

FLEDGE: Ledger-based Federated Learning Resilient to Inference and Backdoor Attacks

Jorge Castillo
jorge.a.castillo01@utrgv.edu
The University of Texas Rio Grande
Valley¹
USA

Phillip Rieger
phillip.rieger@trust.tu-darmstadt.de
Technical University of Darmstadt
Germany

Hossein Fereidooni
hossein.fereidooni@kobil.com
KOBIL GmbH²
Germany

Qian Chen
guenevereqian.chen@utsa.edu
The University of Texas at San
Antonio
USA

Ahmad-Reza Sadeghi
ahmad.sadeghi@trust.tu-
darmstadt.de
Technical University of Darmstadt
Germany

ABSTRACT

Federated learning (FL) is a distributed learning process that uses a trusted aggregation server to allow multiple parties (or clients) to collaboratively train a machine learning model without having them share their private data. Recent research, however, has demonstrated the effectiveness of inference and poisoning attacks on FL. Mitigating both attacks simultaneously is very challenging. State-of-the-art solutions have proposed the use of poisoning defenses with Secure Multi-Party Computation (SMPC) and/or Differential Privacy (DP). However, these techniques are not efficient and fail to address the malicious intent behind the attacks, i.e., adversaries (curious servers and/or compromised clients) seek to exploit a system for monetization purposes. To overcome these limitations, we present a ledger-based FL framework known as FLEDGE that allows making parties accountable for their behavior and achieve reasonable efficiency for mitigating inference and poisoning attacks. Our solution leverages crypto-currency to increase party accountability by penalizing malicious behavior and rewarding benign conduct. We conduct an extensive evaluation on four public datasets: Reddit, MNIST, Fashion-MNIST, and CIFAR-10. Our experimental results demonstrate that (1) FLEDGE provides strong privacy guarantees for model updates without sacrificing model utility; (2) FLEDGE can successfully mitigate different poisoning attacks without degrading the performance of the global model; and (3) FLEDGE offers unique reward mechanisms to promote benign behavior during model training and/or model aggregation.

CCS CONCEPTS

• **Security and privacy** → **Domain-specific security and privacy architectures.**

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

ACSAC '23, December 4–8, 2023, Austin, Texas
© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXXX.XXXXXXX>

KEYWORDS

blockchain, federated learning, homomorphic encryption, security and privacy

ACM Reference Format:

Jorge Castillo, Phillip Rieger, Hossein Fereidooni, Qian Chen, and Ahmad-Reza Sadeghi. 2023. FLEDGE: Ledger-based Federated Learning Resilient to Inference and Backdoor Attacks. In *Proceedings of Annual Computer Security Applications Conference (ACSAC '23)*. ACM, New York, NY, USA, 17 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

In recent times Machine Learning (ML) has gained high popularity and it is used for an increasing number of applications. However, the need to collect a large amount of data to train ML models raises security and privacy concerns in applications where sensitive data (i.e., text messages typed on mobile phones, personal medical information) is constantly stored and manipulated. Federated Learning (FL) allows multiple parties holding private data to collaboratively train ML models. Rather than collecting all data on a central server, each party (client) trains a model locally (local model) and only shares its parameters with the coordinating server that aggregates the parameters from all individual clients. Afterward, this aggregation server distributes the aggregated model (global model) back to the clients for further training rounds.

This privacy feature, in combination with the performance gained by outsourcing the training process from one server to multiple clients, made FL the ideal training framework for different real-world applications, e.g., word suggestions in the mobile keyboard GBoard [37], brain tumor segmentation [57], risk detection on mobile devices [16], or identification of malware infected devices [47].

However, recent work challenges the security and privacy of FL, raising concerns about its practical applicability. For example, it was recently demonstrated that, given the model's parameters or predictions, model inference attacks could extract information about the training data from a model. Therefore by exploiting the model's memorization capabilities (or over-fitting), it is possible to reconstruct samples from the training data or determine if a certain

¹Affiliated with The University of Texas at San Antonio during the research and preparation of this paper.

²Affiliated with Technical University of Darmstadt during the research and preparation of this paper.

sample was used for training [33, 59], negating the privacy gains of FL. The most prominent example of this attack is the Membership Inference Attack [59]. FL is particularly vulnerable to inference attacks executed by a curious server, as it has access to the local models before aggregation. While the aggregation anonymizes the individual clients' contributions and makes inference attacks significantly harder [17], the server's access to individual local models poses a significant threat to the clients' privacy and raises concerns for applications with privacy-sensitive data [57]. To mitigate inference attacks, state-of-the-art defenses use one of the following approaches: Secure Multi-party Computation (SMPC) [2, 17, 54], Differential Privacy (DP) [38, 42], or Multi-Key Homomorphic Encryption (MKHE) [55].

In terms of security, recent work [5, 58] has demonstrated the high vulnerability of FL systems against Poisoning Attacks. In this type of attack, adversaries leverage the distributed property of FL to take control of one or more training clients. Malicious clients can alter their behavior and skew the convergence of the global model. Poisoning attacks can be divided into two categories: untargeted and targeted attacks. Untargeted attacks aim to degrade the utility of the global model [25]. In targeted (or backdoor) attacks [5, 66], however, an adversary guides the global model to a well-defined outcome. One example of an effective backdoor attack is where a malicious model is carefully trained to maintain a high accuracy for the main task but triggers backdoor behavior if specific patterns are detected during inference. Examples for such backdoor behavior include, e.g., injecting advertisement into an FL-based word suggestion system or creating a backdoor that makes an FL-based network intrusion detection system fail to detect network traffic of certain malware. A major threat of these attacks is that models are still black boxes. Thus, in practice, it is still an unsolved problem to determine if a model contains a hidden backdoor. The unique behavior makes backdoor attacks more important compared to other poisoning attacks. State-of-the-art defenses aim to minimize the threat of poisoning attacks using techniques such as model filtering [6, 58] to detect and exclude poisoned models from aggregation, and/or model clipping [48, 53] to limit poisoned updates' impact. These approaches, however, are limited by the underlying assumptions imposed by SMPC, e.g., availability during computations, and more importantly, they fail to address the motivation behind client misbehavior, i.e., accountability. Without crediting the contributions of individual clients, a malicious client may continuously try to poison the model until it is successful in one round. Further, mitigating poisoning and inference attacks at the same time is a complex task and existing approaches [27, 48] are not efficient.

To overcome the limitations of existing solutions, we tackle the following questions: i) how to achieve a privacy-preserving aggregation framework that penalizes malicious intent during the aggregation process, ii) how to discriminate poisoned from benign updates to dynamically reward or punish client's behavioral patterns.

Goals and Contributions. In this paper, we present the design, implementation, and evaluation of FLEDGE, a fully-decentralized crypto-system that provides resiliency to inference and poisoning attacks. FLEDGE is a 3-layer blockchain FL framework powered by smart contracts, where each layer operates specific components, i.e., training clients (client layer), smart contracts (computation layer), and ledger (data layer). Our primary motivation to use smart

contracts is to provide a decentralized and immutable environment to protect the security and privacy of models. By using smart contracts, we achieve reasonable efficiency and are able to mitigate poisoning and inference attacks. FLEDGE leverages blockchain's decentralization to yield high computation availability and ledger immutability (i.e., committed data cannot be changed) to prevent data alterations that could lead to unexpected results such as inaccurate model filtering or incorrect distributions of rewards.

To address the first question and perform privacy-preserving aggregation, we have to consider the following factors: private computation framework and aggregator compensation. To design a privacy-preserving computation platform, FLEDGE introduces the concept of Blockchain Two Contract Computation (BT2C). BT2C is defined as a semi-honest relationship³ between two smart contracts using Homomorphic Encryption (HE), in particular, we rely on the CKKS⁴ encryption scheme [11].

Compared to SMPC and MKHE, BT2C is a decentralized crypto-system based on HE that leverages the blockchain ledger to improve trust among smart contracts, where one contract (Defender) acts as a decryption service and the other contract (Gateway) acts as a computation hub. For our implementation, we develop a *secure decryption* method that includes a compensation algorithm to evaluate a reward for the aggregation service based on its behavior. Our approach operates as a semi-honest cryptographic service such that the Gateway contract receives encrypted models from training clients and performs computations; and the Defender contract evaluates model characteristics (e.g., cosine distance) and provides incentives (i.e., crypto-currency, tokens) for benign behavior.

To address the second question and discriminate poisoned models, we first separate it into the following components: poisoning detection and client compensation. To implement the poisoning detection, FLEDGE calculates the cosine distances between local and global models, and utilizes the Gaussian Kernel Density Estimation (G-KDE) function to divide them into different clusters. Here, a cluster is identified by the location of local minimums⁵. This information is leveraged as a breaking point to separate the distances into different clusters. After benign and malicious clusters have been correctly identified, FLEDGE implements a round-based client compensation algorithm to provide additional incentives to benign training behavior, and to penalize those who attempt model poisoning.

In summary, FLEDGE's contributions are threefold:

(1) FLEDGE offers strong privacy guarantees by operating models in cipher text using the proposed BT2C protocol. Our approach is shown to be resilient against white-box inference attacks with a probability of success of $\frac{1}{m!}$, where m represents the number of ciphers generated per model (Sect. 6.1).

(2) FLEDGE mitigates poisoning attacks using the proposed G-KDE clustering method to analyze the distribution of cosine distances and remove poisoned models. Our extensive evaluation on four public datasets (i.e., Reddit, MNIST, Fashion-MNIST, and CIFAR-10) indicates that FLEDGE is resilient against untargeted and targeted poisoning attacks (Sect. 6.2).

³The semi-honest setting is a well-established security model that dictates how involved parties must adhere to the pre-established protocol

⁴Note that we use the term HE to refer to CKKS in the rest of this paper.

⁵A local minimum is a point on the associated function (e.g., G-KDE) whose value is less than every other point in its vicinity.

(3) FLEDGE relies on our proposed aggregation and training compensation algorithms to offer incentives to benign aggregation services and benign training clients. Our results indicate that the proposed compensation algorithms automatically adjust the rewards to deter malicious intent from the training process (Sect. 6.3).

2 REQUIREMENTS AND CHALLENGES

This section presents the security and privacy requirements that FLEDGE fulfills and the challenges to be tackled in achieving them.

2.1 Privacy for FL

During model submission, clients upload trained models to the aggregation server such that the server generates a new global model (see App. A for details on the FL process). At this point, the server has complete access to each model (e.g., model weights, structure and hyperparameters), which increases the threat of white-box inference attacks. To mitigate the attack, the defender has to satisfy the following requirements:

P1: Utility Retention. The defense must provide resiliency against inference attacks that are executed by curious servers while maintaining the utility of the model, i.e., main task accuracy (MA) remains the same with or without defense. Therefore, the performance of the new global model must not be compromised with the increase of privacy levels.

P2: Computation Availability. The defense must remain available to process and analyze encrypted models⁶. Therefore, every model computation shall not fail due to limited resource availability.

To the best of our knowledge, existing solutions for inference attacks that also preserve model utility rely on frail computation infrastructures, e.g., SMPC [17, 48] or MKHE [55]. Thus, they suffer from the requirement of high availability of the system’s components to use privacy-preserving computations (SMPC and MKHE) [55] and also from a high computation complexity to detect poisoned models when using privacy-preserving computations. Our scheme combines blockchain (see App. B) and HE, in particular, the scheme of Cheon-Kim-Kim-Song (CKKS) [11] (see App. C), to introduce a unique privacy-preserving computation framework, overcoming current limitations. However, the use of blockchain brings additional concerns. Thus, FLEDGE addresses the following challenges:

C1: How to leverage blockchain to improve trust between computation parties.

C2: How to effectively combine HE and blockchain to limit the ledger’s transparency effect and increase privacy to model updates.

2.2 Security for FL

FL is a distributed learning approach that allows numerous clients to participate in the training process through model submissions. An adversary who controls a fraction of the clients can then use their influence to poison the new global model. To mitigate poisoning attacks, the defender has to fulfill the additional requirements:

S1: Effective Poisoning Mitigation. The defense must detect poisoning attempts, e.g., untargeted and targeted attacks, minimize their impact on the global model, and preserve model utility. For example, for targeted (backdoor) attacks, a defense should maintain

⁶An encrypted model is a collection of ciphers that represent encrypted weights.

the backdoor accuracy (BA) at the same level as without the attack. In addition, similar to **P1**, the defense must not negatively affect the training process, e.g., decrease MA by removal of benign models.

S2: Autonomous Behavior. The defense must be flexible to adjust automatically to different strategies without manual configuration.

