# Source Inference Attacks in Federated Learning

**Hongsheng Hu**
University of Auckland
New Zealand
hhu603@aucklanduni.ac.nz

**Zoran Salcic**
University of Auckland
New Zealand
z.salcic@auckland.ac.nz

**Lichao Sun**
Lehigh University
USA
lis221@lehigh.edu

**Gillian Dobbie**
University of Auckland
New Zealand
g.dobbie@auckland.ac.nz

**Xuyun Zhang**[*]
Macquarie University
Australia
xuyun.zhang@mq.edu.au

## Abstract

Federated learning (FL) has emerged as a promising privacy-aware paradigm that allows multiple clients to jointly train a model without sharing their private data. Recently, many studies have shown that FL is vulnerable to *membership inference attacks* (MIAs) that can distinguish the training members of the given model from the non-members. However, existing MIAs ignore the source of a training member, *i.e.*, the information of which client owns the training member, while it is essential to explore source privacy in FL beyond membership privacy of examples from all clients. The leakage of source information can lead to severe privacy issues. For example, identification of the hospital contributing to the training of an FL model for COVID-19 pandemic can render the owner of a data record from this hospital more prone to discrimination if the hospital is in a high risk region. In this paper, we propose a new inference attack called *source inference attack* (SIA), which can derive an optimal estimation of the source of a training member. Specifically, we innovatively adopt the Bayesian perspective to demonstrate that an honest-but-curious server can launch an SIA to steal non-trivial source information of the training members without violating the FL protocol. The server leverages the prediction loss of local models on the training members to achieve the attack effectively and non-intrusively. We conduct extensive experiments on one synthetic and five real datasets to evaluate the key factors in an SIA, and the results show the efficacy of the proposed source inference attack.

## 1 Introduction

Big data and deep learning technologies have enabled us to perform scalable data mining across multiple parties to build powerful prediction models. For example, it will be very appealing for different countries to collaborate, utilizing their medical data records to train prediction models for fighting against the COVID-19 pandemic. However, many countries or regions have issued strong privacy protection laws and regulations, such as GDPR [28], and it is very difficult to straightforwardly collect and combine the data from different parties for a data mining task. To circumvent this major obstacle towards big data mining, a novel machine learning (ML) paradigm named feaderated learning (FL) has recently been proposed, which allows multiple clients coordinated by a central server to train a joint ML model in an iterative manner [22, 10, 11]. In FL, no client can access any training data owned by other clients, leading to a privacy-aware paradigm for collaborative model

---

[*]Corresponding author.

training. Specific to the example mentioned above, FL can greatly facilitate the scenario where many hospitals hope to build a joint COVID-19 diagnosis ML model from their distributed data. A real-life case has been shown in [40], where FL has been successfully adopted to build a promising ML model for COVID-19 diagnosis, with the use of the geographically distributed chest CT (Computed Tomography) data collected from patients at different hospitals.

However, many recent studies [25, 26, 43, 37, 34] have demonstrated that FL fails to provide sufficient privacy guarantees, as sensitive information can be revealed in the training process. In FL, multiple clients send ML model weight or gradient updates derived from local training to a central server for global model training. The communication of model updates renders FL vulnerable to several recently developed privacy attacks, such as property inference attacks [8], reconstruction attacks [12], and membership inference attacks (MIAs) [30]. Among these attacks, MIAs aim to identify whether or not a data record was in the training dataset the model was built on (*i.e.,* a member). This can impose severe privacy risks on individuals. For example, via identifying the fact that a clinical record that has been used to train a model associated with a certain disease, MIAs can infer that the owner of the clinical record has a high chance of having the disease.

However, existing MIAs ignore the source of a training member, i.e., the information of which client owns the training member, while MIAs against FL models distinguish the training members of the model from the non-members. It is essential to explore source privacy in FL beyond membership privacy, because the leakage of such information can lead to further privacy issues. For instance, in the scenario where multiple hospitals jointly train an FL model for COVID-19 diagnosis, MIAs can only reveal who have been tested for COVID-19, but the further identification of the source hospital where the people are from will make them more prone to discrimination, especially when the hospital is in a high risk region or country [5].

In this paper, we propose a novel inference attack called *Source Inference Attack* (SIA) in the context of FL. SIA aims to determine which client owns a training record in FL. In practice, the SIA can be considered as a natural extension beyond MIAs, i.e., after determining which data instances are training members in MIAs, the adversary can further conduct the SIA to identify which client it comes from. To be practical, it is assumed that the adversary is an *honest-but-curious* central server, who knows the identities of clients and receives the updates from them. It is worth noting that the server can infer client-private information without interfering with the FL training nor affecting the model prediction performance. While the adversary can be one of the clients, we argue that it is impractical for her to launch SIAs as she knows little about other clients' identities and can only access the joint models.

