

# MixNN: Protection of Federated Learning Against Inference Attacks by Mixing Neural Network Layers

Antoine Boutet

Univ Lyon, INSA Lyon, Inria, CITI  
antoine.boutet@insa-lyon.fr

Jan Aalmoes

Univ Lyon, INSA Lyon, Inria, CITI  
jan.aalmoes@inria.fr

Thomas Lebrun

Univ Lyon, INSA Lyon, Inria, CITI  
thomas.lebrun@inria.fr

Adrien Baud

Univ Lyon, INSA Lyon, Inria, CITI  
adrien.baud@inria.fr

## Abstract

Machine Learning (ML) has emerged as a core technology to provide learning models to perform complex tasks. Boosted by Machine Learning as a Service (MLaaS), the number of applications relying on ML capabilities is ever increasing. However, ML models are the source of different privacy violations through passive or active attacks from different entities. In this paper, we present MixNN a proxy-based privacy-preserving system for federated learning to protect the privacy of participants against a curious or malicious aggregation server trying to infer sensitive attributes. MixNN receives the model updates from participants and mixes layers between participants before sending the mixed updates to the aggregation server. This mixing strategy drastically reduces privacy without any trade-off with utility. Indeed, mixing the updates of the model has no impact on the result of the aggregation of the updates computed by the server. We experimentally evaluate MixNN and design a new attribute inference attack,  $\nabla\text{SIM}$ , exploiting the privacy vulnerability of SGD algorithm to quantify privacy leakage in different settings (i.e., the aggregation server can conduct a passive or an active attack). We show that MixNN significantly limits the attribute inference compared to a baseline using noisy gradient (well known to damage the utility) while keeping the same level of utility as classic federated learning.

**Keywords:** Machine Learning, Federated Learning, Inference Attacks

## 1 Introduction

The collection of personal data is a subject firmly grounded in public debates. The growing awareness of the population on privacy issues led to stronger regulations on data protection (e.g., GDPR, HIPAA) and contributed to the appearance of new services making privacy an incentive vector such as privacy-based search engine (e.g., Duckduckgo, Qwant), web browsing (e.g., Web Proxy, Tor, Brave), or mailing (e.g., Protonmail). These services rely on infrastructures setup by companies, nonprofit organizations promoting privacy, or are fully peer-to-peer involving devices of end-users.

However, personal and private data is still the fuel of all desires. In this context, Machine Learning (ML) has emerged as a core technology to analyze and provide learning models from large volumes of data and to perform complex tasks such as classifications, predictions or clustering. The success of ML has driven different providers to launch Machine Learning as a Service (MLaaS) engines to make ML operation easier for anyone, without the cost and time to build in-house infrastructures. These new services has led to an ever increasing number of new applications or services relying on ML capabilities in different domains such as computer vision, health analytic and speech recognition to name a few. However, it has been showed that ML models may leak information in the training data [1–3]. The fact that many applications using this technology involve the collection and processing of personal and sensitive data has raised privacy concerns [4].

Despite being popular, the memorization of training data by a ML model is the source of different privacy violations such as membership, property and attribute inference through passive or active attacks. Membership inference [5, 6] refers to the capacity of an adversary to identify if a data point (or the data of an individual) has been used to train the target model. This attack has a serious privacy implication if the model is training with sensitive information (e.g., data from people with certain health status). Property inference [7, 8], in turn, corresponds to the inference by an adversary of the properties of training data such as the features that characterize each class. This property inference can also concern a subset of the training inputs. This ability to learn from training data is desired if the inference is directly related to the main task of the model. By contrast, attribute inference [9] corresponds to the fact that an adversary is able to infer an unintended and undesired attribute not correlated to class’s characteristic feature. Root causes related to these attack surfaces as well as the link between utility (e.g., through model overfitting [10, 11]) and privacy are not well understood.

Recently, Federated Learning (FL) [12] has emerged as promising privacy-by-design alternatives to decentralized learning schemes. In such a collaborative scheme, personal data never leaves the user device. Instead, devices (computing and refining a learning model with their own data) and

a central server (aggregating models) work together to build a global learning model. This new ML scheme has attracted many attention these last years, not only from the research community but also from major Internet companies, suggesting future deployments. For instance, Google envisioned to massively exploit FL in the near future (e.g., through its FLoC API [13]). While the FL scheme is a clear step forward towards enforcing users’ privacy, it still suffers from a large ML-based attack surface including membership, property and attribute inference from participants or from the server. Different protection mechanisms to limit inference capabilities of an adversary have been proposed [14]. For instance, some solutions [15, 16] are based on perturbation in order to reveal only a noisy information to the server, such as differential privacy. However, these solutions significantly damage the accuracy of the model and its capacity to converge. Secure aggregation relying on a cryptographic scheme has been also proposed [17–19]. Similar to MixNN, this solution ensures that the server is only aware of the aggregate of all models, keeping the model of each participant (and the associated inference) private. Beside the overhead of this solution remains low, the underlying cryptographic scheme requires the participation of the server in the protection. We argue that such solutions are not deployed in practice. Indeed, few companies accept to afford the additional cost of the protection. For instance, Private Information Retrieval (PIR) protocols which follow similar cryptographic scheme to protect the profile of users are not widely adopted in practice. Moreover, a curious or malicious server trying to infer information from participants will certainly not adopt such a protection.

In this paper, we present MixNN, a new privacy-preserving service for FL against inference attacks from a curious or malicious aggregation server. To achieve that, MixNN relies on a proxy mixing the layers of the model updates among participants before sending them to the aggregate server. Like Mixnets to ensure anonymity in information routing [20], mixing the layers of the participants’ updates of neural network prevents inference attack without decreasing the accuracy of the aggregated model. This solution, albeit simple, leads to drastically improving the privacy without any trade-off with utility. In addition, MixNN is transparent to the FL service, participants only need to configure a web proxy for the associated traffic. To make the deployment of MixNN easier by anyone (e.g., operate by an individual or non profit organizations willing to protect privacy) and possibly on an untrusted infrastructure, the proxy mixing the neural network layers is running inside an SGX enclave ensuring confidentiality and attestation on its behavior.

To illustrate the capability of MixNN to protect privacy while maintaining the same level of utility, we design  $\nabla\text{SIM}$  a new attribute inference attack exploiting the privacy vulnerability of the Stochastic Gradient Descent (SGD) algorithm.

More precisely, the server infers sensitive attributes of participants from their model updates. These updates represent the gradient vectors which minimize their contribution to the model’s training loss. These gradient vectors are consequently influenced by the local data of the user. Based on auxiliary information on the model updates from participants belonging to each class of sensitive attributes,  $\nabla\text{SIM}$  uses the gradient vector as a fingerprint to infer attributes.  $\nabla\text{SIM}$  can be conducted passively only through the observation of model updates (i.e., the case of a curious and undetectable adversary server), or actively by influencing the model sent to participants to extract more information (i.e., the case of a malicious adversary server modifying the protocol).

Our work takes a quantitative and empirical approach. We implement MixNN and experimentally evaluate it with several datasets and neural networks architectures. We also leveraged  $\nabla\text{SIM}$  to quantify privacy leakage through attribute inference attack on classical federated learning scheme, a baseline using perturbation (i.e., noise) to protect the model updates (widely used in Differential Privacy), and MixNN. We show that MixNN limits the attribute inference compared to other baselines without decreasing the accuracy of the global aggregated model. Moreover, we show that the MixNN proxy introduces only a small latency on the model updates. The code of MixNN and  $\nabla\text{SIM}$  are publicly available.

The remainder of this paper is organized as follows. Section 2 presents background, Section 3 defines the problem and the threat model, Section 4 explains the design and the implementation of MixNN, Section 5 presents the new attribute inference attack, Section 6 reports the evaluation of MixNN, Section 7 reviews related work, and Section 8 concludes this paper.

## 2 Background

In this section, we review background related to Neural Networks 2.1, Federated Learning 2.2, inference attacks 2.3, Mixnet 2.4, and Intel SGX 2.5.