Like existing solutions, FLEDGE leverages the cosine distance between local models and the global model to cluster their scores dynamically. Our approach, however, leverages this information to apply a deterrent to malicious clients, which adds another layer of security. Therefore, FLEDGE addresses the additional challenges:

C3: How to solve the dilemma of preventing the server from analyzing the local models against inference attacks while the server has to inspect the local models to detect/mitigate poisoned models.

C4: How to discriminate poisoned models s.t. malicious clients can be correctly identified to receive disciplinary actions.

C5: How to credit the clients over multiple training rounds to make malicious clients accountable for their attacks.

3 ADVERSARY MODEL AND ASSUMPTIONS

In this section, we describe the threat model and assumptions used for the rest of the paper. We highlight the adversary’s capabilities and main objectives.

3.1 Privacy Threat

Classic FL implementations rely on an aggregation server to compute new global models every training round (see App. A). However, a malicious aggregation server can extract private information from each of the local models, thus, raising privacy concerns.

White-box Inference Attack Goal. Aligned with previous research [17, 48], the *honest-but-curious* aggregator instantiates the attack on local models W_i before aggregation. In other words, an adversary \mathcal{A}^P is aware of any process happening in the aggregator, but remains *honest*, i.e., continues to perform the aggregator’s benign tasks, to avoid detection. However, \mathcal{A}^P is also *curious*, having the ability to infer private information about the training data D_i while processing W_i . Formally, \mathcal{A}^P leverages W_i to learn if a given input x was used as part of D_i , allowing \mathcal{A}^P to extract sensitive information from every local model. Aligned with previous work [3, 17, 27], we focus on inference attacks on the local models, as the aggregation anonymizes the individual contributions.

\mathcal{A}^P Capabilities. We assume \mathcal{A}^P is in full control of the aggregation server s.t. \mathcal{A}^P has access to every local model submitted by clients. We also assume \mathcal{A}^P cannot compromise clients directly or affect any of the training processes.

3.2 Security Threat

Multiple clients are selected to improve model accuracy. This collaboration allows one or more clients to conduct malicious activities in any training round.

Targeted Poisoning Attack Goals. In a targeted poisoning attack, the adversary \mathcal{A}^S has the following goals: poisoning injection and defense evasion. For poisoning injection, \mathcal{A}^S manipulates local model W_i to produce poisoned local model W'_i . W'_i is then used to alter the behavior of global model G_t . In state-of-the-art targeted poisoning attacks (also called backdoor attacks), \mathcal{A}^S guides poisoned global model G'_t to behave normally all the time except when

a specific set of conditions or triggers are present in the input. To achieve its secondary goal, \mathcal{A}^s manipulates W_i' s.t. it remains as close as possible to W_i , e.g., adapting the loss function [5].

\mathcal{A}^s Capabilities. Similar to recent studies [3, 5, 30, 52], we assume \mathcal{A}^s maliciously controls f compromised clients, which should be less than half of the total number of clients n ($f < \frac{n}{2}$). We also assume \mathcal{A}^s cannot observe benign clients' local data or their submitted local updates. To introduce a backdoor into the global model, \mathcal{A}^s can launch a combination between data poisoning [58] and model manipulation attacks[5]. Data poisoning is when \mathcal{A}^s adds *poisoned* data to the existing training sets during model training, e.g., for an image classification task, \mathcal{A}^s can poison an image by drawing a shape into a specific corner. This attack allows \mathcal{A}^s to change model predictions to its desired outputs every time a trigger is identified during inference. In contrast, the model manipulation attack lets \mathcal{A}^s control the training algorithm to alter the convergence point of a model. This attack can be implemented through model scaling, modifying the loss function and/or adapting dedicated hyperparameters.

3.3 Assumptions

FLEDGE provides numerous security and privacy benefits under the following assumptions.

A1: Consensus Protocol is not Compromised. Since blockchain is the underlying platform to exchange information and execute smart contracts, we assume the consensus process to be not compromised.

A2: Non-colluding Servers (Smart Contracts). During the manipulation of encrypted data, deployed smart contracts engage in a semi-honest relationship to enable a privacy-preserving aggregation infrastructure. Therefore, to preserve privacy guarantees, we assume an adversary cannot control both contracts and their storage components, simultaneously.

A3: Clients Perform Encryption. Training clients affiliated to FLEDGE are assumed to have sufficient computational resources to perform encryption.

4 DESIGN

This section first provides a high-level overview of FLEDGE and then describes its components in detail.

4.1 High-level Overview

FLEDGE is designed to fulfill privacy requirements (**P1**, **P2**) and security requirements (**S1**, **S2**). This is achieved by a layered framework that we will detail below. To detect poisoned models and satisfy **S1** and **S2**, FLEDGE uses a Gaussian Kernel Density Estimation (G-KDE) function to partition the received model updates into distinct clusters based on their pairwise distance. Since the cosine determines the angle, it reveals the changes that were applied to the local model. Compared to other metrics such as the Euclidean distance, it is more stable against manipulations. To satisfy **P1** and **S1** FLEDGE uses HE to encrypt models and perform privacy-preserving computations, i.e., private aggregation and/or private distance between models. In addition, FLEDGE leverages blockchain, in particular, smart contracts to meet **P2** and **S2**. Informally, FLEDGE is a 3-layer blockchain framework regulated

by smart contracts that provides FL services to train models on arbitrary learning tasks, e.g., image classification, and/or word prediction. For every new learning task, a session reward is set by the owners of the task s.t. interested parties (clients and/or contracts) who join can be rewarded for their benign efforts. In other words, FLEDGE operates a crediting system to encourage participants to avoid malicious attempts to break the system, e.g., white-box inference or poisoning attacks. To manage every reward, FLEDGE registers training clients by generating unique cryptographic identities via its Membership Service Provider (MSP).⁷ Fig. 1 presents the different layers of FLEDGE.

Client Layer. This is the base layer of FLEDGE, and it is where training clients reside. As discussed in Sect. 3.2 some of the clients can be controlled by the poisoning attacker \mathcal{A}^s .

Computation Layer. This layer illustrates the logical components that enable FLEDGE to operate autonomously. It is formed by two smart contracts, namely Gateway and Defender contracts. The Gateway contract acts as the access gateway where clients submit their local models. Its core functions, e.g., model process, model analysis and model aggregate, provide privacy-preserving computations for encrypted models, addressing C1. The Defender contract, on the other hand, provides support to the Gateway in the form of security and privacy mechanisms, e.g., model privacy and model security, thus, allowing FLEDGE to defend against multiple threats. We define this blockchain infrastructure as Blockchain Two Contract Computation (BT2C). The goal of our BT2C implementation is to enable HE-based computations and secure decryption functions tailored to provide privacy-preserving FL (as defined in Sect. 3.1). Further details about the internal methods of both smart contracts can be found in Steps 2–5 of Sect. 4.2.

Data Layer. The following layer represents the storage components of FLEDGE which constitutes two storage oracles, namely A and B, and a blockchain ledger. The storage oracles (i.e., external databases) are used to manage encrypted models and decryption keys for Gateway and Defender contracts, respectively, and the blockchain ledger stores information about FLEDGE (e.g., model information, session information) using the following transactions types (TT1 – TT7).

Init Transactions (TT1) are generated for each learning task to determine the owner of the task, the encryption keys to be used (i.e., HE public key P_k), the number of rounds T required, and the reward amount R for the full training session.

Storage Transactions (TT2) are generated for every encrypted local model W_i^* submitted by K clients to save the client ID (e.g., wallet address to receive payments), model ID, and the encrypted offset value δ_i^* , where $i \in [1, K]$. The offset δ is a random value generated based on the standard deviation of local model W_i that is injected before model encryption to further obfuscate W_i^* . Further information about the use of δ is found in Step 1 of Sect. 4.2.

Analysis Transactions (TT3) are created for every TT2 to compute the cosine distance between W_i^* and the current encrypted global model G_t^* , where $t \in [1, T]$. TT3 is used to store model ID and its respective score c_i , i.e., cosine distance.

Privacy Transactions (TT4) are generated when a malicious contract

⁷In Fabric, an MSP provides verifiable identities to members of the blockchain network.

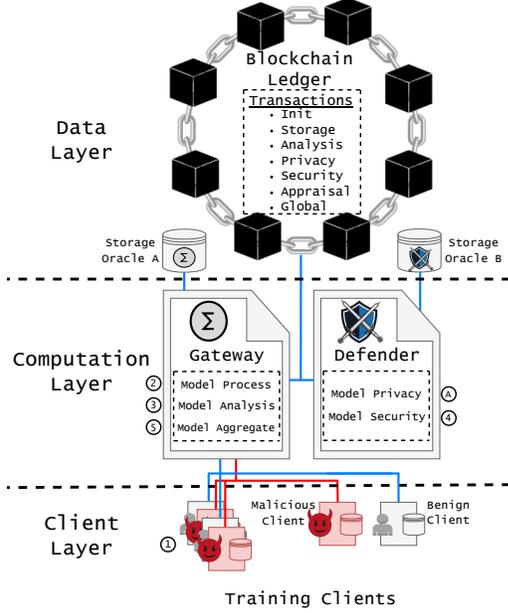


Figure 1: FLEDGE System Overview. Annotated steps illustrate the operation of FLEDGE during a training round t .

(i.e., Gateway) is attempting to break the privacy of W_i^* . TT4 includes the computed aggregation reward R_C for a given training session. Further information for R_C is found in Step A of Sect. 4.2. *Security Transactions (TT5)* are created after FLEDGE evaluates every c_i from TT3 to classify clients into two categories: benign or malicious. This is represented in TT5 as a list of benign IDs and a list of malicious IDs.

Appraisal Transactions (TT6) are determined after TT5 to calculate the round reward R_r for benign clients. Additional details for R_r can be found in Step 4 of Sect. 4.2.

Global Transactions (TT7) are determined after model aggregation has occurred. TT7 includes the global ID, the decrypted weights of the new global model G_{t+1} , and the corresponding encrypted global weights G_{t+1}^* .

4.2 FLEDGE Details

To initialize FLEDGE, the smart contracts seen in the computation layer are deployed and initialized to the blockchain network. The initialization process for the smart contracts is completed after generating TT1. To avoid the possibility of data modification (or forks) at run-time, we refer back to assumption **A1**. Similarly, each smart contract takes into consideration assumption **A1** to protect the integrity of the contracts before/after deployment into the blockchain, and assumption **A2** to prevent an adversary from gaining full control over the system. Finally, to achieve privacy-preserving computations (as defined in Sect. 3.1), clients are bound to assumption **A3** to protect the privacy of local updates.

The annotated steps seen in Fig. 1 illustrate the learning process of FLEDGE during a training round t . Here, we separate the learning process into 6 steps, i.e., 5 main steps (Step 1 – Step 5) and

1 intermediary step (Step A). After multiple clients have joined a learning task, they first download the previous global model G_{t-1} and the corresponding encryption key P_k from TT1 at the Gateway contract. Note that for every P_k , a corresponding secret key S_k is generated and maintained by the Defender contract. Furthermore, every P_k includes an encryption context, which is provided to the clients. The encryption context contains the degree of the polynomial (PolyDeg) used to generate P_k and S_k . This value determines the size of P_k and the size of produced ciphers in terms of bytes. Clients may continue to use the same P_k unless a new public key is required by the system.

Model Encryption (Step 1). Each client i starts to train the model using local data D_i for a predefined number of epochs. After training, clients generate and inject an offset constant δ_i such that $W_i' = W_i + \delta_i$. More specifically, δ_i is generated from the multiplication of two random elements: the model's standard deviation σ_{W_i} after training, and a scaling factor $f_s \in [-100, 100]$ s.t. $f_s \neq 0$. Note that f_s is bounded to $[-100, 100]$ to avoid exceedingly large numbers (positive or negative) as model weights. This is primarily because we are interested in shifting (left or right) the distribution of W_i using the inherently random properties of each local model. Contrary to DP, δ_i is recorded to be used in Step 3. The offset is applied to mask W_i and obfuscate the model during private computations. At this stage, an attacker would require to brute force every δ_i in order to break the privacy of a single local update. Then clients start the encryption process of W_i' and δ_i . Once a client has generated its encrypted local model W_i^* and encrypted offset δ_i^* , the client proceeds to submit them to the Gateway contract for further analysis. By this, FLEDGE addresses **C2**.