Specifically, we innovatively explore the SIA from the Bayesian perspective, and demonstrate that a server can achieve the optimal estimate of the source of a training member in an SIA without violating the FL protocol. To this end, the prediction loss of local models on the training members is utilised to obtain the source information of the training members effectively and non-intrusively. Besides theoretical formulation, we empirically evaluate the SIA in FL trained with one synthetic and five real world datasets, with respect to several FL aspects such as data distributions across clients, the number of clients, and the number of local epochs. The experiment results validate the efficacy of the proposed source inference attack under various FL settings. An important finding is that the success of an SIA is directly relevant to the generalizability of local models and the diversity of the local data.

Our main contribution is multifold, summarized as follows.

- First, we propose the source inference attack (SIA), a novel inference attack in FL that identifies the source of a training member. SIA can further breach the privacy of training members beyond membership inference attacks.

- Second, we adopt the Bayesian perspective to demonstrate that an honest-but-curious central server can fulfil an effective SIA in a non-intrusive manner by optimally estimating the source of a training member, using prediction loss of local models.

- Last, we perform an extensive empirical evaluation on both synthetic and real world datasets under various FL settings, and the results validate the efficacy of the proposed SIA.

We provide all proofs in the full version of our paper, and our source code is available at: `https://github.com/HongshengHu/source-inference-FL`.

## 2 Preliminaries

In this section, we briefly review the background of the federated learning and membership inference attacks.

### 2.1 Federated Learning

Federated learning allows multiple clients to jointly train an ML model in an interactive manner. It is an attractive framework for training ML models without direct access to diverse training data owned by different clients, especially for privacy-sensitive tasks [22, 42, 25, 1]. The *federated averaging* (FedAvg) [22] algorithm is the first and perhaps the most widely used FL algorithm. During multiple rounds of communication between server and clients, a central model is trained. At each communication round, the server distributes the current central model to local clients. The local clients then perform local optimization using their own data. To minimize communication, clients might update the local model for several epochs during a single communication round. Next, the optimized local models are sent back to the server, who average them to allocate a new central model. The performance of the new central model decides the training is either stopped or a new communication round starts. In FL, clients never share data, only their model weights or gradients.

### 2.2 Membership Inference Attacks

Membership inference attacks aim to identify whether a data record was part of the target model's training dataset or not. Shokri et al. [30] present the first MIAs against ML models. Specifically, they demonstrate that an adversary can tell whether a data record has been used to train a classifier or not, solely based on the prediction vector of the data record. Since then, a growing body of work further investigates and explores the feasibility of MIAs on various ML models [13]. Nevertheless, recent works [25, 26] have demonstrated the success of MIAs on FL models. For example, Melis et al. [25] have shown that an adversary can infer whether a specific location profile was used to train an FL model on the FourSquare location dataset with high success rate. Although MIAs can distinguish the training members of the FL model from the non-members, existing inference attacks ignore to further explore which client owns the training member identified by MIAs. In this paper, we fill this gap and show the possibilities of breaching the source privacy of training members.

## 3 Source Inference Attacks

In this section, we formally analyze how an honest-but-curious sever in FL can optimally estimate the source of a training member from the Bayesian perspective.

We focus on the supervised learning of classification tasks. The adversary is the honest-but-curious server who faithfully implements FedAvg while trying to determine where a training data record comes from. Assuming the whole training dataset consists of $n$ i.i.d. data records $z_1, \cdots, z_n$ from a data distribution. Each record is represented as $z = (x, y)$ where $x$ is an input vector and $y$ is the class label. The source status of each record is represented by a $K$-dimensional (assuming there are $K$ clients) multinomial vector $s$ in which one of the elements $s_k$ equals 1, and all remaining elements equal 0. We assume that multinomial source variables $s_1, \cdots, s_n$ are independent, and the training record $z_i$ comes from the client k with the probability $\mathbb{P}(s_{ik} = 1) = \lambda$. Without loss of generality, taking the case of $z_1$, the source inference is defined as follows:

**Definition 1 (Source inference)** *Given local optimized model $\theta_k$, a training record $z_1$, source inference aims to infer the posterior probability of $z_1$ belonging to the client k:*

$$\mathcal{S}(\theta_k, z_1) := \mathbb{P}(s_{1k} = 1 | \theta_k, z_1). \tag{1}$$

For the source inference by *Definition 1*, we want to derive the explicit formula for $\mathcal{S}(\theta_k, z_1)$ from the Bayesian perspective, which establishes the optimal limit that our source inference can achieve. We denote $\tau = \{z_2, \cdots, z_n, s_2, \cdots, s_n\}$ as the set which collects the knowledge about the other training records and their source status. The explicit formula of $\mathcal{S}(\theta_k, z_1)$ is given by the following theorem.

