 |
| | MB Mkrum | 62.50 | 41.31 | 100.0 | 0.19 |
| | MB Bulyan | 51.75 | 45.94 | 99.25 | 0.19 |
| | MB Median | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Trmean | 100.0 | 0.00 | 100.0 | 0.00 |
| | MB Max | 50.25 | 42.25 | 100.0 | 0.06 |
| | MB Sum | 33.75 | 52.69 | 100.0 | 0.00 |
| | Fang Mkrum | 49.17 | 48.40 | 98.61 | 0.00 |
| | Fang Trmean | 19.25 | 38.63 | 100.0 | 0.00 |
| **Average** | | 85.85 | 11.38 | **98.85** | **0.63** |



(a) Evade local anomaly detection

(b) Evade task and domain detection

(c) Evade task detection

(d) Evade domain detection

Figure 23: Accuracy and attack success rate (ASR) of FLTracer against advanced adaptive attacks.

(a) MNIST (SimpleNet)

(b) CIFAR10 (AlexNet))

(c) CIFAR10 (ResNet18)

Figure 24: The accuracy and category accuracy of untargeted attacks compared with clean vanilla in IID settings.



(a) EMNIST (SimpleNet)

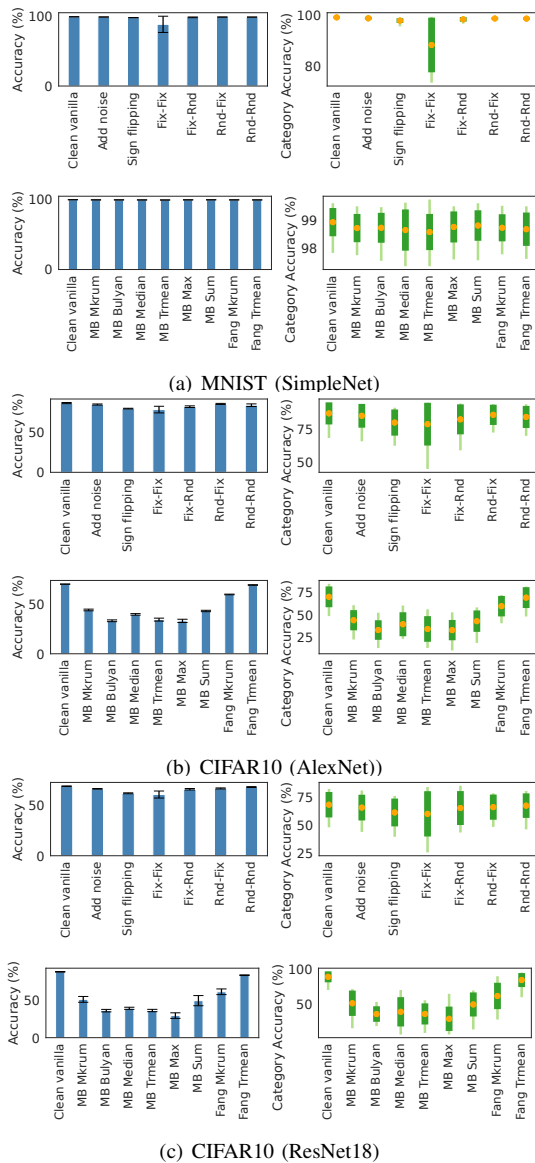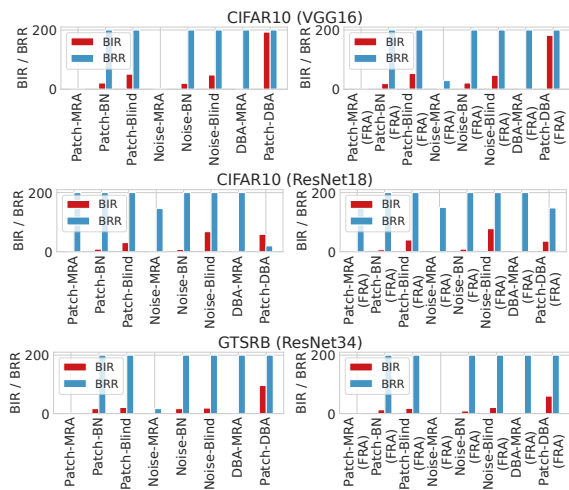(b) CIFAR10 (AlexNet)
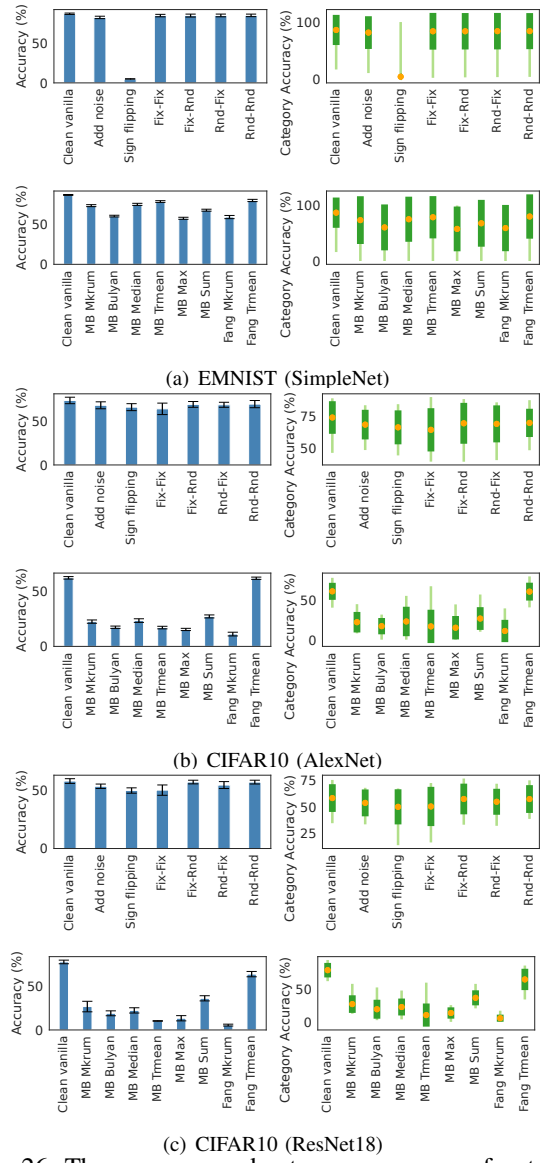
(c) CIFAR10 (ResNet18)