### 2.1 Neural Networks

An ML model is a function  $f(\theta) : X \mapsto Y$  parameterized by a set of parameters  $\theta$ , where  $X$  denotes the input (or feature) space ( $X = x_1, x_2, \dots, x_k$ ), and  $Y$  the output space ( $Y = y_1, y_2$ ). Training an ML model corresponds to find the optimal set of parameters  $\theta$  that fits the training data. This is done by optimizing an objective function (loss) which penalizes the model when it is wrong. For instance, if we consider a classification task trained through a supervised learning, parameters  $\theta$  are updated if the model misclassifies training data.

Neural networks are a family of ML models which have become popular for a variety of ML tasks. A neural network is composed of multiple layers of non-linear mappings from input to intermediate hidden states (or hidden layers) and

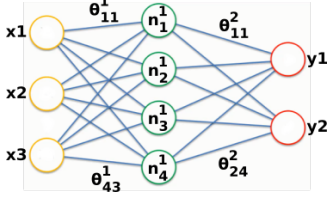


Figure 1. Example of Neural Network.

then to output where each layer transforms the output of the preceding layer to produce input for the next layer. The topology of the connections between layers and the type of considered transformation function are task-dependent and impact the accuracy of the model. For instance, convolutional layers account for locality where each neuron receives a restricted input space while a neuron receives the entire previous layer in a fully connected layer.

A neural network  $f$  is composed of a collection of  $n$  hidden layers ( $f = (l_1, l_2, \dots, l_n)$ ). Each layer  $l_i$  is composed of a set of  $m$  neurons ( $l_i = (n_1^i, n_2^i, \dots, n_m^i)$ ). For input  $X$ , the output of the neural network, can be formally written as:

$$f(\theta) = F_n(F_{n-1}(\dots F_2(F_1(X)))) ,$$

where  $X$  is the input,  $n$  the number of layers,  $F_i$  represents a transformation function of the layer  $l_i$ , and  $\theta$  is the set of floating-point weights associated with each connection between two neurons of different layers. Considering a fully connected neural network (as depicted Figure 1),  $\theta_{ab}^t$  represents the weight connecting the node  $n_a^t$  to the node  $n_b^{t-1}$ . These weights are updated during training according to the method to optimize the objective function. In this work, we consider Stochastic Gradient Descent (SGD) to do this optimization. SGD is an iterative approach where the optimizer receives a batch of training data and updates the model parameters  $\theta$  at each iteration according to both the direction of the gradient of the objective function and a learning rate  $\eta$  which scales the update. Once the gradient is close to zero, the model has converged to a local minimum and the training is finished. The model is evaluated through its accuracy over testing data points not used to train the model. The hyperparameters refer to the set of tunable parameters not related to the neural network (e.g., weights associated to connection) such as the number of training iterations, the size of the training batch or the learning rate.

## 2.2 Federated Learning

Federated Learning (FL) is a collaborative learning scheme to train an ML model [12]. In such a scheme, personal data never leaves the device of participants. Instead, devices train a ML model locally and interact with a central server to build a global learning model.

The iterative-based operating flow of classical FL is depicted Figure 2. Each iteration contains three steps. First, the

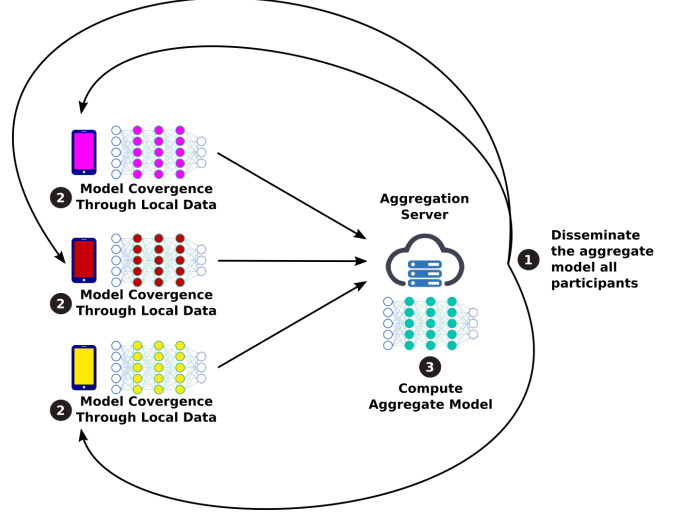


Figure 2. Operating flow of Federated Learning.

aggregation server disseminates a global model to participants (step ① in the figure). Each participant then trains and refines this model with its own data stored locally (step ②). After this local training, each participant holds its own variation of the model sent by the server. Participants then send their updated model parameters to the aggregation server. Finally, the server aggregates all these updates to generate a new global model (step ③) which will be disseminated to participants in the next iteration. Iteratively, the global model maintained on the server converges without requiring access to the personal data of participants.

## 2.3 Inference Attacks

By keeping locally data of the users on their device, FL improves privacy by design. However, FL can disclose sensitive information via model updates that are based on the training data. Indeed, any useful ML model reveals something about the population from which the training data was drawn. Indeed, a classifier model for instance may reveal the features that characterize a given class or help construct data points that belong to this class. The first privacy violation is property inference: identification of the features that characterize each class, making it possible to construct representatives of these classes through model inversion attacks [1]. Another privacy violation is attribute inference [10]: the leak of personal and unintended information (i.e., properties that hold for certain subsets of the training data, but not generically for all class members). The last privacy violation in our setting is membership inference [21]: given an exact data point, determine if it was used to train the model.

Memorization of training data by deep neural networks enables an adversary to conduct all these privacy violations. Firstly, this memorization usually combined with over-fitting of the model are exploited by an adversary to conduct a membership inference attack in order to discriminate if a

user has been part of the training or not [5]. This attack has a serious privacy implication especially if the learning model is related to sensitive information (e.g., presence of a certain pathology). Secondly, as deep-learning models come up with separate internal representations of all kinds of features, some of which are unpredictable and independent of the task being learned, the memorization of the training data can be leveraged by an adversary to infer a sensitive attribute [9]. In addition, due to the distributed nature of FL, passive and active inference attacks can be conducted by any participant or by the server.

Introducing Differential Privacy in the Stochastic Gradient Descent (DP-SGD) [14, 22] has been proposed to reduce the inference capability of an adversary, however this solution significantly damages the accuracy of the model and its capacity to converge [15, 16]. In addition, the noise calibration and the management of the privacy budget is not trivial. Other defenses propose to reduce the overfitting [23] but inherently decrease the utility.

## 2.4 Mixnets

The concept of mixing information to make them indistinguishable or unlinkable is not new. Mix networks (Mixnets) [20] uses this concept to provide a proxy-based anonymity system. This system aims to provide unlinkability between the message sent by an user, and the message received by the destination. More precisely, to prevent traffic analysis attacks, Mixnets route each message of the user through a set of anonymity servers called mixes. Mixes collect and shuffle (or mix) many messages before to route then to the destination. A variety of mixnets have been proposed including Aqua [24], Riffle [25], and Mixminion [26] addressing differently the tradeoff between anonymity, latency and bandwidth.

The limitation of these systems are similar to Tor, it is difficult for a user to determine which edge is uncompromised and powerful adversaries controlling both ends of the circuit can still deanonymize clients.

## 2.5 Intel SGX

The MixNN proxy relies on a Trusted Execution Environment (TEE), which leverages custom microprocessor zones, to enforce isolation, confidentiality and integrity of code and data. Specifically, we use Intel Software Guard Extensions (SGX) [27, 28] which defines the concept of enclave. The memory of an enclave is encrypted and cannot be directly accessed by other system software even by privileged code (e.g., the operating system or hypervisor). Enclaves can be attested to prove that the code running in the enclave is the one intended, and that it is running on a genuine Intel SGX platform. Once attested, enclaves can be provisioned with secret data by using authenticated secure channels. Moreover, enclaves can persist secret data outside the trusted zone by using a sealing mechanism. However, such protection comes with resource constraints. More precisely, only 96 MB out

of the 128 reserved for the enclave can be used by applications. Although virtual and dynamic memory support is available [29–31], it incurs significant overheads in paging (i.e., the sealing and unsealing operations used an encryption key derived from the CPU hardware).