**Theorem 1** *Given local optimized model $\boldsymbol{\theta}_k$, a training record $\boldsymbol{z}_1$, the optimal source inference is given by:*

$$S(\boldsymbol{\theta}_k, \boldsymbol{z}_1) = \mathbb{E}_{\boldsymbol{\tau}} \left[ \sigma \left( \log(\frac{\mathbb{P}(\boldsymbol{\theta}_k | s_{1k} = 1, \boldsymbol{z}_1, \boldsymbol{\tau})}{\mathbb{P}(\boldsymbol{\theta}_k | s_{1k} = 0, \boldsymbol{z}_1, \boldsymbol{\tau})}) + \mu_\lambda \right) \right], \tag{2}$$

where $\mu_\lambda = \log(\frac{\lambda}{1-\lambda})$ and $\sigma(\cdot)$ is the sigmoid function.

We observe that *Theorem 1* does not have the loss $\ell(\cdot)$ form and only relies on the posterior parameter $\boldsymbol{\theta}_k$ in expectation given $\{\boldsymbol{z}_1, \cdots, \boldsymbol{z}_n, \boldsymbol{s}_1, \cdots, \boldsymbol{s}_n\}$ is a random variable. To make $S(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$ more explicit with the loss term, we assume an ML algorithm produced parameters $\boldsymbol{\theta}$ follows a posterior distribution. According to energy based models [15, 6], the posterior distribution of an ML model $\boldsymbol{\theta}$ follows:

$$p(\boldsymbol{\theta} | \boldsymbol{z}_1, \cdots, \boldsymbol{z}_n) \propto e^{-\frac{1}{\gamma} \sum_{i=1}^n \ell(\boldsymbol{\theta}, \boldsymbol{z}_i)}, \tag{3}$$

where $\gamma$ is a temperature parameter controlling the stochasticity of $\boldsymbol{\theta}$. Following this assumption, given $\{\boldsymbol{z}_1, \cdots, \boldsymbol{z}_n, \boldsymbol{s}_1, \cdots, \boldsymbol{s}_n\}$, the posterior distribution of $\boldsymbol{\theta}_k$ follows:

$$p(\boldsymbol{\theta}_k | \boldsymbol{z}_1, \cdots, \boldsymbol{z}_n, \boldsymbol{s}_1, \cdots, \boldsymbol{s}_n) \propto e^{-\frac{1}{\gamma} \sum_{i=1}^n s_{ik} \ell(\boldsymbol{\theta}_k, \boldsymbol{z}_i)}. \tag{4}$$

We further define the posterior distribution of $\boldsymbol{\theta}_k$ given training samples $\boldsymbol{z}_2, \cdots, \boldsymbol{z}_n$ and their source status $\boldsymbol{s}_2, \cdots, \boldsymbol{s}_n$ (i.e., given $\boldsymbol{\tau}$):

$$p_{\boldsymbol{\tau}}(\boldsymbol{\theta}_k) := \frac{e^{-\frac{1}{\gamma} \sum_{i=2}^n s_{ik} \ell(\boldsymbol{\theta}_k, \boldsymbol{z}_i)}}{\int_{\boldsymbol{t}} e^{-\frac{1}{\gamma} \sum_{i=2}^n s_{ik} \ell(\boldsymbol{t}, \boldsymbol{z}_i)} d\boldsymbol{t}}, \tag{5}$$

where the denominator is a constant. The following theorem explicitly demonstrates how to conduct the optimal source inference with the loss term.

**Theorem 2** *Given a local resulting model $\boldsymbol{\theta}_k$, a training record $\boldsymbol{z}_1$, the optimal SIA is given by:*

$$S(\boldsymbol{\theta}_k, \boldsymbol{z}_1) = \mathbb{E}_{\boldsymbol{\tau}} \left[ \sigma \left( g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}}) + \mu_\lambda \right) \right], \tag{6}$$

*where*

$$\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1) := -\gamma \log \left( \int_{\boldsymbol{t}} e^{-\frac{1}{\gamma} \ell(\boldsymbol{t}, \boldsymbol{z}_1)} p_{\boldsymbol{\tau}}(\boldsymbol{t}) d\boldsymbol{t} \right), \tag{7}$$

$$\ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1) := -\gamma \log \left( e^{-\frac{1}{\gamma} \ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1)} \right), \tag{8}$$

$$g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}}) := \frac{1}{\gamma} (\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1) - \ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1)). \tag{9}$$