Figure 26: The accuracy and category accuracy of untargeted attacks compared with clean vanilla in non-IID settings.



Figure 25: The backdoor injection rounds (BIR) and backdoor removal rounds (BRR) of backdoor attacks in IID settings.
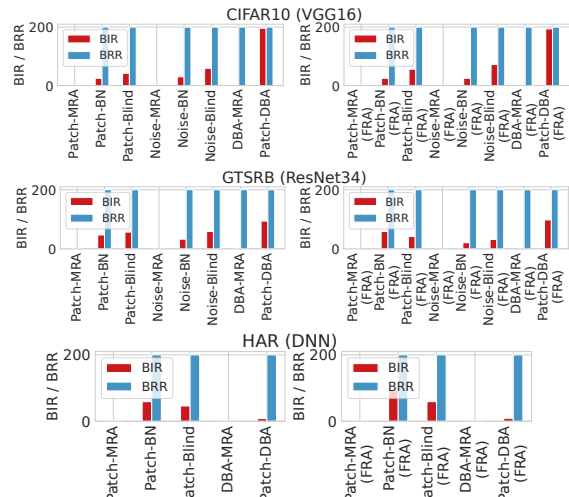


Figure 27: The backdoor injection rounds (BIR) and backdoor removal rounds (BRR) of backdoor attacks in non-IID settings.