## 3 System and adversary model

Before presenting MixNN and our attribute inference attack  $\nabla\text{SIM}$ , we describe our assumptions and the considered threat model. The operating flow of MixNN involves three premises with different level of trust, namely: (i) the client machine; (ii) the MixNN proxy; and (iii) the aggregation server.

First, we assume that the client machine is trusted. This includes the training data and all the computations performed locally. We do not consider malicious users trying to poison the model or to introduce backdoors.

Second, we assume that the MixNN proxy is running inside an Intel SGX enclave on an untrusted node. An adversary is thus not able to compromise the behavior or the data of the proxy but can monitor the node (i.e., honest but curious), possibly physically (e.g., monitoring network traffic, power consumption or memory access patterns). Consequently, an adversary can leverage side channel attacks [32] to infer information. We assume that the SGX enclave has generated a public and private key pair ( $k_{\text{pub}}$  and  $k_{\text{priv}}$ ).

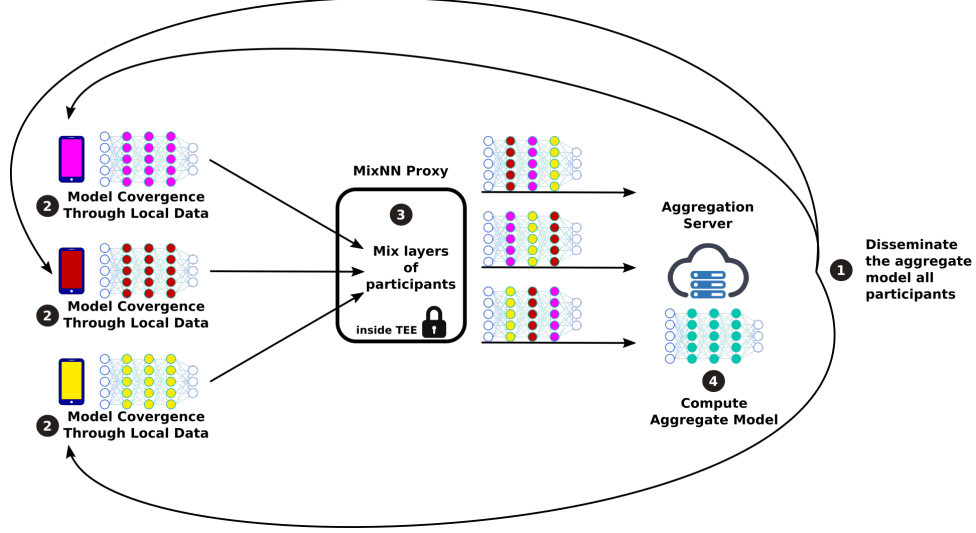
Lastly, we consider a malicious aggregation server. This server builds a model for a main classification task through a federated learning scheme but also aims to infer sensitive attributes from participants. This aggregation server can either conduct a passive or an active attack. Specifically, it can passively follow the FL operational flow to infer sensitive attributes or abuse the protocol by sending a specific model to each participant to amplify the possible inference. We do not consider any mechanism on participants to detect misbehavior of the server. In addition, we consider protected exchanges between participants and the MixNN proxy.

We consider an adversary able to collect or to use a public dataset with similar raw data (including the sensitive attribute) in order to build attack models. Each of these attack models is trained only with data from one specific class of sensitive attribute (e.g., one attack model for activity detection trained only with data from men, and another one trained only through data from women). These specialized attack models allow the adversary to compare for each participant the parameter models (i.e., the direction of the gradient obtained from the local training) to the direction of the gradient which would result in attack models. The attribute inference is computed based on a similarity metric between these directions as detailed Section 5.

## 4 MixNN Framework

In this section, we first present an overview of the MixNN (Section 4.1), the equivalence in terms of utility with a classical





**Figure 3.** MixNN introduces a proxy which receives the parameter updates from each participant, shuffle them to remove attribute footprint before to route them to the aggregation server.

FL (Section 4.2), and then present implementation details (Section 4.3).

#### 4.1 Overview

To avoid inference attacks during the learning process of a service using a federated learning scheme, MixNN operates as depicted in Figure 3. To use MixNN, users have only to configure its system to use a proxy for the associated traffic (e.g., through the configuration of its browser). As such, users seamlessly get protected without changing their habits. More precisely, compared to the classical FL pipeline (Figure 2), all parameter updates will be sent to the MixNN proxy instead of the aggregation server. To secure these updates, they are encrypted with the public key of the enclave (i.e.,  $k_{pub}$ ) to ensure that only the MixNN proxy is able to read and process them. Once loaded in the enclave, the proxy decrypts and stores the parameter updates of each layer in different lists. The proxy then pick at random one update for each layer in the associated list to generate the message containing the parameter updates to send to the aggregation server. Note that the proxy needs to initialize first each list with  $k$  updates before to send updates to the aggregation server.

The rest of the workflow remains unchanged compared to the classical one. The server aggregates the parameter updates to generate a global model which will be disseminated to all participants. Participants will then refine this model locally with their personal data before sending the parameter updates to the MixNN proxy.

The accuracy of the global model remains unchanged with or without using MixNN. Indeed, whether mixed or not, the aggregation of parameter updates of each layer are identical. In contrast, the privacy leakage through the footprint of parameter updates returned by participants is drastically

reduced. Specifically, by receiving an update mixing information from different users (breaking potential attribute footprints), the aggregation server is not able to infer any sensitive attribute. Consequently, MixNN is able to drastically improve privacy without compromising the accuracy of the system (i.e., no trade-off between utility and privacy).

#### 4.2 Utility Equivalence

By design, MixNN provides the same utility than a classical FL scheme. In this section, we prove this equivalence.

Let  $C$  be the number of participants sending their updates to the proxy. We show in this section that whether the participants use MixNN or not, the resulting aggregated model is the same. We assume that the considered MixNN proxy has enough information to send  $L$  updates to the server. Then the proxy creates a sequence  $(M_{ij})$  such that  $\forall (j_1, j_2) \in \{1, \dots, n\}^2$  with  $j_1 \neq j_2$  and  $\forall i \in \{1, \dots, L\} \quad M_{ij_1} \neq M_{ij_2}$ . And also such that  $\forall (i_1, i_2) \in \{1, \dots, L\}^2$  with  $i_1 \neq i_2$  and  $\forall j \in \{1, \dots, n\} \quad M_{i_1 j} \neq M_{i_2 j}$ .

According to previously defined notation for the nodes of a neural network, we define the  $t$ -th layer of the  $c$ -th participant of the proxy by  $(\theta_{\cdot}^t)^c$ . In the following matrix, each line is a model sent by the proxy. We remark that each combination of participant/layer appears once and only once in the matrix. This is a fundamental assumption regarding the equality of accuracy level between traditional FL and MixNN.

$$A = \begin{pmatrix} (\theta_{\cdot}^1)^{M_{11}} & (\theta_{\cdot}^2)^{M_{12}} & \dots & (\theta_{\cdot}^n)^{M_{1n}} \\ (\theta_{\cdot}^1)^{M_{21}} & (\theta_{\cdot}^2)^{M_{22}} & \dots & (\theta_{\cdot}^n)^{M_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ (\theta_{\cdot}^1)^{M_{L1}} & (\theta_{\cdot}^2)^{M_{L2}} & \dots & (\theta_{\cdot}^n)^{M_{Ln}} \end{pmatrix}$$

Now, with the regular FL procedure, the information sent by the participant is:

$$B = \begin{pmatrix} (\theta_{\cdot 1})^1 & (\theta_{\cdot 2})^1 & \dots & (\theta_{\cdot n})^1 \\ (\theta_{\cdot 1})^2 & (\theta_{\cdot 2})^2 & \dots & (\theta_{\cdot n})^2 \\ \vdots & \vdots & \ddots & \vdots \\ (\theta_{\cdot 1})^C & (\theta_{\cdot 2})^C & \dots & (\theta_{\cdot n})^C \end{pmatrix}$$