Due to the limitations of HE, a client is required to first separate W_i' into multiple chunks of data s.t. $W_i' = w'_1 \dots w'_n$, where n is the number of chunks per model. To calculate n , we first determine the capacity for every cipher z , i.e., the maximum amount of elements each cipher is able to contain. We define the capacity to be PolyDeg / 2, e.g., a PolyDeg of 2048 yields a capacity of 1024 elements per cipher. Thus, the number of ciphers required to encrypt W_i' is directly proportional to the number of trainable parameters given PolyDeg. This is further illustrated by Eq. 1:

$$z_1 \dots z_n = \text{Encrypt}(w'_1 \dots w'_n, P_k), n = \frac{\text{len}(W_i')}{\text{PolyDeg}/2} + 1 \quad (1)$$

Algorithm 1: BT2C – Private Cosine Distance

Input: δ^* \leftarrow encrypted offset

G^* \leftarrow encrypted global model

W^* \leftarrow encrypted local model

- 1 $Z_D \leftarrow \text{PrivateDotProduct}(G^* + \delta^*, W^*)$
 - 2 $X_D \leftarrow \text{SecureDecryption}(Z_D) \leftarrow$ defender function
 - 3 $Z_G \leftarrow \text{PrivateMagnitudeSquared}(G^* + \delta^*)$
 - 4 $X_G \leftarrow \text{SecureDecryption}(Z_G)$
 - 5 $Z_L \leftarrow \text{PrivateMagnitudeSquared}(W^*)$
 - 6 $X_L \leftarrow \text{SecureDecryption}(Z_L)$
 - 7 $c \leftarrow 1 - \frac{\sum_{i=1}^n X_{D_i}}{\sqrt{\sum_{i=1}^n X_{G_i}^*} \sqrt{\sum_{i=1}^n X_{L_i}}}$
 - 8 $\text{UpdateScoreToLedger}(c) \leftarrow$ new TT3
-

Algorithm 2: BT2C – Secure Decryption

```

Input :  $z_1, \dots, z_m$   $\leftarrow$  computation ciphers
Output:  $X$   $\leftarrow$  array of decrypted numbers
          $\rho$   $\leftarrow$  array of decrypted model chunks
1  $\delta_1^*, \dots, \delta_K^* \leftarrow$  ReadOffsetFromLedger()  $\leftarrow$  from TT2
2  $S_k \leftarrow$  ReadKeyFromStorage()
3  $t \leftarrow 0.05$   $\leftarrow$  array variation tolerance
4 for each cipher  $i$  in  $[1, m]$  do
5    $\rho_i \leftarrow$  Decrypt( $z_i, S_k$ )
6    $v \leftarrow \left| \frac{\max(\rho_i) - \min(\rho_i)}{\max(\rho_i)} \right|$   $\leftarrow$  compute variation
7   if  $v \leq t$  then
8      $X_i \leftarrow$  Average( $\rho_i$ )
9   else if  $K > 1$  then
10    for each offset  $j$  in  $[1, K]$  do
11       $\delta_j \leftarrow$  Decrypt( $\delta_j^*, S_k$ )
12    end
13     $\rho_i \leftarrow \frac{\rho_i - \sum_{j=1}^K \delta_j}{K}$   $\leftarrow$  offset removal/injection
14  else
15     $R \leftarrow$  ReadRewardFromLedger()  $\leftarrow$  from TT1
16     $s \leftarrow$  CountSessionsFromLedger()  $\leftarrow$  # TT1
17     $\phi \leftarrow$  CountAnomaliesFromLedger()  $\leftarrow$  # TT4
18     $R_C \leftarrow 0.1 * R * e^{-(\phi+1)/s}$   $\leftarrow$  calculating reward
19    UpdateContractRewardToLedger( $R_C$ )  $\leftarrow$  new TT4
20     $\rho_i \leftarrow \emptyset$   $\leftarrow$  empty set
21  end
22 end
23 return  $X$  or  $p$   $\leftarrow$  output type dependent on process

```

For FLEDGE, we have determined a minimum PolyDeg of 4096 is required to successfully compute the desired private functions.

Model Process (Step 2) is the initial function that receives every W_i^* provided by the clients. In this step, the Gateway contract stores W_i^* into storage oracle A to avoid public visibility to any other contract deployed in the network. The storing process saves the ciphers as encoded text into a single document. Every pair of W_i^* and δ_i^* is used to generate and submit a new TT2 to the ledger.

Model Analysis (Step 3) uses the previously submitted TT2 to retrieve the encrypted model from storage and its corresponding encrypted offset. δ_i^* is used to offset the encrypted global model G_{t-1}^* , where G_{t-1}^* can be easily downloaded from TT7 of the previous round. This process aligns encrypted models to compute an accurate cosine distance as given by Alg. 1.

Formally, the private cosine distance function seen in Alg. 1 requires as inputs the encrypted global model G^* , the encrypted local model W^* and the corresponding encrypted offset δ^* . Its goal is to compute the cosine distance score c between G^* and W^* . To calculate the distance, the computation process is segmented into three BT2C rounds. This is to overcome the practical limitations of HE, e.g., inability to compute roots. The first round (lines 1–2) starts by computing the encrypted dot product Z_D between $G^* + \delta^*$ and W^* , where Z_D is a collection of ciphers z_1, \dots, z_n that represent the encrypted value of the dot product operation. Z_D is then delivered to Defender contract to perform *secure decryption*. Note that Z_D

might be in any order to add randomness to the decryption process. This process returns X_D , a collection of numbers x_1, \dots, x_n that represents the results of $\sum G \cdot W$. Similarly, the two remaining rounds (lines 3–4 and 5–6) are used to generate X_G ($\sum G^2$) and X_L ($\sum L^2$), respectively. To finalize the computation process, the values for all three rounds are combined to calculate c and submitted to the ledger (TT3), as seen in line 7–8.

Model Privacy (Step A) relies on TT1, TT2 and previous TT4 to enable the *secure decryption* function seen in Alg. 2. This function includes two unique operations: limitation of data decryption (lines 1–13) and reward adjustment (lines 15–20). The former analyzes the information in every cipher to either return a collection of numbers X (see Step 3) or a collection of model chunks ρ (see Step 5). The latter regulates the contract reward R_C using the information stored in the ledger. The use of R_C enables FLEDGE to incentivize benign aggregation behavior in the framework while penalizing malicious conduct such as attempting to access local models. Note that in FLEDGE, R_C is set to be 10% (max) of the session reward R by default.

Alg. 2 requires as inputs only the computation ciphers z_1, \dots, z_m , where m is the number of submitted ciphers in the BT2C round. Note that for secure decryption, m is independent from the number of ciphers n in an encrypted model such that $m \leq n$.⁸ Our approach returns decrypted data, which is represented by an array of numbers X or an array of model chunks ρ . To initiate the decryption process and address C3, the Defender contract retrieves (lines 1–2) every encrypted offset $\delta_1^*, \dots, \delta_K^*$ from TT2, and the corresponding secret key S_k from storage oracle B. The variation tolerance value in line 3 is set to $t = 0.05$ (or 5%) as it is required to discriminate summation operations, e.g., $\sum G \cdot W$, from model operations, e.g., model aggregation. More specifically, summation operations are determined by a low array variation, which indicates that all elements are the same or closely related⁹; and model operations are characterized with high variations as these are represented by distinct values within the decrypted array. After decryption (line 5), the data is analyzed w.r.t. the array variation factor v (line 6), where v is defined as the absolute percent difference between the maximum and minimum elements within ρ_i . At this step, if $v \leq t$, the elements inside ρ_i are considered to be the result of a summation operation, thus, generating X_i to represent its average as shown in line 8. Otherwise, the function proceeds to consider ρ_i as a model operation, where ρ_i is then adjusted by every δ_i and divided by the number of available models K (from TT2) to complete the aggregation process $\sum_{i=1}^K \frac{W_i}{K}$ as defined by lines 9–13.

Lines 15–20, however, are used to assess K since aggregation must be performed with a minimum of $K = 2$ models. If $K \leq 1$, it adjusts R_C to penalize the contract for attempting to bridge the confidentiality of the first local model as this would not be obfuscated by multiple offset constants. This subroutine is performed by gathering the current session reward R and the number of successful training sessions s from TT1 (lines 15–16), and the number of privacy anomalies ϕ previously registered from TT4 (line 17). Note that ϕ

⁸A practical example is when the aggregation contract (Gateway) divides its computations into multiple rounds s.t. $Z_m \in Z_n$ to prevent the Defender from potentially accessing full data.

⁹Array symmetry comes natively in HE to maintain optimum conditions during computations.

Algorithm 3: Poisoning Defense

Input: (c_1, \dots, c_K) \triangleleft distance scores

- 1 $f \leftarrow 2000$ \triangleleft resolution factor for smooth curves
- 2 $(x_1, \dots, x_f), (y_1, \dots, y_f) \leftarrow \text{GaussianKDE}([c_1, \dots, c_K], f)$ \triangleleft compute gaussian kernel density estimation
- 3 $(l_1, \dots, l_N) \leftarrow \text{LocalMinimums}([y_1, \dots, y_f])$ \triangleleft l_n is the index of local minimum found in y
- 4 $G \leftarrow \{[x_1, x_{l_1}], \dots, [x_{l_{N-1}}, x_{l_N}], [x_{l_N}, x_f]\}$ \triangleleft group set based on local minimums
- 5 $M \leftarrow N + 1$ \triangleleft maximum number of available groups
- 6 **for** each group m in $[1, M]$ **do**
- 7 **for** each score i in $[1, K]$ **do**
- 8 **if** $c_i \in G_m$ **then**
- 9 $g_m \leftarrow i$ \triangleleft append model index i to a group
- 10 **end**
- 11 **end**
- 12 **end**
- 13 $\text{UpdateGroupsToLedger}(g)$ \triangleleft new TT5. g_1 is closest to G_{t-1}
- 14 $R \leftarrow \text{ReadRewardFromLedger}()$ \triangleleft from TT1
- 15 $T \leftarrow \text{ReadTotalNumberOfRoundsFromLedger}()$ \triangleleft from TT1
- 16 $R_C \leftarrow \text{ReadContractRewardFromLedger}()$ \triangleleft from TT4
- 17 $R_\tau \leftarrow \frac{R - R_C}{T * \text{len}(g_1)}$ \triangleleft training reward
- 18 $\text{UpdateTrainingRewardToLedger}(R_\tau)$ \triangleleft new TT6

increases every time $K \leq 1$. To finalize the penalization process, the new R_C is generated and added to the ledger (TT4) as seen in lines 18–19, and ρ is set to empty to avoid leaking information (line 20). **Model Security (Step 4)** collects the cosine distance scores c_i from TT3, applies the proposed clustering technique to filter poisoned models, and determines the client reward R_τ to promote benign training behavior. To remove poisoned updates and address **C4**, our clustering method uses the Gaussian Kernel Density Estimation (G-KDE) function to identify the number of data distributions within the distance scores c_i , and selects models according to their assigned distribution. If models are determined to be malicious, they are removed from the aggregation process and penalized by receiving a reward of 0 for that round. These leftover rewards are divided evenly among other clients, thus, addressing **C5**. Additional information related to the use of G-KDE to generate clusters is discussed in App. D.

The poisoning defense is presented in Alg. 3. The defense requires distance scores c_1, \dots, c_K as input to generate model clusters (or groups). To produce an accurate representation of the distribution between scores, we set a resolution factor f of 2000 (line 1), where f denotes the number of data points used to fit the G-KDE function. Note that we empirically found that $f = 2000$ provides the necessary resolution to find local minimums. In line 2, c_1, \dots, c_K and f are used to compute G-KDE, thus, generating 2000 (x, y) data points, where x is bounded between $\min(c_1, \dots, c_K)$ and $\max(c_1, \dots, c_K)$, and y is the density estimation obtained from the process. Density values y_1, \dots, y_f are used in line 3 to calculate the location or index of every local minimum l_1, \dots, l_N found in the distribution of scores. These locations are used to generate a group set G that contains $N + 1$ (M) groups of models (line 4). At this stage,

Algorithm 4: BT2C – Private Aggregation

Input: W_1^*, \dots, W_N^* \triangleleft selected models

- 1 $Z \leftarrow W_1^*$ \triangleleft encrypted base model
- 2 **for** each update i in $[2, N]$ **do**
- 3 $Z \leftarrow \text{Add}(Z, W_i^*)$
- 4 **end**
- 5 $G_t \leftarrow \text{SecureDecryption}(Z)$ \triangleleft defender function
- 6 $P_k \leftarrow \text{ReadKeyFromLedger}()$ \triangleleft from TT1
- 7 $G_t^* \leftarrow \text{Encrypt}(G_t, P_k)$
- 8 $\text{UpdateGlobalToLedger}(G_t^*, G_t)$ \triangleleft new TT7

each score is allocated inside specific groups g to separate benign models from malicious (lines 6–12). These groups are committed to the ledger as part of TT5 in line 13. Finally, to calculate R_τ , the process retrieves R and the number of training rounds T from TT1, the current R_C from TT4, and the group closest to G_{t-1} (g_1) are combined by the defense as seen in lines 14–17. The updated R_τ is committed to the ledger as a new TT6 (line 18).