The term $g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}})$ in Equation 9 is the gap between $\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1)$ and $\ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1))$. Since $\boldsymbol{\tau}$ is a training set that does not contain any information about $\boldsymbol{z}_1$, $p_{\boldsymbol{\tau}}$ corresponds to a posterior distribution of the parameters of an ML model that was trained without seeing $\boldsymbol{z}_1$. Note that $\ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$ is the local model $\boldsymbol{\theta}_k$'s evaluation of the loss on the training record $\boldsymbol{z}_1$. Comparing Equation 7 and Equation 8, we can easily find that $\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1)$ is the expectation of the loss $\ell(\cdot, \boldsymbol{z}_1)$ over the typical models that have not seen $\boldsymbol{z}_1$. Thus, we can interpret $g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}})$ as the difference between $\boldsymbol{\theta}_k$'s loss on $\boldsymbol{z}_1$ and other models' (trained without $\boldsymbol{z}_1$) average loss on $\boldsymbol{z}_1$.

In FL, the malicious server can implement an SIA in each communication round. The server receives the updated local models from each client and conducts the SIA to identify whether $\boldsymbol{z}_1$ belongs to the client k. Let us qualitatively analyze $S(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$ in *Theorem 2*. $S(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$ has two important terms $g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}})$ and $\mu_\lambda$, which decide the posterior probability. In FL, $\ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$ represents the local updated model $\boldsymbol{\theta}_k$'s loss on $\boldsymbol{z}_1$. $\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1)$ represents the average loss of $\boldsymbol{z}_1$ under the local models $\boldsymbol{\theta}_1, \cdots, \boldsymbol{\theta}_{k-1}, \boldsymbol{\theta}_{k+1}, \cdots, \boldsymbol{\theta}_K$ that are updated without $\boldsymbol{z}_1$. Note that $\ell(\cdot)$ is a loss function which measures the performance of a model on a data record. If $\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1) \approx \ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$, which means the client k behaves almost the same as other clients on $\boldsymbol{z}_1$, then $g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}}) \approx 0$. Since $\sigma(\mu_\lambda) = \lambda$, the posterior probability $S(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$ is equal to $\lambda$. Thus, we have no information gain on $\boldsymbol{z}_1$ beyond prior knowledge. In FL, the prior knowledge is $\mathbb{P}(s_{ik} = 1) = \lambda = \frac{1}{K}$. In this case, the source inference is equal to a *random guess*. However, if $\ell_{p_{\boldsymbol{\tau}}}(\boldsymbol{z}_1) > \ell(\boldsymbol{\theta}_k, \boldsymbol{z}_1)$, that is, the client k performs better than other clients on $\boldsymbol{z}_1$, $g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}})$ becomes positive. When $g(\boldsymbol{z}_1, \boldsymbol{\theta}, p_{\boldsymbol{\tau}}) > 0$, $\mathbb{P}(s_{1k} = 1 | \boldsymbol{\theta}_k, \boldsymbol{z}_1) > \lambda$

**Algorithm 1** FEDSIA The $K$ clients are indexed by $k$; $B$ is the local mini-batch size; $E$ is the number of local epochs; $\eta$ is the learning rate; $z_1$ is a training data.

---

1: **Server executes**
2: initialize $\theta_0$ // initialize weights
3: $m \leftarrow max(C \cdot K, 1)$
4: **for** each round $t = 1$ to $T$ **do**
5:     $S_t \leftarrow$ (random set of m clients)
6:     **for** each client $k \in S_t$ **do**
7:         $\theta_t^k \leftarrow$ **ClientUpdate**$(\theta_{t-1}^k)$
8:         Compute $\ell_k(\theta_t^k, z_1)$ // calculate local loss on $z_1$
9:     **end for**
10:    $i \leftarrow argmin(\ell_1(\boldsymbol{\theta}_1, z_1), \cdots, \ell_m(\boldsymbol{\theta}_m, z_1))$ // source
11:    $\theta_t \leftarrow \sum_k \frac{n^{(k)}}{n}\theta_t^k$ // update central model
12: **end for**
13: **ClientUpdate**$(\theta)$
14:    $\mathcal{B} \leftarrow$ (split $D_k$ into batches of size $B$)
15:    **for** each local epoch $i$ from 1 to $E$ **do**
16:       **for** batch $b \in \mathcal{B}$ **do**
17:          $\theta \leftarrow \theta - \eta\nabla\ell(b; \theta)$ // mini-batch gradient descent
18:       **end for**
19:    **end for**
20:    **return** $\theta$ // return model to central server

---

and thus we gain non-trivial source information on $z_1$. Moreover, since $\sigma(\cdot)$ is non decreasing, smaller $\ell(\boldsymbol{\theta}_k, z_1)$) indicates a higher probability that $z_1$ belonging to the client k.