We note  $Agr : \mathcal{M}(C \times n) \rightarrow \mathcal{M}(1 \times n)$  the aggregation function which makes the mean of the columns. We show that  $Agr(A) = Agr(B)$

$$Agr(A) = \left( \frac{1}{L} \sum_{i=1}^L (\theta_{\cdot 1})^{M_{i1}}, \frac{1}{L} \sum_{i=1}^L (\theta_{\cdot 2})^{M_{i2}}, \dots, \frac{1}{L} \sum_{i=1}^L (\theta_{\cdot n})^{M_{in}} \right)$$

$$Agr(B) = \left( \frac{1}{C} \sum_{c=1}^C (\theta_{\cdot 1})^c, \frac{1}{C} \sum_{c=1}^C (\theta_{\cdot 2})^c, \dots, \frac{1}{C} \sum_{c=1}^C (\theta_{\cdot n})^c \right)$$

We make the additional assumption that  $L = C$  which means that the MixNN proxy waits for the  $C$  participants to send their updates before mixing. Which gives us that

$$Agr(A) = Agr(B) \iff$$

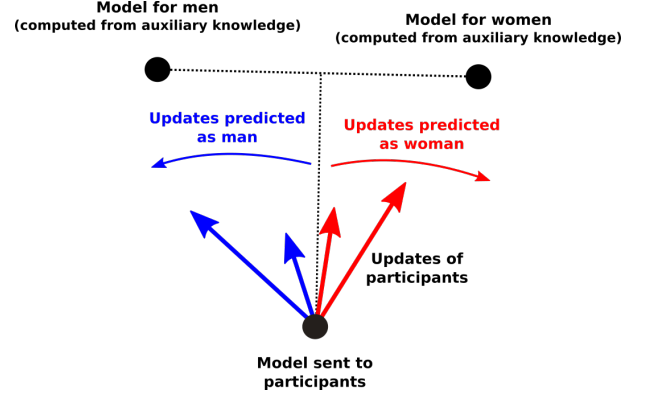
$$\left[ \forall l \in \{1, \dots, L\} \quad \sum_{c=1}^C (\theta_{\cdot l})^{M_{cl}} = \sum_{c=1}^C (\theta_{\cdot l})^c \right]$$

Which is true since our assumption on  $(M_{ij})$  gives us that  $\varphi : \{1, \dots, C\} \rightarrow \{1, \dots, C\} \quad c \mapsto M_{cl}$  is a bijective mapping.

### 4.3 Implementation

MixNN is implemented inside an Intel SGX enclave to protect its behaviors and confidentially even if it is deployed on an untrusted node. In addition, the parameter updates are encrypted by participants with the public key *pub* of the enclave. Each update received by the MixNN proxy, once decrypted, is split by layer and the parameters associated to each layer are stored in different lists. The size of these lists (noted  $k$ ) and the memory allocation according to the considered neural network models are initialized at the creation of the enclave. The  $k$  first parameter updates are used to fill out the different lists. Once these lists are full, for each further received parameter update, the MixNN proxy picks at random and removes one element in each list to build a parameter update to send to the aggregation server. The empty element in each list is then filled out with information coming from the incoming update.

To avoid side-channel attacks against SGX [32], the cost (i.e., the execution time) to process an update is constantly the same. Depending on the considered model, the size of a model can be important and not fit into the memory limit of the enclave (96MB), requiring encrypted storage outside the enclave. To avoid side-channel attack based on memory access, ORAM mechanisms (e.g., ZeroTrace [33]) can be adopted to carry out secure and oblivious access of data.



**Figure 4.**  $\nabla\text{SIM}$  infers attributes according to the gradient vector returned by participants (i.e., the parameter updates) and the learning models representative to each class of sensitive attributes (background knowledge).

The associated overhead is negligible in our context where updates are sent only periodically.

## 5 $\nabla\text{SIM}$ : A Similarity-Based Attribute Inference Attack

We design a new attribute inference attack,  $\nabla\text{SIM}$ , exploiting the privacy vulnerability of the SGD algorithm. Specifically,  $\nabla\text{SIM}$  is based on a similarity metric measuring the gradient vector returned by participants to minimize its contribution to the model’s training loss (i.e., the parameter update). Figure 4 illustrates how the inference attack works. During one round, the gradient vector returned by a participant reflects how its local data have influenced the global model disseminated by the aggregation server at the beginning of the round. This gradient vector is used as a fingerprint to infer the sensitive attribute. This fingerprint can be amplified if the attack is conducted during multiple rounds. Indeed, in this case, the influence of the user data on the evolution of the global model is more observable.

More formally,  $\nabla\text{SIM}$  measures the cosine similarity between the gradient vector returned by a participant and a reference gradient vector representative to a specific sensitive class of attribute (e.g., representative to men or women). This reference gradient vector is computed through the background knowledge of the adversary. As described in Section 3, we consider an adversary able to learn a classification model (similar to the one used for the main task) trained only with participants with a specific sensitive attribute.  $\nabla\text{SIM}$  evaluates for which of the representative models of each class of sensitive attribute the returned gradient vector is closest.

$\nabla\text{SIM}$  can be active or passive. In its passive form, the curious aggregation server follows the standard FL learning round and compares the model update returned by participants to the models representative to the sensitive attributes

computed from auxiliary knowledge. In contrast, in its active form, the aggregation server sends to participants the model calculated for being equidistant from the models associated to the sensitive attributes (e.g., the model for men and women). In this case, the server is malicious and actively modifies the protocol to conduct the inference attack.

## 6 Evaluation

We now report the results in terms of utility (Section 6.2) and privacy (Section 6.3) provided by MixNN under the considered experimental setup (Section 6.1). We also analyse the robustness of MixNN against an aggregation server trying to defeat the protection. (Section 6.4) as well as its performance from a systems perspective (Section 6.5).

Our results show that MixNN efficiently reduces the information leakage through an attribute inference attack without compromise on the accuracy of the model. We also show that MixNN introduces a negligible end-to-end latency.

### 6.1 Experiment Setup

In this section we presents the experimental setup used to evaluate MixNN, which includes datasets, metrics, the comparative baseline we compared against, and the considered methodology.

**6.1.1 Dataset.** We used two image recognition benchmark datasets (CIFAR10 and LFW) and two motion datasets for activity recognition (MotionSense and MobiAct) to assess MixNN.

**CIFAR10** is a major image classification benchmarking dataset where the data records are composed of 60,000  $32 \times 32$  RGB images where each record is mapped to one of 10 classes of common objects such as airplane, bird, cat, dog. There are 50,000 training images and 10,000 test images. The main task is the classification of the images. We artificially define 20 participants split into three groups with different preferences. We define 3 types of preference which corresponds to specific and non overlapping categories of images. The dataset is slightly balanced, two groups gather 6 participants and the last one gathers 8 participants. The profile of the participant is composed of 80% of images corresponding to its preferred classes, and the remaining 20% is composed of random images from other classes. The sensitive attribute is the preferences of the user.

**MotionSense** [34] contains data captured from an accelerometer (i.e., acceleration and gravity) and gyroscope at a constant frequency of 50Hz collected with an iPhone 6s kept in the front pocket. Overall, a total of 24 participants have performed six activities (i.e., going downstairs, going upstairs, walking, jogging, sitting and standing) during 15 trials in the same environment and conditions. The main classification task is the activity detection and the sensitive attribute is the gender of the users.

**MobiAct** [35] records the motion data from 58 subjects during more than 2500 trials, all captured with a Samsung Galaxy S3 in the front pocket. This dataset includes signals recorded from the accelerometer and gyroscope at 20Hz. We only used the trials corresponding to the same activities as MotionSense in order to do the evaluation with the same settings. Similar to MotionSense dataset, the main classification task is the activity detection and the sensitive attribute is the gender of the users.