Model Aggregate (Step 5) selects the models defined by TT5 from storage oracle A to compute a new global model G_t . The private aggregation (Alg. 4) uses a single BT2C computation round to create G_t .

Formally, Alg. 4 requires as inputs every selected local model W_1^*, \dots, W_N^* , where N is the number of admitted models selected in Step 4. Models are simply added (lines 1-4) into a single encrypted model Z , where $Z = (z_1, \dots, z_n)$. In line 5, Z is submitted to Defender contract to complete the aggregation process, which returns a collection of arrays denoted as G_t . The new model is then encrypted (lines 6–7) using the available P_k to produce the new encrypted global model G_t^* . G_t and G_t^* are compiled into a new TT7 and committed into the ledger (line 8) to prepare FLEDGE for the next training round.

5 EXPERIMENTAL SETUP

The following sections illustrate our testbed, and describe the datasets and models used during evaluation. Note that a detailed list of evaluation metrics is provided in App. E.

Experimental Testbed. We simulate a generic blockchain using Hyperledger Fabric (HLF) to illustrate the practicality of our approach for other blockchain implementations. For experimental setup, we abstract away the complexities introduced by the consensus protocol, and instead focus on the computational entanglements added by FLEDGE. This is because FLEDGE relies solely on smart contracts rather than the underlying blockchain platform. To instantiate the simulation environment, we deploy docker containers on a Windows PC with Intel Core i7-9750H and 32 GB RAM. The blockchain test network is formed by a single ordering node operated by single organization with two peers transacting under a single communication channel.

To fit multiple encrypted models within a single block, we increase the block size to 100MB. Note that this is 100 times larger than a Bitcoin block (1MB) [45]. We implement the gateway and defender contracts (chaincodes) using NodeJS and the Node-SEAL library. Node-SEAL is a NodeJS wrapper library used to interface

Table 1: Dataset description for different learning tasks.

| Application | IC | | | WP |
|-----------------|-------|---------|----------------------------|---------|
| Datasets | MNIST | Fashion | CIFAR-10 | Reddit |
| #Records | 70K | 70K | 60K | 20.6M |
| Model | CNN | CNN | ConvMixer _{256/3} | LSTM |
| #params | ~ 23K | ~ 29K | ~ 234K | ~ 20M |
| #ciphers | 12 | 15 | 115 | ~ 10.1K |

with Microsoft SEAL [10], an efficient and open-source HE library available in C++. Our HE setup uses a PolyDeg of 4096 to encrypt local models. To evaluate models, we use Pytorch in an Ubuntu 20 server with 2 AMD EPYC 7302, 480 GB RAM, and 6 NVIDIA A100 (40 GB RAM each).

Datasets and Models. To assess FLEDGE, we use two popular FL applications: word prediction (WP) [36, 38], and image classification (IC) [12, 31, 57]. Note that every model used for evaluation has been pre-trained to reach an acceptable accuracy level. Tab. 1 describes the dataset types (Dataset), the rounded number of records per dataset (#Records), the AI models used for training (Model), the number of trainable parameters (#params) and the number of ciphers (#ciphers) found per model. We use smaller models with fewer trainable parameters for IC datasets, compared with WP, to evaluate how model complexity impacts FLEDGE.

Word Prediction (WP). We use the Reddit dataset as example of WP for Natural Language Processing (NLP) applications, e.g., the real-world FL application G-Board [37]. The dataset contains over 20M records of reddit users’ posts from November 2017. Following previous work [5, 53] we use a 2-layer LSTM model.

Image Classification (IC). We selected three popular IC datasets of different image complexity: MNIST, Fashion-MNIST (or Fashion for short) and CIFAR-10. They all consist of 10 evenly divided categories, where MNIST contains 70K handwritten digits, Fashion has 70K images of articles of clothing (i.e., shoe, dress, shirt), and CIFAR-10 has 60K pictures of objects (i.e., frog, airplane, car). For MNIST and Fashion, we use a simple CNN model comprised of 1 and 2 CNN layers, respectively. We customized the ConvMixer model [62] to train CIFAR-10 with a width of 256.

Backdoor Attacks. Aligned with earlier work [3, 5, 48] we use the constrain-and-scale attack of Bagdsaryan *et al.* [5]. Note that we focus on adaptive attacks, e.g., adversary adapts the loss function using the same metric as the defensive strategy. In other words, our adversary model leverages the *cosine distance* in an attempt to evade our defense. For the Reddit dataset, the adversary aims to make the model predicting the word "delicious" after the trigger "pasta from astoria tastes" [5]. The CIFAR-10 backdoor shall make all cars in front of a striped background being classified as birds [5]. For MNIST and Fashion MNIST, the backdoor forces models to predict the number 0, and t-shirt/top, respectively, when the image trigger is detected. The trigger is simply a white rectangle located at the bottom left corner of each poisoned image.

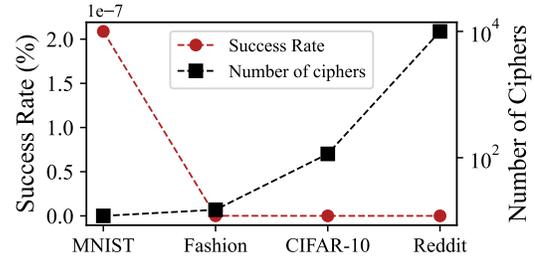


Figure 2: Probability of success for white-box inference attack w.r.t. model complexity.

6 EXPERIMENTAL RESULTS

The following sections evaluate the privacy of FLEDGE under a naive setup, the security aspect of FLEDGE against poisoning attacks, and the behavior of the reward system in FLEDGE for aggregation services and clients. We also provide a run-time performance analysis of FLEDGE in App. F, which is used to illustrate the increased complexity of our learning process.

6.1 White-box Inference Attack Resiliency

Evaluation Baseline. To evaluate the privacy of FLEDGE, we step outside **A2** to explore a limited collaboration between the Gateway and Defender contracts. In this scenario, the Gateway is in full control of the attacker and attempts aggregation when there is only one model in FLEDGE ($K = 1$), disregarding its potential reward R_C . The Defender, however, is only partially compromised allowing the attacker to observe ρ_i during secure decryption such that the attacker can reverse the offset from local model W_1 .

Effectiveness of FLEDGE. We evaluate the effectiveness of the obfuscation techniques implemented in FLEDGE to prevent white-box inference attacks. To breach the privacy of W_1 , a Defender requires to find the correct order of ciphers, since this is random for every BT2C computation round. We define such a brute force

Table 2: Effectiveness of FLEDGE against multiple poisoning attacks in terms of Backdoor Accuracy % (BA) and Main Task Accuracy % (MA).

| Poisoning Attack | Dataset | No Defense | | FLEDGE | |
|-------------------------|----------|------------|------|--------|------|
| | | BA | MA | BA | MA |
| Untargeted [23] | Reddit | - | 15.8 | - | 22.7 |
| | MNIST | - | 91.5 | - | 98.3 |
| | Fashion | - | 41.1 | - | 90.0 |
| | CIFAR-10 | - | 28.9 | - | 83.0 |
| Constrain-and-Scale [5] | Reddit | 100 | 22.6 | 0.0 | 22.7 |
| | MNIST | 98.0 | 87.7 | 0.4 | 98.3 |
| | Fashion | 100.0 | 69.3 | 2.4 | 90.6 |
| | CIFAR-10 | 100.0 | 66.1 | 0.0 | 83.8 |
| DBA [66] | Reddit | 100.0 | 22.6 | 0.0 | 22.7 |
| | MNIST | 82.6 | 77.2 | 0.1 | 98.3 |
| | Fashion | 99.7 | 36.7 | 1.0 | 98.3 |
| | CIFAR-10 | 85.2 | 67.4 | 2.1 | 83.8 |

Table 3: Comparison of FLEDGE and five state-of-art defenses’ efficiency to mitigate backdoor. BA refers to Backdoor Accuracy % and MA refers to Main Task Accuracy %.

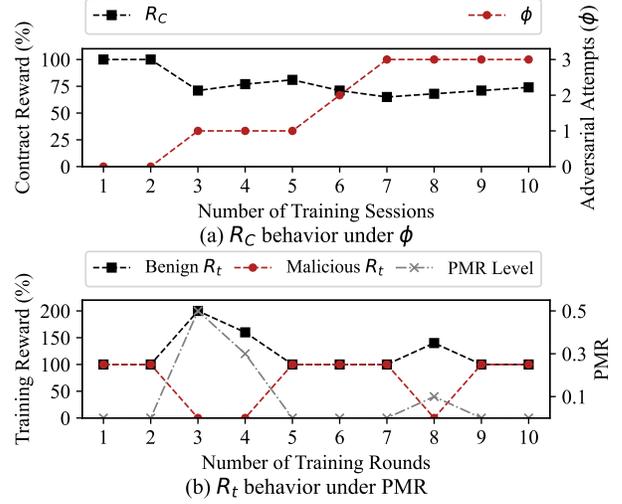
| Defenses | Reddit | | MNIST | | Fashion | | CIFAR-10 | |
|----------------|------------|-------------|------------|-------------|------------|-------------|------------|-------------|
| | BA | MA | BA | MA | BA | MA | BA | MA |
| Benign Setting | 0.0 | 22.7 | 0.5 | 98.3 | 3.7 | 90.9 | 0 | 83.9 |
| No Defense | 100.0 | 22.7 | 98.0 | 87.7 | 100.0 | 69.2 | 100.0 | 66.1 |
| Krum [6] | 100.0 | 22.6 | 0.6 | 98.3 | 2.8 | 90.1 | 0.0 | 83.0 |
| FoolsGold [18] | 0.0 | 22.7 | 0.5 | 98.3 | 3.0 | 90.7 | 0.0 | 83.6 |
| Auror [58] | 100.0 | 22.5 | 0.5 | 98.3 | 2.5 | 90.9 | 0.0 | 83.9 |
| AFA [43] | 100.0 | 22.6 | 83.1 | 94.2 | 97.9 | 87.3 | 100.0 | 66.5 |
| DP [38] | 77.0 | 22.0 | 26.5 | 97.3 | 52.2 | 88.6 | 60.0 | 76.6 |
| FLEDGE | 0.0 | 22.7 | 0.4 | 98.3 | 2.4 | 90.6 | 0.0 | 83.8 |

process to be equivalent to $m!$, where m is the number of ciphers for the BT2C round. Therefore, the probability of an attacker finding the right combination of ciphers to generate the correct model W_1 is given by $1/m!$. Fig. 2 illustrates the potential success rate ($1/m!$) of an attacker to extract W_1 , and its contrast w.r.t. m for every dataset. Here, we observe that the probability of success decreases from $2.08e-7\%$ to 0% as the number of ciphers increases from 12 (MNIST model) to $\sim 10.1K$ (Reddit model), i.e., large models provide better resiliency. Therefore, FLEDGE is $1/m!$ resilient to white-box inference attacks even during a limited collaboration outside **A2**.

6.2 Poisoning Mitigation

Evaluation Baseline. We set $PMR=0.5$, $non-IID=0.7$, $PDR=0.5$ and $\alpha=0.7$ as baseline parameters for untargeted and targeted attacks (unless otherwise indicated). PMR (or Poisoned Model Rate) indicates the influence level of an attacker to the system, i.e., PMR of 0.5 denotes an attacker maliciously controls 50% of training clients. $Non-IID$ data (or non-Independent and Identically Distributed) represents the number of training samples dispersed to a client that belongs to a specific class within a pre-defined group, i.e., $non-IID$ of 0.7 means that clients should use 70% of training data from their given class while the rest 30% is from the remaining classes. We follow the approach in [15] to prepare each dataset according to the number of output classes. PDR (or Poisoned Data Rate) determines the fraction of poisoned samples with respect to benign samples during training, i.e., PDR of 0.5 defines 50% of training data from a target class are poisoned samples. A higher PDR increases the success rate of the attacks. Similarly, the regularization term α , as defined by Bagdasaryan *et al.* [5], balances the loss function of client models aiming at limiting the distance between local and global models. A high value of α allows the attacker to increase its success rate at the cost of visibility.