We conclude that the smaller loss of client k's local model on a training record $z_1$, the higher posterior probability that $z_1$ belongs to the client k. This motivates us to design the SIA in FL such that the client whose local model has the smallest loss on a training record should own this record. Moreover, if the client's local model's behavior on its local training data is different from that of other clients, our attack will always achieve better performance than random guess. We give more empirical evidence in Section 4. Based on the conclusion above, we propose FEDSIA as described in Algorithm 1, an FL framework based on FedAvg [22] that allows an honest-but-curious server to implement SIAs without violating the FedAvg protocol.

## 4 Experiments

### 4.1 Datasets and Model Architectures

In the experiments, we evaluate SIAs on six datasets, i.e., Synthetic, Location[2], Purchase[3], CHM-NIST[4], MNIST[5], and CIFAR-10[6]. Among them, Synthetic is a synthetic i.i.d. dataset, which allows us to manipulate data heterogeneity more precisely. We follow the same generation setup as described in [18, 19]. Location, Purchase, CHMNIST, MNIST, and CIFAR-10 are realistic datasets which are widely used for evaluating privacy leakage on ML models [30, 14, 8, 37]. For MNIST and CIFAR-10, we use the training dataset and testing dataset given. For the rest of the datasets, we randomly select 80% samples as the training records and use the remaining 20% samples as the testing records.

We consider deep neural networks (DNN) as the collaborative models for the classification tasks. In particular, for MNIST, CHMNIST, CIFAR-10, we use a convolutional neural network with two 5x5 convolution layers (the first with 32 channels, the second with 64, each followed with 2x2 max pooling), two fully connected layers with 512 and 128 units and ReLu activation, and a final softmax

---

[2]https://sites.google.com/site/yangdingqi/home/foursquare-dataset

[3]https://www.kaggle.com/c/acquire-valued-shoppers-challenge/data

[4]https://www.kaggle.com/kmader/colorectal-histology-mnist

[5]http://yann.lecun.com/exdb/mnist/
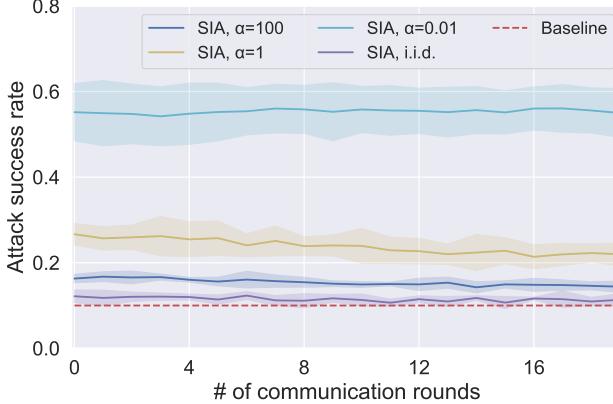
[6]https://www.cs.toronto.edu/ kriz/cifar.html

Figure 1: ASR of source inference on the Synthetic dataset in various data distribution settings. The x-axis represents the number of communication rounds. The y-axis represents ASR.

output layer. For Synthetic, Location, and Purchase, we use a fully-connected neural network with 1-hidden layer with 200 units each using ReLu activations. For each client in FL, we set a local mini-batch size of 12 for all the experiments. For all models, we use SGD with the learning rate of 0.01. Our DNN architecture does not necessarily achieve the highest classification accuracy for the considered datasets, as our goal is not to attack the best DNN architecture. Our goal is to show that SIAs can identify which local client a training record comes from when the DNN classifier is trained in a federated manner.

In our experiment, we randomly select 100 training records from each client as the target training examples of which the server wants to identify the source. We set the fraction of the clients $C$ to 1 in FL to simplify our experiments as we ignore the efficiency of the FL training when analyzing the privacy leakage. We consider *attack success rate* (ASR) as the evaluation metric for the source inference. The ASR is defined as the fraction of the target records' where the source status are correctly identified by the server. We consider a trivial attack of *randomly guessing*, which randomly selects a client as the source of the target training record as the performance baseline of an SIA. For all the learning tasks, we train the central model for 20 rounds, which is enough for the central model to converge. We record ASR during each communication round and report the highest ASR. All experiments are implemented using PyTorch with a single GPU NVIDIA Tesla P40.