**Labeled Faces in the Wild (LFW)** [36] contains face images for face recognition with 13,233 total samples with images for 5,749 people. The dataset additionally has attributes such as age, race, gender, smile, facial hair, glasses etc. The main classification task is smile detection and the sensitive attribute is the gender of the users.

For CIFAR10, MotionSense and MobiAct datasets, we use a neural network composed of two convolutional layers and three fully connected layers for the classification task. For LFW, in turn, we use a more complex architecture provided by Facebook, named Deep Face [37]. This neural network is composed of multiple convolutional, locally connected, maxpooling, and fully connected layers.

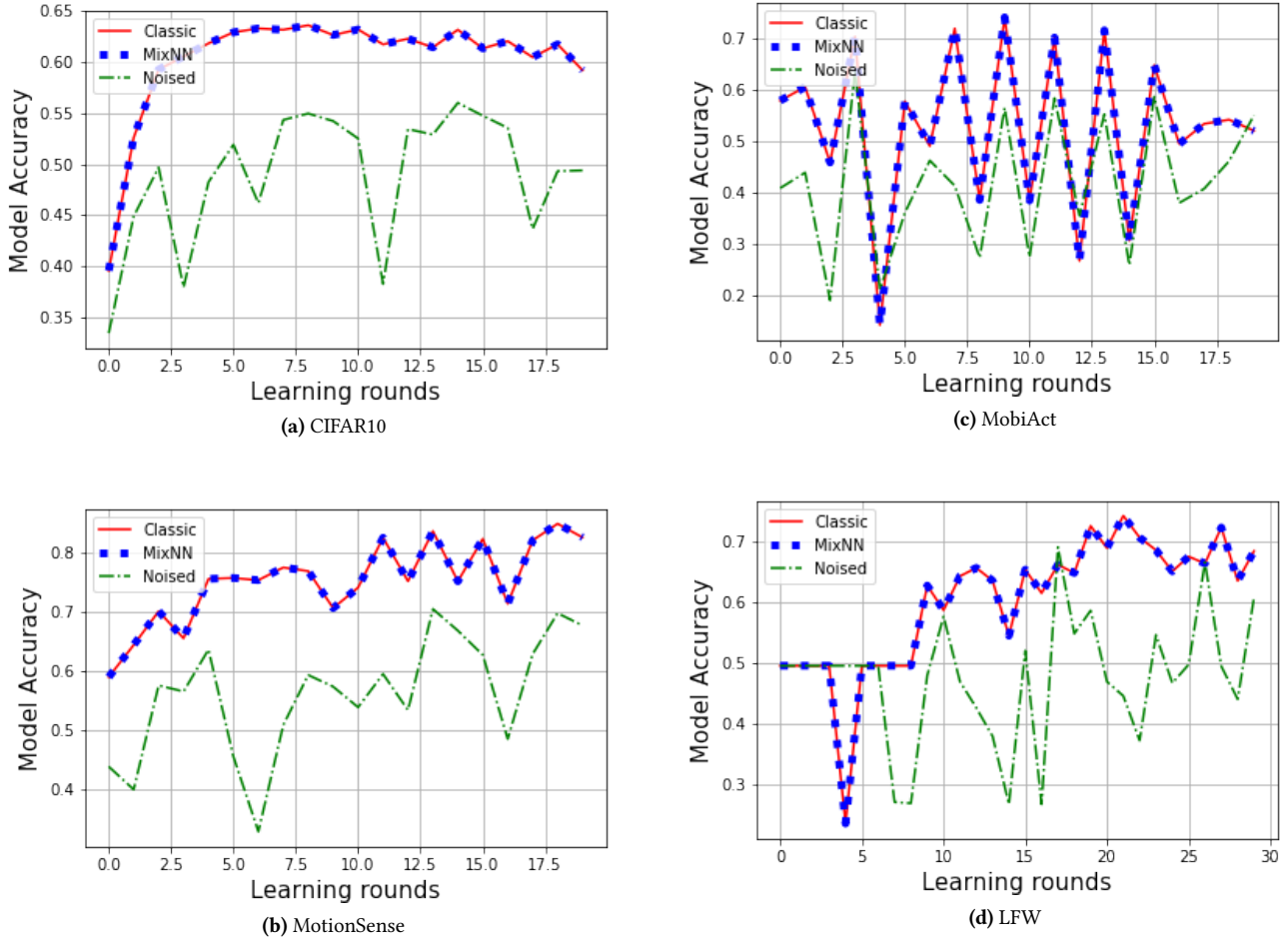
**6.1.2 Evaluation Metrics.** We evaluate MixNN through three complementary dimensions: utility, privacy and system performance.

To evaluate the utility of the target model, we consider the classification accuracy for the main task (e.g., the activity detection), noted *Model Accuracy*, measuring the ratio of number of correct predictions to the total number of predictions made.

We use  $\nabla\text{SIM}$  to conduct the inference attack and we use the classification accuracy of the sensitive attribute to estimate the success of the attribute inference, noted *Inference Accuracy*. The value of this metric indicates a data leakage according to the number of classes and if the dataset is balanced over all classes. For instance, with a balanced dataset over the gender, an accuracy above 50% indicates a data leakage through attribute inference attack. This indicates that the adversary is able to identify the gender of a participant with an accuracy higher than random guess.

To evaluate the behavior of MixNN from a systems perspective, we consider the end-to-end latency which is the time spent by the proxy to route the parameter updates to the aggregation server.

**6.1.3 Baselines.** We compare the utility and privacy provided by MixNN against a comparative approach using noisy gradient widely used in Differential Privacy studies [14, 38]. We consider an implementation based on an introduction of Gaussian noise to the updates computed through a classical local training such as using in local differential privacy [14].



**Figure 5.** MixNN provides the same utility than a standard FL scheme, noisy gradient however decreases significantly the utility and slows down the convergence.

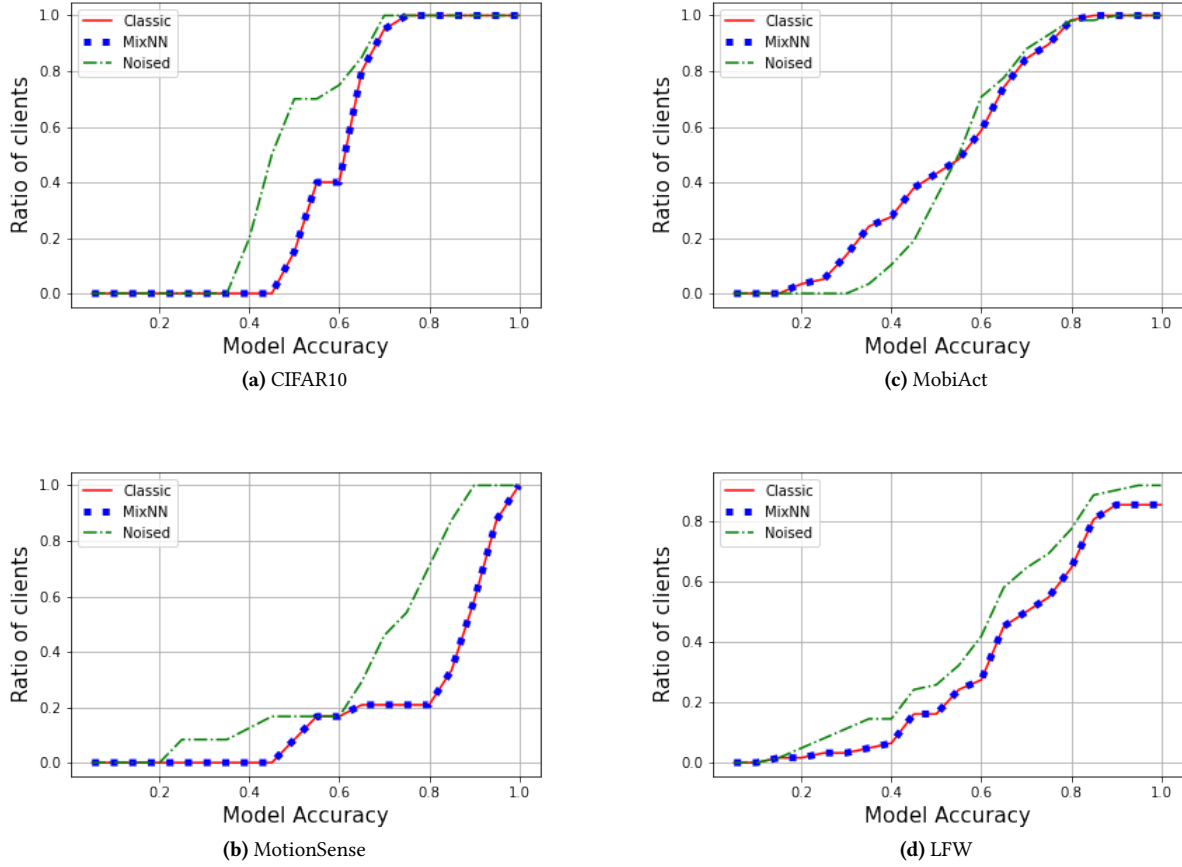
**6.1.4 Methodology.** The dataset is split between training and testing, with 5/6 of trials used for training and validation and 1/6 for testing. For CIFAR10, the federated learning model is trained on 3 local epochs for a size of data batch of 32 samples on each learning rounds, the server aggregates 16 users on each of the 10 learning rounds. For MotionSense (and MobiAct), the training is (respectively) done on 2 (and 3) local epochs for batches of 256 (and 64) samples for each of the 20 learning rounds, and the server aggregates 20 users for MotionSense and 40 users for MobiAct. For LFW, the training is done on 2 local epochs for batches of 16 samples for each of the 30 learning rounds, and the server aggregates 20 users. For every datasets, we use the "Adam" optimizer proposed by Tensorflow. We use 5-fold cross-validation in which the testing set is randomly generated from 1/5 of the users. Reported results correspond to average over 5 repetitions of each experiment. The experiments have been computed on a Laptop DELL, intel core i7 (i7-6600U) and 4G of RAM,