Effectiveness of FLEDGE. We evaluate the resiliency of FLEDGE against different poisoning attacks such as untargeted poisoning [23], constrain-and-scale [5] and DBA [66]. The experimental results illustrated in Tab. 2 show that FLEDGE successfully mitigates these attacks without sacrificing benign performance (MA). For untargeted poisoning, the adversary successfully degrades model performance when no defenses are in place, reaching as low as 15.8% (22.7% original) for Reddit, and 28.9% (83.9% original) for CIFAR-10.

**Figure 3: Behavior of rewards (a) R_C , (b) R_t in FLEDGE.**

During constrain-and-scale attacks, the adversary is able to inject a backdoor into the model with almost 100% accuracy. Similarly, for DBA, the backdoor is injected into the global model with 80+%. These attacks, however, are not effective against FLEDGE as $BA \approx 0$ for every evaluated dataset.¹⁰ Moreover, FLEDGE maintains or even increases MA. Note that for the rest of the evaluation, we focus on targeted (or backdoor) attacks since they are the most sophisticated type of poisoning attacks.

Comparison to Existing Work. Tab. 3 compares the effectiveness of FLEDGE with five state-of-the-art defense approaches [6, 18, 36, 43, 58]. Notably, several defenses such as Krum [6] cannot handle non-IID scenarios. FoolsGold is the most resilient defense that mitigates backdoors for all four datasets as its BAs are very close to 0, but still a little higher than FLEDGE’s BA rates. Similar to FoolsGold, the other four defenses’ BA rates are much higher than or equal to FLEDGE’s while MA rates are much lower than or equal to ours. Auror [58] works well for the image datasets, where the clients’ local datasets overlap and show similar distributions but fails for the intrinsic non-IID Reddit dataset. Therefore, FLEDGE is shown to provide the most resilient defense to mitigate state-of-the-art backdoors. Appendix H and I provide further experiment results for WP and IC tasks respectively, showing FLEDGE’s performance for different PDRs, PMRs, further attack strategies, and IID rates.

6.3 Reward Analysis

The reward systems implemented in FLEDGE are an additional layer of security designed to discourage malicious actions during the learning process. However, a defensive strategy is only as good as its weakest component. In other words, FLEDGE’s reward mechanisms are constrained by the efficiency of its white-box inference resiliency (see Sect. 6.1) and its poisoning defense (see Sect. 6.2). Thus, the following section investigates the rewarding mechanism

¹⁰In some applications, misclassifications are counted in favor of the BA if $MA < 100\%$. For this reason, the BA is greater than 0%, e.g., 2.4% for Fashion, although the aggregated model does not contain the backdoor.

behind FLEDGE. Fig. 3 shows the behavior of the contract reward R_C and the training reward R_T .

Effect of ϕ in R_C . We assume FLEDGE has received the first local model W_1 ($K = 1$), and that the adversary has control over the Gateway contract. Fig. 3a illustrates how R_C is adjusted by the Defender every time an attempt to access private models (ϕ) is registered during *secure decryption*. In other words, ϕ represents the number of TT4 in FLEDGE, where a new TT4 is generated every time *secure decryption* is attempted for $K \leq 1$ as seen in Step A of Sect. 4.1. Note that we have simulated three attempts (Session 3, 5 & 7) to access W_1 , which forces the Defender to adjust R_C . In particular, we observe that R_C is severely affected by ϕ in comparison to the number of training sessions (s). This indicates that to increase R_C , the contract must behave normally for a large number of training sessions such that the ratio ϕ/s approximates 0.

Effect of PMR in R_T . For this experiment, we vary the PMR to $\{0.1, 0.3, 0.5\}$ to observe the behavior of the training reward R_T . Fig. 3b shows the benign R_T and malicious R_T under different PMR levels. Note that the process to determine R_T for each client is based on the number of malicious models found during a training round. In other words, the Defender forces malicious clients to transfer their potential rewards to benign ones, i.e., benign clients get 140% (Round 8), 160% (Round 4) and 200% (Round 3) for PMR of 0.1, 0.3 and 0.5, respectively. This process promotes benign training behavior since R_T is reduced to 0 for every client detected as malicious during a training round.

7 SECURITY AND PRIVACY ANALYSIS

The following sections provide a security and privacy analysis to further explore the resiliency against white-box inference and poisoning attacks under different adversary configurations. We also discuss the robustness of FLEDGE against clients randomly dropping from the learning process in App. G.

7.1 FLEDGE Privacy Analysis

FLEDGE uses a decentralized crypto-system maintained by Gateway and Defender contracts. This allows FLEDGE to manage and analyze ciphers. FLEDGE adopts a semi-honest security model such that only one (out of 2) entity could be compromised at the time as discussed in assumption A2. Therefore, to undermine the privacy of FLEDGE according to Sect. 3.1, an adversary \mathcal{A}^P may formulate one of the following scenarios.

\mathcal{A}^P Compromises Gateway Contract. If \mathcal{A}^P maliciously controls the Gateway contract, \mathcal{A}^P would have access to every encrypted model. However, \mathcal{A}^P cannot decrypt any model directly as the private key is only held by the Defender contract. To access local updates, \mathcal{A}^P may try the following strategies.

Direct Decryption Request. \mathcal{A}^P directly requests decryption of encrypted models by submitting the appropriate ciphers to Defender. This initial approach yields ineffective results as the decryption process follows *secure decryption* (Step A in Sect. 4.2), where this process can identify the type of computation (i.e., summations or model operations) by analyzing data composition after decryption. To mitigate this attack, our approach identifies each cipher as a model operation and returns decrypted arrays with injected random noises δ . Consequently, attempts to decrypt local updates result

in random shifts to the distribution of each model. Therefore, this defense can effectively obfuscate local updates when an adversary attempts to decrypt them directly.

Reverse Engineer from Computations. A sophisticated \mathcal{A}^P may consider to reverse engineer encrypted local updates from the result of specific computations, i.e., $G_i^* + \delta_i^* - W_i^*$, W_i^{*2} . However, this approach is also found ineffective since any type of model operation is constrained by *secure decryption*, resulting in data arrays being indirectly affected by δ .

FLEDGE Aware Decryption Request. \mathcal{A}^P attempts decryption of the first ($K = 1$) encrypted model committed to FLEDGE during a training round to bypass the security measures imposed by *secure decryption*. Put differently, \mathcal{A}^P aims to extract the first local model before other δ values skew the results. However, this approach is also ineffective because the Defender contract is aware of the number of models currently present in FLEDGE. As a consequence, the Defender contract leverages that information to adjust the contract reward R_C to penalize the attempt and address curious behavior as illustrated in Sect. 6.3.

\mathcal{A}^P Compromises Defender Contract. For the following scenario, \mathcal{A}^P attempts to visualize local weights during *secure decryption* as Defender holds the private key. To achieve this goal, \mathcal{A}^P requires the assistance of Gateway contract as the latter is the one that holds every encrypted model. Put differently, \mathcal{A}^P needs the Gateway to send encrypted models rather than encrypted computations, which breaks assumption A2. Therefore, FLEDGE is resistant to \mathcal{A}^P given assumption A2 holds.

\mathcal{A}^P Under Limited Collaboration. In this scenario, \mathcal{A}^P is aware of FLEDGE's limitation and aims to retrieve the first local model ($K = 1$) as described in Sect. 6.1. However, our evaluation showed that FLEDGE is also resilient to this scenario, since \mathcal{A}^P 's probability of success reaches $\sim 0\%$ as defined by $1/m!$.

7.2 FLEDGE Security Analysis

FLEDGE efficiently filters state-of-the-art backdoor injections with the defense deployed in Defender contract. To bypass FLEDGE, an adversary \mathcal{A}^S should inject a poisoned model such that FLEDGE cannot distinguish benign models from poisoned ones. The following elaborate the methodologies that could be used by \mathcal{A}^S to achieve this goal.

\mathcal{A}^S manipulates α . \mathcal{A}^S continuously monitors and adjusts client's loss function to reduce the distance (i.e., cosine or L2 norm) between the client model and the current global model, a larger value of α means more aggressive attacks. Sect. sec:eval-backdoorMitigation demonstrates that FLEDGE can eliminate poisoning attempts efficiently under different values of α for WP and IC applications, respectively.

\mathcal{A}^S manipulates PMR. \mathcal{A}^S would minimize its visibility by increasing the level of control over (or the number of) malicious clients. However, this approach cannot defeat FLEDGE as we empirically proved that FLEDGE maintains high performance w.r.t. changes in PMR for WP (App. H) and IC learning tasks (App. I).

\mathcal{A}^S manipulates PDR. \mathcal{A}^S can also limit the number of poisoned samples by decreasing the PDR value during training to generate less suspicious models. As a result, backdoors (BA) will be reduced.

Additionally, regardless of PDR, FLEDGE continuously filter poisoned models efficiently as shown in App. I.

8 RELATED WORK

Privacy-Preserving Defenses. Multiple approaches have been proposed to protect the privacy of the clients' training data. Passerat *et al.* use a blockchain for privately aggregating the individual models [51]. Ryffel *et al.* propose a framework that eases the use of Secure-Multi-Party Computation (SMPC) for secure aggregation [54], while Fereidooni *et al.* rely only on SMPC [17]. Sav *et al.* use Multiparty Homomorphic Encryption (HE) for collaboratively training a model [55]. Bonawitz *et al.* propose a multi-party-computation scheme based on Shamir's secret sharing [7]. However, as this approach prevents the server from analyzing the local models, it also prevents analyzing them to identify poisoned models. FLEDGE uses Blockchain Two Contract Computation (BT2C) to engage in decentralized privacy-preserving computations based on smart contracts and HE (see Step A in Sect. 4.2). FLEDGE raises accountability by including a reward system that promotes benign contract behavior (see evaluation in Sect. 6.1 and Sect. 6.3).

Poisoning Defenses. Existing defenses against data and model poisoning attacks aim at distinguishing malicious and benign model updates [28, 63] utilizing filter-based approaches (i.e., clustering techniques). However, all of these defenses make specific assumptions about the distribution of benign and malicious data or the characteristics of injected models, causing them to fail if any of these assumptions do not hold. Moreover, such defenses cannot detect stealthy attacks, e.g., where the adversary constrains their poisoned updates within benign update distribution such as constrain-and-Scale attack [5]. Yin *et al.* [68] and Guerraoui *et al.* [20] propose to change aggregation rules to mitigate the effect of malicious model updates. They utilize median parameters from all local models as the global model parameters. Other approaches validate the local models or the aggregated model [61, 64, 69]. However, they cannot detect stealthy backdoors that have only minimal impact on the Main Task Accuracy (MA). Weak Differential Privacy (DP) techniques [39, 46, 65] have also been used to mitigate the effects of poisoning attacks. DP-based defense dilutes the impact of poisoned models by clipping model weights and adding noise to individual clients' model updates. DeepSight provided an efficient filtering that works even in non-IID data scenarios [53] but does not credit the individual contributions and requires a central server than can inspect the model updates. Kalapaaking *et al.* use smart-contract to verify the client-side training process [26] but cannot prevent attackers from manipulating the training data [49]. The use of blockchain implementations [13, 24, 32, 70] have managed to provide defenses against poisoning attacks, however, these solutions only consider untargeted attacks and/or rely on specific assumptions about the distribution of training data. In contrast, FLEDGE effectively removes poisoned models by instantiating a filtering approach based on Gaussian Kernel Density Estimation (G-KDE) function (see Step 4 in Sect. 4.2). The poisoning resiliency of our solution is empirically evaluated on Sect. 6.2, which demonstrates that FLEDGE is an effective solution to mitigate poisoning attacks. In addition, our solution allows us to treat the underlying problem of poisoning attacks, i.e., client training misbehavior, by imposing

monetary deterrents for detected poisoning attempts. This is illustrated in Sect. 6.3.