## 4.2 Factors in Source Inference Attack

**Data Distribution.** The training data across clients are usually non-i.i.d. (heterogeneity) in FL. That is, a client's local data can not be regarded as samples drawn from the overall data distribution. If the training data is more heterogeneous, each local optimized model will be more different during the FL training, which benefits SIAs. Intuitively, an SIA is more effective when the degree of data heterogeneity increases. To simulate heterogeneity of training data, we follow the method used in [38, 1, 20] and use a Dirichlet distribution to divide the training records. The degree of data heterogeneity is controlled by a hyperparameter $\alpha$ ($\alpha > 0$) of the Dirichlet distribution. In general, the reverse of the magnitude of $\alpha$ reflects the degree of data heterogeneity.

**Number of Local Epochs.** In each communication round of FL, the client locally runs SGD on the current central model using its entire training dataset for several epochs and then submits the optimized model to the server. Recent studies [31, 4] have demonstrated that ML models are prone to memorize their training data. Intuitively, if a client updates the model on its local dataset with more epochs in each communication round, its local resulting model remembers the information of the local dataset better, which benefits SIAs.

## 4.3 Source Inference on Synthetic Dataset

We first conduct experiments on Synthetic to investigate how data distribution affects SIAs, because synthetic data allows us to manipulate the heterogeneity of the training data precisely. Without loss

6

(a) SIAs under various $\alpha$.
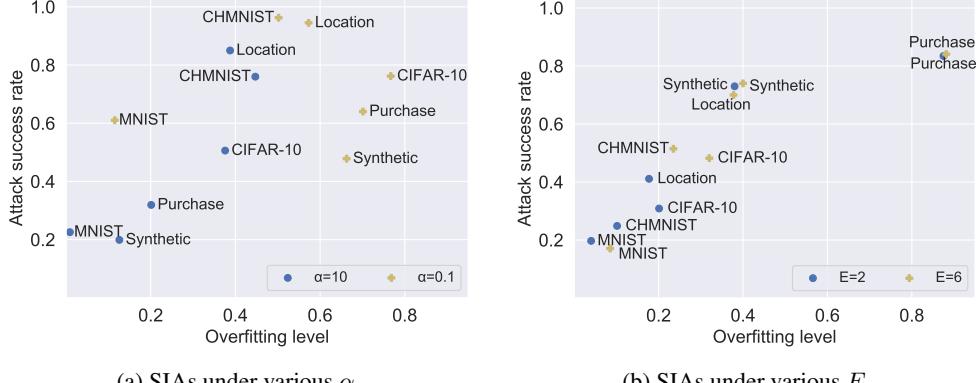
(b) SIAs under various $E$.

Figure 2: The overfitting levels affect the performance of SIAs in FL, where the x-axis represents different overfitting levels and the y-axis represents ASR. (a) We fix $E$ and $K$ for the FL model on the same dataset and only change $\alpha$ from 10 to 0.1. (b) We fix $K$ and $\alpha$ for the FL model on the same dataset and only change $E$ from 2 to 6.

of generality, we assume there are 10 clients and $E = 5$. Fig. 1 depicts the ASR of SIAs in each communication round during the FL training. We observe that our proposed SIAs always perform better than the random guessing baseline. This serves as empirical evidence for our theoretical analysis that random guess is the lower bound of our optimal source inference. The attacker performs better when the local data changes from i.i.d. to non-i.i.d., and the ASR increases as the heterogeneity of data increases.

## 4.4 Source Inference on All Dataset

We have demonstrated that SIAs are effective on synthetic data in both i.i.d. and non-i.i.d. settings. Now we use real datasets to further validate the effectiveness of SIAs and investigate the factors affecting the performance of SIAs. The SIAs leverages the local models' different prediction loss on the training examples. Intuitively, if the local model is overfitted, it will perform much better on its training members than other data, i.e., distinguishable prediction loss between the local training data and other data. We link the level of non-i.i.d. and the number of local epochs to overfitting to study how the two factors affect the performance of SIAs.

Fig. 2 shows the SIAs' ASR of different FL models under different overfitting levels. The overfitting level of the FL model here is calculated as the average of all local models' generalization gap. As we can see, increasing the level of non-i.i.d. across clients (i.e., changing $\alpha$ from 10 to 0.1) increases the ASR of SIAs in all models, as increasing the level of non-i.i.d. will inevitably increase the level of overfitting. However, when we increase the number of local epochs from 2 to 6, the ASRs of SIAs on CHMNIST, CIAFR-10, Location increases while on MNIST, Synthtic, Purchase the ASRs does not vary much. This is because changing local epochs does not increase the overfitting level of all models.

Success of SIAs is directly related to the generalizability of the local models and the diversity of the local training data. If a local model generalizes well to inputs beyond its local training members, it will not leak too much source information about its local data. Moreover, if the local training set fails to represent the overall training data distribution, the local model leaks significant information about its local data and the ASR of SIAs remains high. Recent works [42, 19, 17] have demonstrated that the non-i.i.d. of training data in FL has brought statistical heterogeneity challenges for model convergence guarantees. In this paper, we show another harm of non-i.i.d.: the leakage of source privacy for local data.