using TensorFlow version 2.4. The noise introduced consists on adding a Gaussian noise  $\mathcal{N}(0, 1)$  on each scalars of the neural network weights. The attack models are trained for 5 learning rounds of the previous architecture and use 4/5 users as background knowledge.

## 6.2 No compromise with utility

In this section, we evaluate the capacity of MixNN to protect privacy without compromising the utility. We compare the accuracy performance for the main classification task provided by MixNN against a classic FL scheme (i.e., without MixNN proxy) and a baseline using noisy gradient such as using in local differential privacy. Figure 5 reports the accuracy according to the learning round for all datasets. First, the results show that the same level of accuracy is provided by a standard FL scheme and MixNN with an accuracy growing according to the learning rounds. This result is expected





**Figure 6.** Using noisy gradient decreases the utility for all participants.

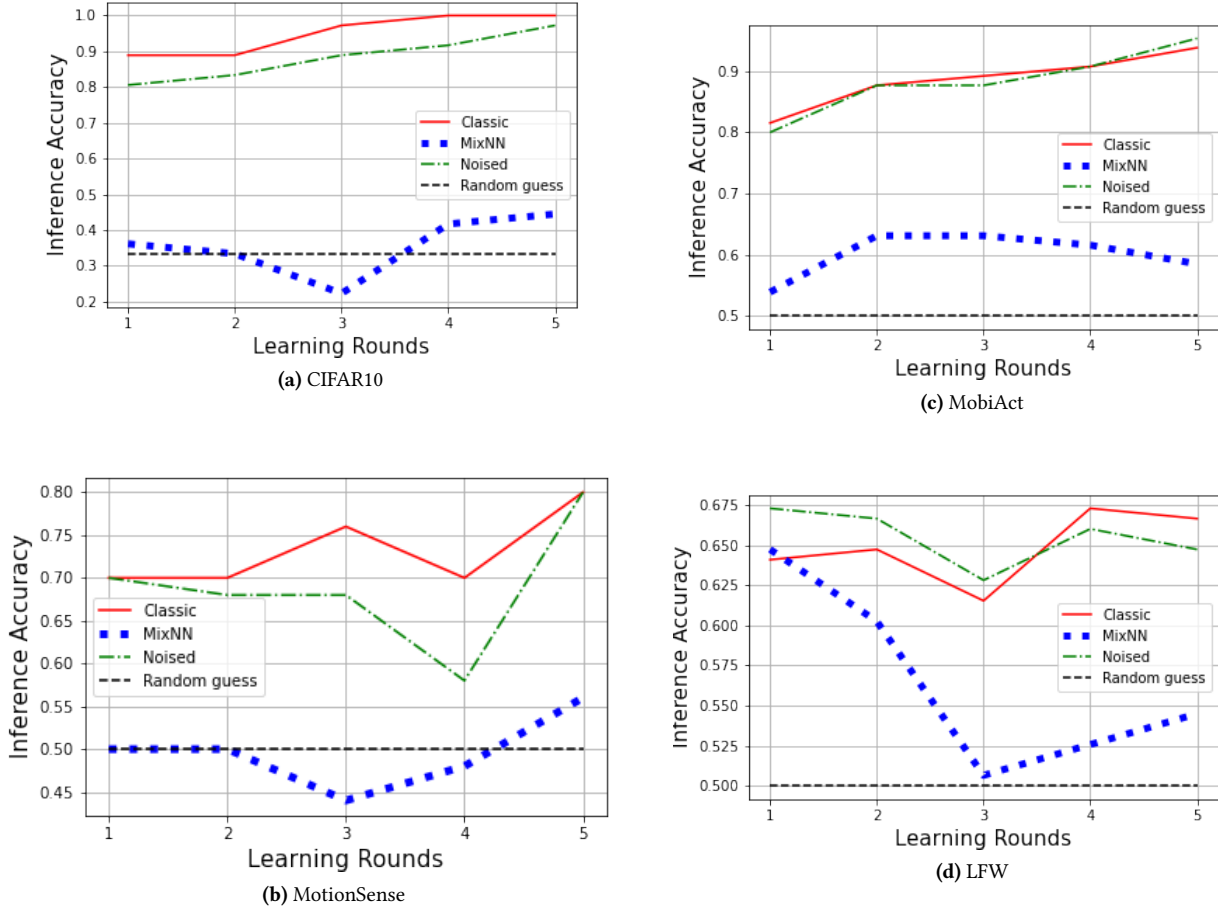
due to the aggregation equivalence of both approaches. Second, the results show that noisy gradient provides 10% lower accuracy on average and slows down the convergence. Figure 6 reports the cumulative distribution of the accuracy over the population of participants at the learning round 6. Results show that most of the participants have an accuracy with noisy gradient smaller than MixNN for all datasets (on average 0.56 for noisy gradient against 0.68 for MixNN).

### 6.3 Prevent information leakage

In this section, we evaluate the privacy leakage through the sensitive attribute attack  $\nabla_{\text{SIM}}$  through an active attack (i.e., which represents the worst case where the aggregation server is malicious and sends a calibrated model to participants to amplify the inference). Figure 7 reports for all datasets the accuracy of the inference for MixNN, Classical FL and noisy gradient according to a growing number of learning rounds. First, results show that without any protection, the server can infer the sensitive attribute of participant with a quasi perfect accuracy (100% of accuracy after 4 rounds for CIFAR10, and around 80%, 94%, and 66% of accuracy after 5 rounds for MotionSense, MobiAct, and LFW

dataset, respectively). This means that the gradient vector returned by participants and exploited by  $\nabla_{\text{SIM}}$  provides an efficient footprint to infer attributes. Second, results show that MixNN successfully limits the privacy leakage with an inference accuracy close to a random guess (0.33 for CIFAR10, and around 0.5 for MotionSense, MobiAct and LFW). The precision of the inference tends to increase with the number of local learning round but becomes stable once the model is converged. Finally, results show that noisy gradient leaks less information than a standard FL scheme but much more information than MixNN for all datasets (on average around 65% more inference accuracy).