Hybrid Defenses. A number of existing works recently focused on poisoning attacks while preserving the privacy of the individual model updates. Two works implemented Krum [6] using SMPC [17, 27]. However, Krum focuses on untargeted poisoning attacks and fails to effectively mitigate backdoor attacks (see Sect. 6.2), and SMPC is vulnerable to attacks that limit availability (i.e., DoS attacks). Similarly, FLAME [48] leverages DP with model filtering to provide an efficient defense against backdoors. However, FLAME also relies on SMPC to provide privacy-preserving computations. Several approaches utilize Trusted Execution Environments (TEE) to realize a privacy-preserving poisoning detection [21, 40, 41]. However, requiring TEEs restricts their application to a few scenarios as mobile devices often do not have a TEE, while FLEDGE does not make any assumption about the hardware. Biscotti [56], BEAS [42] and the work of Lu *et al.* [35] are blockchain-based frameworks that target secure and private FL. The first two approaches use Multi-Krum [6] (a Krum variant) to reduce the impact of backdoors in the system, where Multi-Krum also focuses on untargeted poisoning. In particular, Biscotti uses DP and Shamir secrets to preserve privacy of local updates during aggregation. BEAS and the work of Luet *et al.*, on the other hand, leverages DP to obfuscate model weights and provide resiliency against inference attacks. In comparison, FLEDGE provides (1) strong privacy guarantees to client models by obfuscating them using our BT2C protocol (see Sect. 6.1), (2) an effective defense using G-KDE functions to filter different poisoning attacks (see Sect. 6.2), and (3) compensation algorithms to automatically adjust the reward and deter malicious behavior (see Sect. 6.3). Liu *et al.* proposed a smart-contract-based FL framework that utilizes a server-side validation dataset to detect untargeted poisoning attacks [34]. However, assuming a server-side dataset is not practical [53], while FLEDGE detects poisoning attacks without making such an assumption and also includes a reward mechanism to penalize malicious clients.

9 CONCLUSION AND FUTURE WORK

In this paper, we illustrated the current research gaps facing FL systems in terms of security and privacy. To fill in those gaps, we propose FLEDGE, a 3-layer blockchain FL framework powered by smart contracts and HE. Our proposed HE infrastructure, BT2C, enables the analysis and operation of model weights in cipher text through the use of a semi-honest collaboration between smart contracts. BT2C is shown to provide resiliency against white-box inference attacks with a decreased adversarial probability of success of $\frac{1}{m!}$, where m is the number of ciphers inside an encrypted model. Furthermore, our extensive evaluation shows that FLEDGE also offers high resiliency against numerous poisoning strategies (BA \approx 0) with minimal impact on benign accuracy. Our solution to both white-box inference and poisoning attacks allows us to effectively embed incentives or penalizations to both aggregation contracts (e.g., Gateway) and training clients in an effort to minimize adversarial intent via behavior accountability.

Limitations. Although the use of blockchain technology and HE enable the security properties in FLEDGE, they also contribute to

the growth in computation costs of its learning process. For instance, training round latency increases from 15.86s (MNIST model) to 76.6s (CIFAR10 model). This indicates that FLEDGE has an increase in latency of approximately five-times (4.82) for a model that is ten-times larger (10.17). Put differently, FLEDGE offers limited scalability given that its learning process slows down as both model complexity and the number of clients/models increases in the system. Furthermore, the reward mechanism embedded in FLEDGE is directly related to the performance of its defensive strategies. For example, an attacker capable of avoiding FLEDGE's poison defense (i.e., G-KDE Defense) may continue to receive credit even though its model negatively contributes to the learning process.

Future Work. A formal in-depth analysis aimed into the scalability of FLEDGE, e.g., transaction fees, storage costs, computation costs and communication costs, is needed to generate additional insights into the limitations and efficiency of FLEDGE, specifically those imposed by the use of different blockchain platforms.

ACKNOWLEDGMENT

This work was supported in part by the U.S. Department of Energy/National Nuclear Security Administration (DOE/NNSA) #DE-NA0003985, Intel through the Private AI Collaborate Research Institute (<https://www.private-ai.org/>), BMBF and HMWK within ATHENE, as well as from the OpenS3 Lab, the Hessian Ministry of Interior and Sport as part of the F-LION project, following the funding guidelines for cyber security research, the Horizon Europe framework program of the European Union under grant agreement No. 101093126 (ACES). We extend our appreciation to KOBIL GmbH for their support and collaboration throughout the course of this project. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of any of these funding agencies.

REFERENCES

- [1] Abbas Acar, Hidayet Aksu, A Selcuk Uluagac, and Mauro Conti. 2018. A survey on homomorphic encryption schemes: Theory and implementation. *ACM Computing Surveys (Csur)* 51, 4 (2018), 1–35.
- [2] Abbas Acar, Z Berkay Celik, Hidayet Aksu, A Selcuk Uluagac, and Patrick McDaniel. 2017. Achieving secure and differentially private computations in multi-party settings. In *IEEE Symposium on Privacy-Aware Computing (PAC)*. IEEE.
- [3] Sebastien Andreina, Georgia Azzurra Marson, Helen Möllering, and Ghassan Karame. 2021. BaFFLe: Backdoor Detection via Feedback-based Federated Learning. In *ICDCS*.
- [4] Elli Androulaki, Artem Barger, Vita Bortnikov, Christian Cachin, Konstantinos Christidis, Angelo De Caro, David Enyeart, Christopher Ferris, Gennady Laventman, Yacov Manevich, et al. 2018. Hyperledger fabric: a distributed operating system for permissioned blockchains. In *EuroSys conference*.
- [5] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. 2020. How to backdoor federated learning. In *AISTATS*. PMLR.
- [6] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. 2017. Machine learning with adversaries: Byzantine tolerant gradient descent. *Advances in Neural Information Processing Systems (NIPS)* (2017).
- [7] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. 2017. Practical secure aggregation for privacy-preserving machine learning. In *CCS*.
- [8] Joppe W Bos, Kristin Lauter, and Michael Naehrig. 2014. Private predictive analysis on encrypted medical data. *Journal of biomedical informatics* 50 (2014).
- [9] Hao Chen, Ilaria Chillotti, and Yongsoo Song. 2019. Improved bootstrapping for approximate homomorphic encryption. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer.
- [10] Hao Chen, Kim Laine, and Rachel Player. 2017. Simple encrypted arithmetic library-SEAL v2. 1. In *International Conference on Financial Cryptography and Data Security*. Springer.
- [11] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yongsoo Song. 2017. Homomorphic encryption for arithmetic of approximate numbers. In *International conference on the theory and application of cryptology and information security*. Springer, 409–437.
- [12] Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaram. 2014. Project adam: Building an efficient and scalable deep learning training system. In *USENIX OSDI*.
- [13] Harsh Bimal Desai, Mustafa Safa Ozdayi, and Murat Kantarcioglu. 2021. Blockfla: Accountable federated learning via hybrid blockchain architecture. In *ACM conference on data and application security and privacy*.
- [14] Junfeng Fan and Frederik Vercauteren. 2012. Somewhat practical fully homomorphic encryption. *Cryptology ePrint Archive* (2012).
- [15] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. 2020. Local Model Poisoning Attacks to {Byzantine-Robust} Federated Learning. In *USENIX Security*.
- [16] Hossein Fereidooni, Alexandra Dmitrienko, Phillip Rieger, Markus Miettinen, Ahmad-Reza Sadeghi, and Felix Madlener. 2022. Fedcri: Federated mobile cyber-risk intelligence. In *NDSS*.
- [17] Hossein Fereidooni, Samuel Marchal, Markus Miettinen, Azalia Mirhoseini, Helen Möllering, Thien Duc Nguyen, Phillip Rieger, Ahmad-Reza Sadeghi, Thomas Schneider, Hossein Yalame, et al. 2021. SAFElearn: Secure aggregation for private Federated learning. In *IEEE Security and Privacy Workshops (SPW)*. IEEE.
- [18] Clement Fung, Chris JM Yoon, and Ivan Beschastnikh. 2020. The limitations of federated learning in sybil settings. In *International Symposium on Research in Attacks, Intrusions and Defenses (RAID)*.
- [19] Arthur Gervais, Ghassan O Karame, Karl Wüst, Vasileios Glykantzis, Hubert Ritzdorf, and Srdjan Capkun. 2016. On the security and performance of proof of work blockchains. In *CCS*.
- [20] Rachid Guerraoui, Sébastien Rouault, et al. 2018. The Hidden Vulnerability of Distributed Learning in Byzantium. In *ICML*. PMLR.
- [21] Hanieh Hashemi, Yongqin Wang, Chuan Guo, and Murali Annavaram. 2021. Byzantine-robust and privacy-preserving framework for fedml. *arXiv preprint arXiv:2105.02295* (2021).
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.
- [23] Md Tamjid Hossain, Shafkat Islam, Shahriar Badsha, and Haoting Shen. 2021. Desmp: Differential privacy-exploited stealthy model poisoning attacks in federated learning. In *International Conference on Mobility, Sensing and Networking (MSN)*. IEEE.
- [24] Gaofeng Hua, Li Zhu, Jinsong Wu, Chunzi Shen, Linyan Zhou, and Qingqing Lin. 2020. Blockchain-based federated learning for intelligent control in heavy haul railway. *IEEE Access* 8 (2020), 176830–176839.
- [25] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. 2021. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning* 14, 1–2 (2021), 1–210.
- [26] Aditya Pribadi Kalapaaking, Ibrahim Khalil, and Mohammed Atiquzzaman. 2023. Smart Policy Control for Securing Federated Learning Management System. *IEEE Transactions on Network and Service Management* (2023).
- [27] Youssef Khazbak, Tianxiang Tan, and Guohong Cao. 2020. MLGuard: Mitigating poisoning attacks in privacy preserving distributed collaborative learning. In *ICCCN*. IEEE.
- [28] Yein Kim, Huili Chen, and Farinaz Koushanfar. 2022. Backdoor Defense in Federated Learning Using Differential Testing and Outlier Detection. *arXiv preprint arXiv:2202.11196* (2022).
- [29] Sunny King and Scott Nadal. 2012. Ppcoin: Peer-to-peer crypto-currency with proof-of-stake. *self-published paper, August* 19, 1 (2012).
- [30] Kavita Kumari, Phillip Rieger, Hossein Fereidooni, Murtuza Jadhwal, and Ahmad-Reza Sadeghi. 2023. BayBFed: Bayesian Backdoor Defense for Federated Learning. In *IEEE S&P*.
- [31] Huimin Li, Phillip Rieger, Shaza Zeitouni, Stjepan Picek, and Ahmad-Reza Sadeghi. 2023. FLAIRS: FPGA-Accelerated Inference-Resistant & Secure Federated Learning. *International Conference on Field-Programmable Logic and Applications (FPL)* (2023).
- [32] Yuzheng Li, Chuan Chen, Nan Liu, Huawei Huang, Zibin Zheng, and Qiang Yan. 2020. A blockchain-based decentralized federated learning framework with committee consensus. *IEEE Network* 35, 1 (2020).
- [33] Pengrui Liu, Xiangrui Xu, and Wei Wang. 2022. Threats, attacks and defenses to federated learning: issues, taxonomy and perspectives. *Cybersecurity* 5, 1 (2022).
- [34] Yi Liu, Jialiang Peng, Jiawen Kang, Abdullah M Ilyyasu, Dusit Niyato, and Ahmed A Abd El-Latif. 2020. A secure federated learning framework for 5G networks. *IEEE Wireless Communications* 27, 4 (2020), 24–31.
- [35] Yunlong Lu, Xiaohong Huang, Yueyue Dai, Sabita Maharjan, and Yan Zhang. 2019. Blockchain and federated learning for privacy-preserved data sharing in industrial IoT. *IEEE Transactions on Industrial Informatics* 16, 6 (2019), 4177–4186.
- [36] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*. PMLR.
- [37] Brendan McMahan and Daniel Ramage. 2017. Federated learning: Collaborative Machine Learning without Centralized Training Data. Google AI.