## 5 Discussion

Many works [29, 9, 23] suggest differential privacy [7] can be used against the inference attacks due to its theoretical guarantee of privacy protection. Here, we test the differential privacy as a

Table 1: Source Inference Defense via Differential Privacy.

| Dataset | FL without DP | | | FL with DP | | | |
|---------|------------------------------|-----------------------------|-------------------|------------------------------|-----------------------------|-------------------|------------|
| | Accuracy$_{train}$ | Accuracy$_{test}$ | ASR$_{SIA}$ | Accuracy$_{train}$ | Accuracy$_{test}$ | ASR$_{SIA}$ | $\epsilon$ |
| Location | 97.4% | 71.2% | 75.1% | 12.7% | 12.3% | 55.1% | 1.31 |
| CIFAR-10 | 97.7% | 68.3% | 55.2% | 10.1% | 10.0% | 24.4% | 1.48 |

defense technique against SIAs in FL. In this experiment, we evaluate the defense approaches on Location and CIFAR-10, as shown in Table 1. In the experimental setting, we set $\alpha = 10$, $K = 2$, $E = 5$ for Location, and $\alpha = 10$, $K = 5$, $E = 5$ for CIAFR-10. From the results, we can see that the ASRs drop from 75.1% to 55.1% on Location, and 55.2% to 24.4% on CIFAR-10, while applying differential privacy. However, when differential privacy can defend the SIA, it also hurts the performance of the model on its tasks, where the model utility drops from 71.2% to 12.3% on Location, and 68.3% to 10.0% on CIFAR-10. In this case, we can conclude that vanilla DP is not a effective solution for SIA in FL, which provides future research opportunities.

## 6 Related work

### 6.1 Inference Attacks in FL

Macahan et al. [22] first propose the federated learning framework that can mitigate the privacy leakage of model training with limited, unbalanced, massively, or even non-IID data among distributed devices, such as mobile phones [27], healthcare data [39]. The motivation is to share the model weights instead of the private data for better privacy protection. However, recent works [25, 26, 43, 37, 41] investigate several privacy attacks in FL, including property inference attacks [8], reconstruction attacks [12], and membership inference attacks [30, 36]. MIAs in FL allows a malicious participant or server to distinguish the training members of the trained model from the non-members. Melis et al. [25] first explore MIAs in FL and demonstrate that an adversary can infer whether a specific location profile was used to train an FL model on FourSquare location dataset with 0.99 precision and perfect recall. Nasr et al. [26] suggest an adversary can actively craft his updated model to extract more membership information about other clients. For training members of the FL model, the existing inference attacks fail to explore which client owns them. The source inference attacks proposed in this paper fill this gap.

### 6.2 Privacy Defenses in FL

To enhance privacy protection, differential privacy and other privacy protection mechanisms, *e.g.,* secure aggregation, have been recently applied to federated learning [21, 2, 9, 24, 3, 32]. Previous works mostly focus on either the centralized differential privacy mechanism that requires a central trusted party [9, 24], or local differential privacy, in which each user perturbs its updates randomly before sending it to an untrusted aggregator [35, 33]. These privacy-preserving approaches have been evaluated effectively for inference and other attacks [9, 23, 3, 16, 33] in FL. However, no protection approaches have been explored for SIAs. As discussed in the last section, applying differential privacy in FL is not an effective solution, since it suffers from the trade-off between model utility and defense performance of SIAs, providing future research opportunities.

## 7 Conclusion

In this paper, we propose a new inference attack named source inference attack in the context of FL, which enables a malicious server to infer the source of a training example between clients. We derive an optimal attack strategy formally that the malicious server is able to gain non-trivial source information of the training members by evaluating the local model's loss. We evaluate SIAs in FL with many real datasets and different settings. The extensive experimental results demonstrate the effectiveness of SIAs in practice.

# References

[1] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov. How to backdoor federated learning. In *AISTATS*. PMLR, 2020.

[2] A. Bhowmick, J. Duchi, J. Freudiger, G. Kapoor, and R. Rogers. Protection against reconstruction and its applications in private federated learning. *CoRR, arXiv:1812.00984*, 2018.

[3] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth. Practical secure aggregation for privacy-preserving machine learning. In *CCS*, 2017.

[4] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *USENIX Security*, 2019.

[5] D. Devakumar, G. Shannon, S. S. Bhopal, and I. Abubakar. Racism and discrimination in covid-19 responses. *The Lancet*, 2020.