As described Section 5,  $\nabla_{\text{SIM}}$  leverages background knowledge to build a model representative to men and a model representative to women. The attack then measures the distance between these representative models and the model updates returned by participants, the closest distance indicating a predominance towards a sensitive attribute. We evaluate now the impact of this background knowledge on the accuracy of the inference. Figure 8 depicts the inference accuracy according to a growing ratio of data used as background



**Figure 7.** MixNN better prevents attribute leakage compared to using noisy gradient.

knowledge to build the representatives models of the sensitive attributes (the quality of a model is usually correlated to the volume of data used in learning, the larger, the better). As expected, a model built with more background knowledge is more representative to individuals with a specific sensitive attribute and consequently improves the accuracy of the inference for both a classic FL and a solution using noisy gradient. However, results show that MixNN successfully protects participants against inference attack regardless the quantity of background knowledge used by the aggregation server.

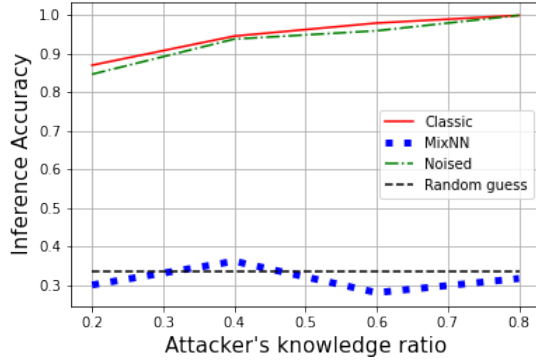
#### 6.4 Robustness of the protection

MixNN shuffles model updates sent by participants. A malicious aggregation server could then try to break the protection by enumerating through possible combinations of the shuffled updates to "reconstruct" back the original update for instance. Figure 9 reports for all datasets the cumulative distribution over all participants of the number of neighbors who have a gradient very close (i.e., in a radius of 0.5 using

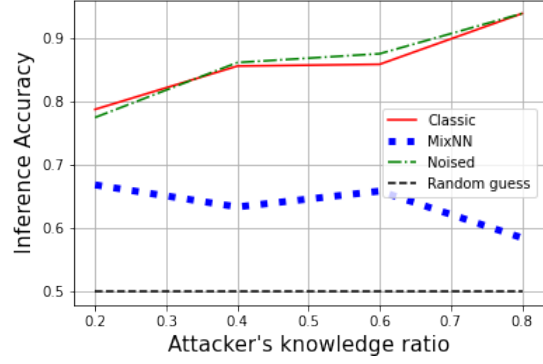
euclidean distance). All participants have at least a few other alter egos with very close gradients making it difficult for a malicious aggregation server to retrieve and distinguish all pieces of the gradient (i.e., the layers of the neural network) coming from the same participant once mixed by MixNN.

#### 6.5 System performance

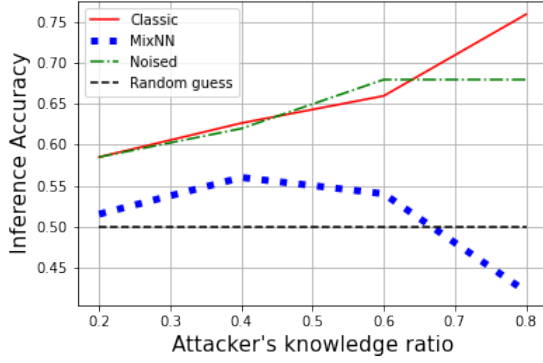
We implemented a MixNN proxy inside an SGX enclave to evaluate its system performance. To do that, the enclave receives the update parameters in an encrypted binary format and then decrypts and stores each layer of the updates in the trusted memory of the enclave. We ran the experiment using the model designed for CIFAR10 (i.e., two convolutional layers and three fully connected layers). Each update consumes 26.9MB inside the enclave and is processed in 0.19s (0.17s for the decryption and 0.02s for the storage). The mixing operation spends only 0.03s on average in our experiment. This processing time and the memory consumption are dependent on the size of the model. Using a model with three convolutional layers and three fully connected layers slightly



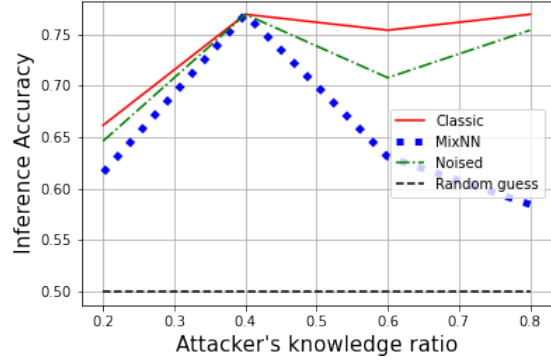
(a) CIFAR10



(c) MobiAct

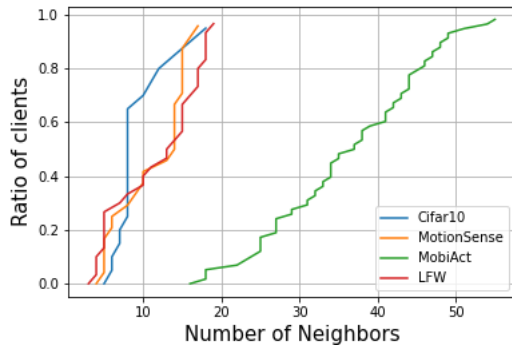


(b) MotionSense



(d) LFW

**Figure 8.** With more background knowledge, the malicious aggregation server infers more information with a standard FL scheme or using noisy gradient. This background knowledge has only a small impact on the protection of MixNN.



**Figure 9.** Many participants have very close model updates making it difficult for a malicious aggregation server to retrieve and distinguish all pieces of the gradient coming from the same participant once mixed by MixNN.

increases this time to 0.22s and 51.3MB. However, the constant processing time over all updates for a given model

(i.e., the MixNN proxy waits to receive a constant number of model updates before to mix them) reduces the surface for side channel attacks. As the learning round of classical FL scheme is usually not conducted at high rate (e.g., only when the device is plugged in and has a WiFi connection), this short delay introduced by MixNN is negligible.

## 7 Related work

The incentive behind using FL is to collectively build a learning model with better accuracy than if each user trained a model with their own data. The goal is to improve the accuracy as much as possible but several dimensions have an influence. The standard FL scheme [39] learns one global model and replicates it locally on every client. However heterogeneity of data across user devices can severely degrade performance of standard federated averaging for ML learning applications, especially for atypical users. Indeed, one unique model cannot cope with the heterogeneity of data and provide the best utility for all users [40]. To address this data heterogeneity, several approaches have been proposed such as local adaptation [16, 41] and clustering [42]. Specifically,

the clustering mechanism proposed in [42] also leverages a similarity metrics between the model updates sent back by participants (similar to MixNN) to cluster the population.

[5] designs passive and active inference algorithms for federated learning. However, this work only targets membership inference. In addition, while this attack also exploits the privacy vulnerabilities of the SGD algorithm, authors used a neural network to classify if a participant is a member of the training data or not a member. [38] also exploits the gradient exchange to infer private training data of participants. To do that, authors iteratively optimize "dummy" inputs and labels to minimize the distance between dummy gradients and real gradients. Once the optimization finished, the dummy data is close to the private training data. [43], in turn, investigates the guarantees of differentially private SGD but via data poisoning attacks.

Running MixNN in an Intel SGX enclave improves trust and confidentiality through an isolated execution environment. However, this TEE is still vulnerable to side channel attacks [32, 44]. The most common countermeasure is to use data oblivious algorithms. The objective of this technique is to eliminate the link between the nature of data inputs and the execution of the program (e.g., through the execution time or memory footprints). To achieve that, the obfuscation technique consists to hide potential patterns by making them all uniform regardless of the considered data. To reduce its inherent cost, the considered data oblivious algorithm needs to be chosen carefully according to the application [45].