- [38] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. 2017. Learning differentially private recurrent language models. *ICLR* (2017).
- [39] Lu Miao, Wei Yang, Rong Hu, Lu Li, and Liusheng Huang. 2022. Against Backdoor Attacks In Federated Learning With Differential Privacy. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE.
- [40] Fan Mo, Hamed Haddadi, Kleomenis Katevas, Eduard Marin, Diego Perino, and Nicolas Kourtellis. 2021. PPFL: privacy-preserving federated learning with trusted execution environments. In *Annual international conference on mobile systems, applications, and services*.
- [41] Arup Mondal, Yash More, Ruthu Hulikal Rooparagunath, and Debayan Gupta. 2021. Flatee: Federated learning across trusted execution environments. *arXiv preprint arXiv:2111.06867* (2021).
- [42] Arup Mondal, Harpreet Virk, and Debayan Gupta. 2022. BEAS: Blockchain Enabled Asynchronous & Secure Federated Machine Learning. *arXiv preprint arXiv:2202.02817* (2022).
- [43] Luis Muñoz-González, Kenneth T Co, and Emil C Lupu. 2019. Byzantine-robust federated machine learning through adaptive model averaging. *arXiv preprint arXiv:1909.05125* (2019).
- [44] Michael Naehrig, Kristin Lauter, and Vinod Vaikuntanathan. 2011. Can homomorphic encryption be practical?. In *ACM workshop on Cloud computing security workshop*.
- [45] Satoshi Nakamoto. 2019. *Bitcoin: A peer-to-peer electronic cash system*. Technical Report. Manubot.
- [46] Mohammad Naseri, Jamie Hayes, and Emiliano De Cristofaro. 2022. Local and central differential privacy for robustness and privacy in federated learning. *NDSS* (2022).
- [47] Thien Duc Nguyen, Samuel Marchal, Markus Miettinen, Hossein Fereidooni, N. Asokan, and Ahmad-Reza Sadeghi. 2019. DfOT: A Federated Self-learning Anomaly Detection System for IoT. In *ICDCS*.
- [48] Thien Duc Nguyen, Phillip Rieger, Huili Chen, Hossein Yalame, Helen Möllering, Hossein Fereidooni, Samuel Marchal, Markus Miettinen, Azalia Mirhoseini, Shaza Zeitouni, et al. 2022. FLAME: Taming Backdoors in Federated Learning. *USENIX Security* (2022).
- [49] Thien Duc Nguyen, Phillip Rieger, Markus Miettinen, and Ahmad-Reza Sadeghi. 2020. Poisoning attacks on federated learning-based IoT intrusion detection system. In *NDSS Workshop on Decentralized IoT Systems and Security*.
- [50] Diego Ongaro and John Ousterhout. 2015. The raft consensus algorithm. (2015).
- [51] Jonathan Passerat-Palmbach, Tyler Farnan, Robert Miller, Marielle S Gross, Heather Leigh Flannery, and Bill Gleim. 2019. A blockchain-orchestrated federated learning architecture for healthcare consortia. *arXiv preprint arXiv:1910.12603* (2019).
- [52] Phillip Rieger, Torsten Krauß, Markus Miettinen, Alexandra Dmitrienko, and Ahmad-Reza Sadeghi. 2024. CrowdGuard: Federated Backdoor Detection in Federated Learning. *NDSS* (2024).
- [53] Phillip Rieger, Thien Duc Nguyen, Markus Miettinen, and Ahmad-Reza Sadeghi. 2022. DeepSight: Mitigating Backdoor Attacks in Federated Learning Through Deep Model Inspection. *NDSS* (2022).
- [54] Theo Ryyffel, Andrew Trask, Morten Dahl, Bobby Wagner, Jason Mancuso, Daniel Rueckert, and Jonathan Passerat-Palmbach. 2018. A generic framework for privacy preserving deep learning. *arXiv preprint arXiv:1811.04017* (2018).
- [55] Sinem Sav, Apostolos Pyrgelis, Juan R Troncoso-Pastoriza, David Froelicher, Jean-Philippe Bossuat, Joao Sa Sousa, and Jean-Pierre Hubaux. 2021. POSEIDON: privacy-preserving federated neural network learning. *NDSS* (2021).
- [56] Muhammad Shayam, Clement Fung, Chris JM Yoon, and Ivan Beschastnikh. 2020. Biscotti: A blockchain system for private and secure federated learning. *IEEE Transactions on Parallel and Distributed Systems* 32, 7 (2020), 1513–1525.
- [57] Micah J Sheller, G Anthony Reina, Brandon Edwards, Jason Martin, and Spyridon Bakas. 2018. Multi-institutional deep learning modeling without sharing patient data: A feasibility study on brain tumor segmentation. In *International MICCAI Brainlesion Workshop*. Springer.
- [58] Shiqi Shen, Shruti Tople, and Prateek Saxena. 2016. Auror: Defending against poisoning attacks in collaborative deep learning systems. In *Annual Conference on Computer Security Applications*.
- [59] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In *IEEE S&P*. IEEE.
- [60] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [61] Dimitris Stripelis, Marcin Abram, and Jose Luis Ambite. 2022. Performance Weighting for Robust Federated Learning Against Corrupted Sources. *arXiv preprint arXiv:2205.01184* (2022).
- [62] Asher Trockman and J Zico Kolter. 2022. Patches Are All You Need? *arXiv preprint f:2201.09792* (2022).
- [63] Yongkang Wang, Dihua Zhai, Yufeng Zhan, and Yuanqing Xia. 2022. RFLBAT: A Robust Federated Learning Algorithm against Backdoor Attack. *arXiv preprint arXiv:2201.03772* (2022).
- [64] Binhan Xi, Shaofeng Li, Jiachun Li, Hui Liu, Hong Liu, and Haojin Zhu. 2021. BatFL: Backdoor Detection on Federated Learning in e-Health. In *IEEE/ACM International Symposium on Quality of Service (IWQoS)*. IEEE.
- [65] Chulin Xie, Minghao Chen, Pin-Yu Chen, and Bo Li. 2021. Crfl: Certifiably robust federated learning against backdoor attacks. In *ICML*. PMLR.
- [66] Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. 2019. Dba: Distributed backdoor attacks against federated learning. In *ICLR*.
- [67] Xun Yi, Russell Paulet, and Elisa Bertino. 2014. Homomorphic encryption. In *Homomorphic encryption and applications*. Springer.
- [68] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. 2018. Byzantine-Robust Distributed Learning: Towards Optimal Statistical Rates. In *ICML*. PMLR.
- [69] Lingchen Zhao, Shengshan Hu, Qian Wang, Jianlin Jiang, Chao Shen, Xiangyang Luo, and Pengfei Hu. 2020. Shielding collaborative learning: Mitigating poisoning attacks through client-side detection. *IEEE Transactions on Dependable and Secure Computing* 18, 5 (2020).
- [70] Sicong Zhou, Huawei Huang, Wuhui Chen, Pan Zhou, Zibin Zheng, and Song Guo. 2020. Pirate: A blockchain-based secure framework of distributed machine learning in 5g networks. *IEEE Network* 34, 6 (2020).

APPENDIX

A BACKGROUND ON FEDERATED LEARNING

Federated Learning (FL) [36] is a round-based machine learning protocol with the purpose of providing a collaborative learning surface between K clients and a central aggregation server \mathcal{S} . For a training round $t \in [1, T]$ in FedAVG [36], each client downloads an initial deep learning model from \mathcal{S} known as the global model G_{t-1} . Clients begin training using its local data D_i and send local updates W_i , where $i \in [1, K]$ to be averaged by \mathcal{S} resulting in a new global model G_t , such that $G_t = \sum_{i=1}^K \frac{n_i}{n} W_i$, where $n_i = \|D_i\|$, and $n = \sum_{i=1}^K n_i$. The process continuously repeats for T rounds until the model performance reaches an appropriate/target level. Similar to previous work [5, 58], we weight all clients equally such that $G_t = \sum_{i=1}^K \frac{W_i}{K}$.

B BACKGROUND ON BLOCKCHAIN

Blockchain is a decentralized storage system sustained over a vast network of computers (peers) with the common goal to serve as an immutable record [45]. The record is formed from a series of cryptographically linked *blocks*, where each block contains multiple transactions (or data samples). The transaction order is defined by a systematized ordering protocol. The protocol is also known as the consensus process (i.e., proof of work (PoW) [19] or proof of stake (PoS) [29]) that uses a subset of peers to approve and validate new blocks of data. Informally, the blockchain participants broadcast data received by consensus nodes as transactions. These nodes gather different transactions to formulate a new block during the consensus process. After consensus is reached, the block is distributed to every other participant so that it can be validated and appended to the ledger.

Smart contracts of blockchain systems allow authorized users to retrieve any data from and submit any data to the underlying blockchain system. They provide a verifiable and transparent environment dedicated to build trust among untrusted parties without an intermediary. The use of smart contracts in the FL context enables a decentralized storage to merge and analyze client models in an automated fashion. Currently, blockchains can be categorized as either a public or private.¹¹ For our solution we consider a private blockchain to have additional control over the visibility of the ledger as only specific entities are selected to join the network.

¹¹Public (or permissionless) blockchains allow anyone to join and verify the system, whereas private (or permissioned) blockchains do not.

Hyperledger Fabric (HLF) offers high degree of flexibility and scalability [4] while utilizing a modular approach that is easy-to-use. We use HLF as the blockchain platform in this study because of the following reasons.

- *Scalability.* Compared with public blockchains that require every consensus node to validate newly appended blocks, Fabric controls the number of consensus nodes to improve its throughput and significantly reduce the communication latency.
- *Leader-based consensus process.* HLF uses Raft [50], a fault-tolerant algorithm that requires nodes to elect a leader to process and distribute data. This leader-based approach reduces data processing times and makes the blockchain systems time/computational efficient.

C BACKGROUND ON HOMOMORPHIC ENCRYPTION

Homomorphic Encryption (HE) [1, 67] is a special public key encryption. It allows the application of mathematical functions on encrypted data so that data can remain confidential while being processed. HE has two main branches, namely fully homomorphic encryption (FHE) and partially homomorphic encryption (PHE). FHE allows users to compute *any* mathematical function while PHE only allows *some* functions to be applied [9, 14]. Prominent examples of PHE are Cheon-Kim-Kim-Song (CKKS) [11] and Brakerski/Fan-Vercauteren (BFV) [14] schemes.

CKKS Encryption Scheme. Computing only addition and multiplication operations is a prominent example of practical PHE [9]. Unlike other PHE (i.e., BFV [14]), CKKS allows floating point operations, which makes it applicable for encrypting data and models for FL systems. In typical CKKS implementations, the encryption process begins when a client first generates a key pair, a secret-key S_k and a public-key P_k . The client uses P_k to encrypt its data D_i before sharing it to an untrusted server S_u to perform private computations on encrypted data D_i^* . Computed data is then returned to the client for decryption with S_k , unveiling the real result after computation. Previous research [8, 44] demonstrated the practicality of HE in applications where privacy is preferred over efficiency.

D INTUITION FOR G-KDE CLUSTERING

Poisoning attempts are known to generate a sufficiently large gap between the distance scores of benign models and malicious updates such that our G-KDE clustering is able to filter malicious models regardless of the obfuscation technique used by an adversary (e.g., manipulation of loss function [5]). This hypothesis is validated and shown in Fig. 4, where four subfigures illustrate the distance score distribution found for both word prediction tasks (WP) on the left, and image classification tasks (IC) on the right. Fig. 4a uses a histogram to prove the existence of a gap within distance scores of WP tasks. Similarly, Fig. 4b displays a similar gap between the malicious and benign scores obtained from an IC task. These gaps indicate that it is possible to discriminate models based on where they reside within the data distribution regardless of model complexity, making this approach generic enough to be applicable as a poisoning defense. Fig. 4c and d illustrate how the use of G-KDE

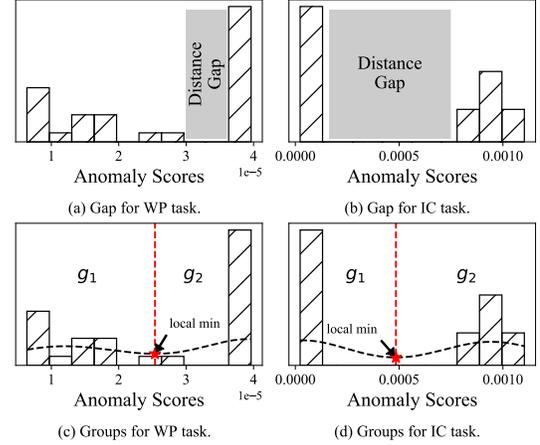


Figure 4: Distribution analysis of anomaly scores for WP and IC datasets. (a) A gap between malicious and benign models for WP. (b) A gap between malicious and benign models for IC. (c) Groups found during WP analysis. (d) Groups found during IC analysis.

functions finds multiple distributions inside the distance scores, thus, generating two groups, e.g., g_1 (benign) and g_2 (malicious).

E EVALUATION METRICS

We consider the following metrics to define the effectiveness of poisoning attacks:

- *Main Task Accuracy (MA).* It indicates the accuracy level of a model when performing its main (or benign) task. The goal of an adversary is to minimize its impact such that MA is not degraded during malicious training.
- *Backdoor Accuracy (BA).* It is used to measure the accuracy of the backdoor injected into a model. The goal of the attacker is to maximize BA during a training round.
- *True Positive Rate (TPR).* It indicates the accuracy of the defenses for detecting poisoned models such that True Positives (TP) are models correctly identified as malicious, and False Negatives (FN) are those mislabeled as benign: $TPR = TP / (TP + FN)$. A high TPR suggests poisoned models are being removed.
- *True Negative Rate (TNR).* It determines the accuracy of the defense for detecting benign models such that True Negatives (TN) are models correctly labeled as benign, and False Positives (FP) are those incorrectly classified as malicious: $TNR = TN / (TN + FP)$. TNR increases as the defense removes less benign models.