[6] Y. Du and I. Mordatch. Implicit generation and modeling with energy-based models. In *NeurIPS*, 2019.

[7] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*. Springer, 2006.

[8] K. Ganju, Q. Wang, W. Yang, C. A. Gunter, and N. Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. In *CCS*, 2018.

[9] R. C. Geyer, T. Klein, and M. Nabi. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*, 2017.

[10] C. He, S. Li, J. So, X. Zeng, M. Zhang, H. Wang, X. Wang, P. Vepakomma, A. Singh, H. Qiu, et al. Fedml: A research library and benchmark for federated machine learning. *arXiv preprint arXiv:2007.13518*, 2020.

[11] C. He, K. Balasubramanian, E. Ceyani, C. Yang, H. Xie, L. Sun, L. He, L. Yang, S. Y. Philip, Y. Rong, et al. Fedgraphnn: A federated learning benchmark system for graph neural networks. 2021.

[12] B. Hitaj, G. Ateniese, and F. Perez-Cruz. Deep models under the gan: information leakage from collaborative deep learning. In *CCS*, 2017.

[13] H. Hu, Z. Salcic, G. Dobbie, and X. Zhang. Membership inference attacks on machine learning: A survey. *arXiv preprint arXiv:2103.07853*, 2021.

[14] B. Jayaraman and D. Evans. Evaluating differentially private machine learning in practice. In *USENIX Security*, 2019.

[15] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning. *Predicting structured data*, 2006.

[16] D. Li and J. Wang. Fedmd: Heterogenous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, 2019.

[17] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 2020.

[18] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith. Federated optimization in heterogeneous networks. In *MLSys*, 2020.

[19] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang. On the convergence of fedavg on non-iid data. In *ICLR*, 2019.

[20] T. Lin, L. Kong, S. U. Stich, and M. Jaggi. Ensemble distillation for robust model fusion in federated learning. *arXiv preprint arXiv:2006.07242*, 2020.

[21] L. Lyu, H. Yu, X. Ma, L. Sun, J. Zhao, Q. Yang, and P. S. Yu. Privacy and robustness in federated learning: Attacks and defenses. *arXiv preprint arXiv:2012.06337*, 2020.

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

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

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

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

[26] M. Nasr, R. Shokri, and A. Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *S&P*. IEEE, 2019.

[27] W. Pan and L. Sun. Global knowledge distillation in federated learning. *arXiv preprint arXiv:2107.00051*, 2021.

[28] G. D. P. Regulation. Regulation eu 2016/679 of the european parliament and of the council of 27 april 2016. *Official Journal of the European Union*, 2016.

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

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

[31] C. Song, T. Ristenpart, and V. Shmatikov. Machine learning models that remember too much. In *CCS*, 2017.

[32] L. Sun and L. Lyu. Federated model distillation with noise-free differential privacy. *arXiv preprint arXiv:2009.05537*, 2020.

[33] L. Sun, J. Qian, X. Chen, and P. S. Yu. Ldp-fl: Practical private aggregation in federated learning with local differential privacy. In *IJCAI*, 2021.

[34] S. Truex, L. Liu, M. E. Gursoy, L. Yu, and W. Wei. Towards demystifying membership inference attacks. *arXiv preprint arXiv:1807.09173*, 2018.

[35] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou. A hybrid approach to privacy-preserving federated learning. In *12th ACM Workshop on Artificial Intelligence and Security*, 2019.

[36] Y. Wang and L. Sun. Membership inference attacks on knowledge graphs. *arXiv preprint arXiv:2104.08273*, 2021.

[37] Z. Wang, M. Song, Z. Zhang, Y. Song, Q. Wang, and H. Qi. Beyond inferring class representatives: User-level privacy leakage from federated learning. In *INFOCOM*. IEEE, 2019.

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

[39] X. Xu, H. Peng, L. Sun, M. Z. A. Bhuiyan, L. Liu, and L. He. Fedmood: Federated learning on mobile health data for mood detection. *arXiv preprint arXiv:2102.09342*, 2021.

[40] Y. Xu, L. Ma, F. Yang, Y. Chen, K. Ma, J. Yang, X. Yang, Y. Chen, C. Shu, Z. Fan, et al. A collaborative online ai engine for ct-based covid-19 diagnosis. *medRxiv*, 2020.

[41] C. Yang, H. Wang, K. Zhang, L. Chen, and L. Sun. Secure deep graph generation with link differential privacy. *arXiv preprint arXiv:2005.00455*, 2020.

[42] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.

[43] L. Zhu, Z. Liu, and S. Han. Deep leakage from gradients. In *NeurIPS*, 2019.