## 8 Conclusion

We presented MixNN, a proxy-based privacy-preserving framework to prevent attribute inference attacks conducted from a curious or malicious aggregation server exploiting the model updates. MixNN breaks the attribute footprint leaked in the model updates by mixing layers between multiple participants. As this mixing strategy does not impact the result of the model aggregation performed by the server, the privacy improvement of MixNN does not compromise the utility of the model learned collaboratively. We also designed  $\nabla$ SIM, a new attribute inference attack which exploits the privacy vulnerability of the SGD algorithm. This attack can be passive or active. We experimentally evaluated MixNN with a standard image recognition bench-mark dataset and compared it against a state-of-the-art baseline using local differential privacy. Results show MixNN provides the same model accuracy than a classical FL scheme (i.e., the same utility) while providing a better protection against attribute inference attack (i.e., a better privacy).

## References

- [1] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *CCS*, page 1322–1333, 2015.
- [2] Liwei Song and Prateek Mittal. Systematic evaluation of privacy risks of machine learning models. *arXiv:2003.10595*, 2020.
- [3] Benjamin Zi Hao Zhao, Aviral Agrawal, Catisha Coburn, Hassan Jameel Asghar, Raghav Bhaskar, Mohamed Ali Kaafar, Darren Webb, and Peter Dickinson. On the (in)feasibility of attribute inference attacks on machine learning models. *arXiv:2103.07101*, 2021.
- [4] Emiliano De Cristofaro. An overview of privacy in machine learning. *arXiv:2005.08679*, 2020.
- [5] Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *S&P*, 2019.
- [6] Bargav Jayaraman, Lingxiao Wang, David Evans, and Quanquan Gu. Revisiting membership inference under realistic assumptions. *arXiv:2005.10881*, 2020.
- [7] Karan Ganju, Qi Wang, Wei Yang, Carl A. Gunter, and Nikita Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. In *CCS*, page 619–633, 2018.
- [8] Wanrong Zhang, Shruti Tople, and Olga Ohrimenko. Dataset-level attribute leakage in collaborative learning. *arXiv:2006.07267*, 2020.
- [9] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting unintended feature leakage in collaborative learning. *arXiv:1805.04049*, 2018.
- [10] Congzheng Song and Vitaly Shmatikov. Overlearning reveals sensitive attributes. *arXiv:1905.11742*, 2020.
- [11] Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In *CSF*, pages 268–282, 2018.
- [12] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konecny, Stefano Mazzocchi, H Brendan McMahan, et al. Towards federated learning at scale: System design. *MLSys*, 2019.
- [13] Google Research and Ads. Evaluation of cohort algorithms for the flocc api, 2020.
- [14] Mohammad Naseri, Jamie Hayes, and Emiliano De Cristofaro. Toward robustness and privacy in federated learning: Experimenting with local and central differential privacy. *arXiv:2009.03561*, 2021.
- [15] Laurens van der Maaten and Awni Hannun. The trade-offs of private prediction. *arXiv:2007.05089*, 2020.
- [16] Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging federated learning by local adaptation. *arXiv:2002.04758*, 2020.
- [17] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for federated learning on user-held data. *arXiv:1611.04482*, 2016.
- [18] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-preserving machine learning. In *CCS*, page 1175–1191, 2017.
- [19] James Henry Bell, Kallista A. Bonawitz, Adrià Gascón, Tancrede Lepoint, and Mariana Raykova. Secure single-server aggregation with (poly)logarithmic overhead. In *CCS*, page 1253–1269, 2020.
- [20] David L. Chaum. Untraceable electronic mail, return addresses, and digital pseudonyms. *Commun.*, 24(2):84–90, February 1981.
- [21] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. *arXiv:1610.05820*, 2017.
- [22] *Deep Learning with Differential Privacy*, 2016.
- [23] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. ML-leaks: Model and data independent membership inference attacks and defenses on machine learning models. *arXiv:1806.01246*, 2018.
- [24] Stevens Le Blond, David Choffnes, Wenxuan Zhou, Peter Druschel, Hitesh Ballani, and Paul Francis. Towards efficient traffic-analysis resistant anonymity networks. *SIGCOMM Comput. Commun. Rev.*,



- 43(4):303–314, August 2013.
- [25] Albert Kwon, David Lazar, Srinivas Devadas, and Bryan Ford. Riffle: An efficient communication system with strong anonymity. *Proc. Priv. Enhancing Technol.*, 2016(2):115–134, 2016.
  - [26] G. Danezis, R. Dingleline, and N. Mathewson. Mixminion: design of a type iii anonymous remailer protocol. In *2003 Symposium on Security and Privacy, 2003.*, pages 2–15, 2003.
  - [27] Victor Costan and Srinivas Devadas. Intel sgx explained. *IACR Cryptology ePrint Archive*, 2016(086):1–118, 2016.
  - [28] Ittai Anati, Shay Gueron, Simon Johnson, and Vincent Scarlata. Innovative technology for cpu based attestation and sealing. In *HASP*, 2013.
  - [29] Feng Chen, Chenghong Wang, Wenrui Dai, Xiaoqian Jiang, Noman Mohammed, Md Momin Al Aziz, Md Nazmus Sadat, Cenk Sahinalp, Kristin Lauter, and Shuang Wang. Presage: Privacy-preserving genetic testing via software guard extension. *BMC medical genomics*, 10(2):48, 2017.
  - [30] Feng Chen, Shuang Wang, Xiaoqian Jiang, Sijie Ding, Yao Lu, Jihoon Kim, S Cenk Sahinalp, Chisato Shimizu, Jane C Burns, Victoria J Wright, et al. Princess: Privacy-protecting rare disease international network collaboration via encryption through software guard extensions. *Bioinformatics*, 33(6):871–878, 2016.
  - [31] Frank McKeen, Ilya Alexandrovich, Ittai Anati, Dror Caspi, Simon Johnson, Rebekah Leslie-Hurd, and Carlos Rozas. Intel® software guard extensions (intel® sgx) support for dynamic memory management inside an enclave. In *HASP*, 2016.
  - [32] Ferdinand Brasser, Urs Müller, Alexandra Dmitrienko, Kari Kostiaainen, Srdjan Capkun, and Ahmad-Reza Sadeghi. Software grand exposure: SGX cache attacks are practical. In *WOOT*, 2017.
  - [33] Sajin Sasy, Sergey Gorbunov, and Christopher W. Fletcher. Zerotracer: Oblivious memory primitives from intel SGX. In *NDSS*, 2018.
  - [34] J. L. Reyes-Ortiz. *Smartphone-based human activity recognition*. Springer, 2015.
  - [35] George Vavoulas, Charikleia Chatzaki, Thodoris Malliotakis, M. Pediaditis, and M. Tsiknakis. The mobiaact dataset: Recognition of activities of daily living using smartphones. *ICT4AWE*, 2016.
  - [36] Gary B. Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst, October 2007.
  - [37] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deep-face: Closing the gap to human-level performance in face verification. In *CVPR*, 2014.
  - [38] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. *arXiv:1906.08935*, 2019.
  - [39] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. *arXiv:1602.05629*, 2017.
  - [40] Antoine Boutet, Carole Frindel, Sébastien Gambs, Théo Jourdan, and Claude Rosin Ngueveu. Dysan: Dynamically sanitizing motion sensor data against sensitive inferences through adversarial networks. In *AsiaCCS*, 2021.
  - [41] Manoj Ghuhane Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. *arXiv:1912.00818*, 2019.
  - [42] Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-agnostic distributed multi-task optimization under privacy constraints. *arXiv:1910.01991*, 2019.
  - [43] Matthew Jagielski, Jonathan Ullman, and Alina Oprea. Auditing differentially private machine learning: How private is private sgd? *arXiv:2006.07709*, 2020.
  - [44] Stephan van Schaik, Andrew Kwong, Daniel Genkin, and Yuval Yarom. SGAXe: How SGX fails in practice. <https://sgaxeattack.com/>, 2020.
  - [45] A K M Mubashwir Alam, Sagar Sharma, and Keke Chen. Sgx-mr: Regulating dataflows for protecting access patterns of data-intensive sgx applications. *arXiv:2009.03518*, 2020.