F RUNTIME-PERFORMANCE

To assess the latency aspects of FLEDGE during any given round, we evaluate the following processes in terms of latency: model encryption, model upload, model filtering and model aggregation. Model encryption determines the time required to fully produce an encrypted model. Model upload defines the average time an encrypted model is committed and analyzed by FLEDGE such that

it produces the distance score that corresponds to the uploaded encrypted model. Model filtering evaluates how long the poison defense takes to analyze anomaly scores to remove models from aggregation. Model aggregation determines the time it takes to produce a new global model. We perform experiments using MNIST, Fashion, and CIFAR-10 datasets.

Effect of Model Complexity. For the following experiment we set the number of clients to be 50 to assess the impact of model complexity on system latency. Table 4 compiles the preliminary evaluation of FLEDGE. For model encryption, we found there is negligible impact on latency with values of 0.12, 0.14, 0.56 and 1.03 seconds for each of the models evaluated. Note that model encryption behaves as $O(n^2)$ as the number of trainable parameters increases, i.e., $\sim 20.4M$ params are encrypted (10K ciphers) in $\sim 110s$. Moreover, model upload shows an average latency of 5.23, 5.85, and 17.57 seconds for MNIST, Fashion, and CIFAR-10 tasks. The increased latency is mainly attributed to the distance score computation, where we use three BT2C rounds to compute the cosine distance. The values displayed for model filtering are close to 2s in each evaluated task. The reason is because this is the only process where HE is not used, as the Defender contract only analyzes the anomaly scores previously computed during model upload. In contrast, model aggregation is determined to be the most computationally expensive process in FLEDGE, showing values of 8.44, 9.73, and 56.09 seconds, respectively. As a consequence, FLEDGE successfully completes a single training round in 15.86s (MNIST), 17.7s (Fashion), and 76.6s (CIFAR-10). Therefore, the use of FLEDGE is best suited to protect the privacy of constrained environments, where local clients do not have the computational resources (e.g., GPU) to train a robust model with millions of parameters such as ResNet [22] and VGG [60]. However, they can operate lite models, i.e., lite-CNN and/or ConvMixer, to achieve good model performance while delegating every intensive computational task to blockchain via smart contracts.

Effect of Number of Clients. In this experiment, we narrow our focus to the model aggregation process since it was determined to be the most computationally expensive process in FLEDGE. To further inspect this process, we vary the number of clients to {10, 30, 50}. Fig. 5 shows the impact of model aggregation for the distinct learning tasks using a different number of clients. Here, we can visualize that not only does model aggregation scales according to model complexity, but it also shows a linear dependency w.r.t. the number of clients. Put differently, model aggregation is the major contributor of latency in FLEDGE as other components (i.e., model encryption, model upload) are only defined by the computational requirements of HE.

G ROBUSTNESS AGAINST CLIENT DROPOUTS

The goal of the training client is to train and submit its local update to the Gateway contract. Uploading the model can be seen as an automatic operation, either the model is completely received and further processed or it is ignored. Once a local model is successfully uploaded, a TT2 is formulated to generate a record in the ledger. At this stage, the model is now considered to be part of FLEDGE. Otherwise, the encrypted model is ignored during the training round. After the record for a model is part of the ledger, a client

Table 4: Training round efficiency of FLEDGE for the following processing tasks in seconds: Model Encryption, Model Upload, Model Filtering, and Model Aggregation.

| Dataset | MNIST | Fashion | CIFAR-10 |
|--------------------------|-------|---------|----------|
| Model Encryption | 0.12 | 0.14 | 1.03 |
| Model Upload | 5.23 | 5.85 | 17.57 |
| Model Filtering | 2.07 | 1.98 | 1.91 |
| Model Aggregation | 8.44 | 9.73 | 56.09 |
| Total | 15.86 | 17.7 | 76.6 |

dropout does not affect FLEDGE anymore. FLEDGE takes advantage of the visibility of blockchain to consistently account for the number of encrypted models currently present. This allows our approach to perform the aggregation using only the available models. Thus, making it robust against dynamic dropouts.

H BACKDOOR EVALUATION FOR WP TASK

The following section evaluates the effect of PMR and α in FLEDGE for the WP task. We analyze the behavior of FLEDGE during 10 training rounds. The attacker attempts to inject backdoored models in different training rounds to compromise the global model. We analyze FLEDGE in terms of MA, BA, TPR and TNR, and draw a comparison with a typical *No Defense* scenario using a FL environment with 30 training clients. The results of the direct comparison are shown in Fig. 6, and are elaborated as follows.

Effect of PMR Rate. In this experiment, we modify the control ratio of the attacker in the system by setting PMR values of {0.1, 0.3, 0.5} (or 6, 9 and 15 malicious clients). In Fig. 6a, we observe the ability of the adversary to successfully inject backdoors (BA=100% for *No Defense*) into the model regardless on the number of malicious clients it controls. In contrast, FLEDGE is demonstrated to effectively remove the threat of backdoors (BA=0%) while preserving the benign accuracy of MA=22.7%. FLEDGE achieves a TNR and TPR values of 100% in the process, which denotes that FLEDGE is able to defend against multiple malicious clients.

Effect of α Rate. Further, we modified the level of intensity for backdoor injection such that $\alpha = \{0.2, 0.5, 0.9\}$. Fig. 6b shows that BA and MA (*No Defense*) are not affected by α with an average value of 100% and 22.6%, respectively. FLEDGE effectively mitigates every backdoor attempt (BA=0%, TPR=100%). TNR is 86.66% when $\alpha = 0.2$,

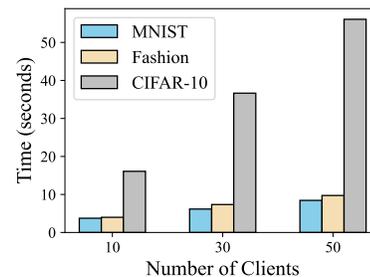


Figure 5: Effect of number of training clients on Model Aggregation process.

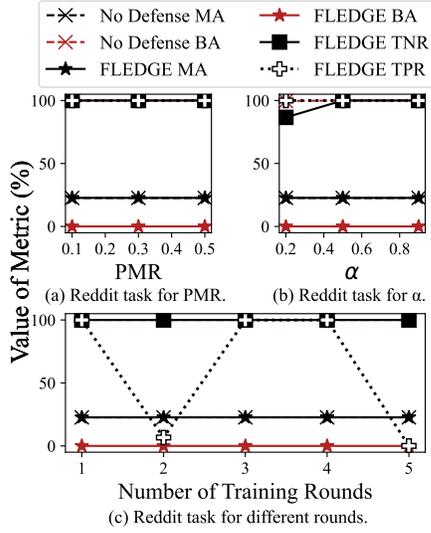


Figure 6: Effect of backdoors on evaluation metrics for WP task. Effect of (a) PMRs, (b) α , (c) multiple injections.

and this value increases to 100% when $\alpha=\{0.5, 0.9\}$. FLEDGE filters 2 out of 15 (13.34%) benign models when defending against $\alpha=0.2$, and it continues to provide an effective defense against different levels of aggression when α value increases.

Effect of Backdoor Injections for Multiple Rounds. For this experiment, we set PMR=0.5 and the number of rounds to be 5, where rounds 1,3 and 4 are poisoned, and rounds 2 and 5 are benign. This is because we are interested in testing the efficacy of FLEDGE w.r.t. multiple attack rounds. Fig. 6c illustrates the operation of FLEDGE under multiple attack rounds. Note that for all 5 rounds, FLEDGE continues to maintain BA=0%, whereas the *No Defense* setup has BA=100%. Further, TNR values are always 100%, however, TPR values differ depending if the round is benign (TPR \approx 0%) or malicious (TPR=100%). This behavior indicates that FLEDGE adjusts its filtering mechanics to only remove poisoned models and avoid benign models.

I BACKDOOR EVALUATION FOR IC TASK

In this section, we evaluate the effect of PMR, PDR, non-IID rate and α on IC tasks during a single training round. We illustrate the difference in behavior for FLEDGE and *No Defense* system during multiple attack settings. We use MA, BA, TPR, and TNR as metrics for comparison. For the purposes of demonstration, we focus on CIFAR-10 as it is the most complex IC task. We select 50 as the total number of clients. The results of the experiments are illustrated in Fig. 7.

Effect of PMR Rate. The following experiment changes the influence of the attacker over the system by setting PMR values of $\{0.1, 0.3, 0.5\}$. Fig. 7a shows the evaluation for PMR. At the *No Defense* setup, we clearly observe the attacker negatively affects MA from 83.9% to 66.1% to reach an appropriate BA level (BA=100%). Nevertheless, FLEDGE efficiently mitigated this, since FLEDGE continues to filter malicious updates (BA \approx 0%), reaching TNR and TPR of 100%. Consequently, this elevates MA to a benign level,

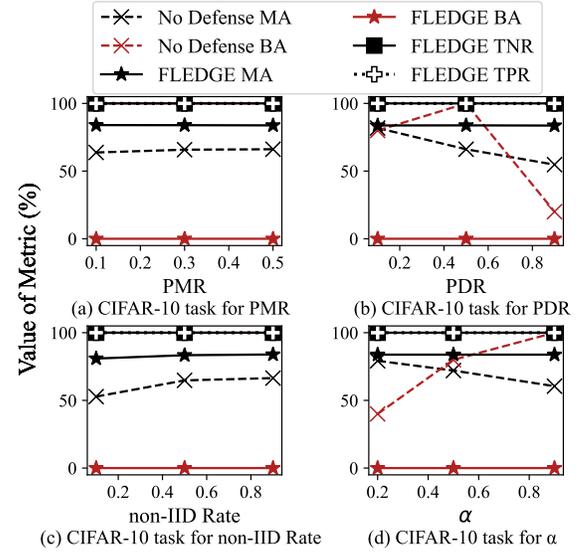


Figure 7: Effect of backdoors on evaluation metrics for IC task. Effect of (a) PMRs, (b) PDRs, (c) non-IID rates, (d) α .

i.e., MA=83.94%. Therefore, we further demonstrate that FLEDGE is resistant to backdoor injections where the adversary is able to change the number of clients it controls.

Effect of PDR Rate. In this experiment, we evaluate the influence of PDR in FLEDGE by setting PDR= $\{0.1, 0.5, 1\}$ as illustrated in Fig. 7b. Similar to previous experiments, Fig. 7b shows the evaluation metrics for different PDR values. Here, we observe the following BA behaviors. First, it shows a mild drop in performance for PDR=0.1, i.e., BA=80% and MA=81.96%. Second, it shows a significant degradation in performance when PDR=1, i.e., BA=20% and MA=54.8%. This is, however, filtered by FLEDGE which effectively reduces BA to 0% with TNR=100% and TPR=100%. Therefore, we consider FLEDGE to be resistant to changes in PDR.

Effect of non-IID Data. For the following experiment, we aim to analyze FLEDGE under different data concentrations such that non-IID is set to $\{0, 0.5, 1\}$. Fig. 7c show that backdoors continue to be effective for every non-IID value, with BA of 100%. These results also show that MA increases from 52.8% to 66.5% for non-IID of 0 and 1, respectively. However, FLEDGE minimizes the impact of backdoors (BA \approx 0%) for every non-IID setting with TNR and TPR of 100%. This indicates that FLEDGE is highly resilient to changes caused by different non-IID rates.

Effect of α Rate. For the next experiment, we set α to be $\{0.2, 0.5, 0.9\}$ in order to test the performance of FLEDGE under different intensity levels. Similar to previous experiments, Fig. 7d shows the evaluation metrics. In here, we observe how BA is directly proportional to α , i.e., BA of 40% ($\alpha=0.2$) increases to 100% ($\alpha=0.9$). Consequently, this also has an impact on MA, yielding reduced values of 60.5% when $\alpha=0.9$. FLEDGE removed poisoned models (TPR=100%) while preserving benign ones (TNR=100%) such that

BA=0% and MA is returned to its benign level, i.e., 83.8%. Hence, we further determine that FLEDGE is robust against changes in α .