



# Compression Boosts Differentially Private Federated Learning

Raouf Kerkouche, Gergely Ács, Claude Castelluccia, Pierre Genevès

## ► To cite this version:

Raouf Kerkouche, Gergely Ács, Claude Castelluccia, Pierre Genevès. Compression Boosts Differentially Private Federated Learning. EuroS&P 2021 - 6th IEEE European Symposium on Security and Privacy, Sep 2021, Vienna, Austria. pp.1-15. hal-03066941

HAL Id: hal-03066941

<https://hal.archives-ouvertes.fr/hal-03066941>

Submitted on 15 Dec 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Compression Boosts Differentially Private Federated Learning

Raouf Kerkouche  
*Privatics team*  
Univ. Grenoble Alpes, Inria  
38000 Grenoble, France  
raouf.kerkouche@inria.fr

Gergely Ács  
*Crysys Lab*  
BME-HIT  
Budapest, Hungary  
acs@crysys.hu

Claude Castelluccia  
*Privatics team*  
Univ. Grenoble Alpes, Inria  
38000 Grenoble, France  
claude.castelluccia@inria.fr

Pierre Genevès  
*Tyrex team*  
Univ. Grenoble Alpes, CNRS,  
Inria, Grenoble INP, LIG  
38000 Grenoble, France  
pierre.geneves@cnrs.fr

**Abstract**—Federated Learning allows distributed entities to train a common model collaboratively without sharing their own data. Although it prevents data collection and aggregation by exchanging only parameter updates, it remains vulnerable to various inference and reconstruction attacks where a malicious entity can learn private information about the participants’ training data from the captured gradients. Differential Privacy is used to obtain theoretically sound privacy guarantees against such inference attacks by noising the exchanged update vectors. However, the added noise is proportional to the model size which can be very large with modern neural networks. This can result in poor model quality. In this paper, compressive sensing is used to reduce the model size and hence increase model quality without sacrificing privacy. We show experimentally, using 2 datasets, that our privacy-preserving proposal can reduce the communication costs by up to 95% with only a negligible performance penalty compared to traditional non-private federated learning schemes.

**Index Terms**—Federated Learning, Compressive Sensing, Differential Privacy, Compression, Denoising, Bandwidth Efficiency, Scalability.

## 1. Introduction

Traditional training of machine learning models usually requires the centralization of the user-held data. This limitation is considerably penalizing especially when the data is sensitive such as medical data. To deal with this problem, federated learning protocols have been proposed [1], [2] to collaboratively train a common model without sharing any private training data held by individual parties. In federated learning, each entity trains a common model using its own training data, and share only the gradients (i.e., model update) with each other through a central server. The server updates the common model with the shared gradients, and re-distributes the updated model to the clients for further training. This process repeats until the convergence of the common model.

However, sharing gradients computed by individual parties can leak information about their private training

data. Several recent attacks have demonstrated that a sufficiently skilled adversary, who can capture the model updates (gradients) sent by individual parties, can infer whether a specific record [3] or a group property [4] is present in the dataset of a specific party. Moreover, complete training samples can also be reconstructed purely from the captured gradients [5].

Differential privacy (DP) [6] has become a de facto privacy model which provides a formal privacy guarantee for any participant. It guarantees that the common model is roughly independent of any single client’s training data, and depends only on the characteristics that are shared among multiple parties’ training data<sup>1</sup>. DP can be achieved by adding Gaussian noise to the shared model updates. In addition, secure aggregation protocols [7] allow parties to add noise to the model update in a distributed manner, which increases robustness against byzantine attacks and requires less noise than other decentralized perturbation approaches such as randomized response [8] used in local differential privacy [9].

However, the norm of the added noise is proportional to the model size (i.e., the number of model parameters or weights). Indeed, the noise is added to every coordinate value of the gradient (update) vector including those which have very small magnitude and would anyway not improve convergence. In other words, adding noise to sparse model updates can slow down convergence significantly, or result in poor model quality [10].

In this paper, we propose to first lossily compress the gradients using compressive sensing [11]–[13] and then add noise to the compressed gradient vector. The noisy compressed vectors are then transferred to the server for aggregation. This approach has several benefits. First, compressed gradients are less sparse and also shorter than the original gradient vectors. This allows to add less noise to the compressed gradients, which eventually yields faster convergence with more accurate models than the uncompressed noisy gradients. Second, compressive sensing is linear, which means that the sum of compressed gradients

1. If client-level DP is considered. See Section 3.2.1 for more clarification.

equals the compressed sum of the gradients. Therefore, compressive sensing can be smoothly integrated with secure aggregation; the server can only access the aggregated compressed vectors which is identical to the compressed aggregation. Finally, by decreasing the size of the model updates, communication costs are reduced and bandwidth is saved. This is crucial with resource constrained parties training large models which is not uncommon nowadays.

The main contributions of this paper are summarized as follows:

- We use a slightly modified version of compressive sensing to compress sparse model updates in federated learning. Our protocol, called FL-CS allows to save bandwidth and reduce communication costs by transferring only the low frequency components of the gradient vector to the server (instead of some random frequency components like in traditional compressive sensing). The server can reconstruct the approximated sparse gradient vector by efficiently solving a convex quadratic optimization problem. This approach provides more accurate reconstruction than simply applying the inverse Fourier transform on the low frequency components. Our approach is scalable to large gradient vectors and is almost as accurate as the vanilla federated learning protocol, referred to FL-STD, without any compression, still incurring much smaller communication cost.
- We propose a privacy-preserving extension of FL-CS, called FL-CS-DP, by adding Gaussian noise to the compressed gradients. In FL-CS-DP, participants inject Gaussian noise in a distributed manner so that the sum of the noisy compressed vectors is differentially private. In addition, secure aggregation guarantees that the server (or any other third party) can only learn the noisy compressed aggregate owing to the linear compression scheme. Reconstructing the approximated gradients is an instance of Basis Pursuit Denoising (or LASSO), which can be solved with efficient solvers that provide large accuracy despite the added Gaussian noise. We show that FL-CS-DP produces more accurate models than FL-STD-DP, that is, the differentially private variant of the vanilla federated learning protocol without any compression. Therefore, compression boosts the accuracy of differentially private federated learning and also reduces bandwidth cost by more than 60% with early stopping [14].
- We evaluate our proposals on real datasets, a private medical dataset of 1.2 millions of US hospital patients and the public Fashion-MNIST dataset. We show that FL-CS-DP reduces its bandwidth cost with more than 60% compared to FL-STD-DP, meanwhile suffering negligible performance loss compared to uncompressed federated learning without any privacy guarantee (FL-STD).

## 2. Background

### 2.1. Federated Learning (FL-STD)

In federated learning [1], [2], multiple parties (clients) build a common machine learning model from union of their training data without sharing them with each other. At each round of the training, a selected set of clients retrieve the global model from the parameter server, update the global model based on their own training data, and send back their updated model to the server. The server aggregates the updated models of all clients to obtain a global model that is re-distributed to some selected parties in the next round.

In particular, a subset  $\mathbb{K}$  of all  $N$  clients are randomly selected at each round to update the global model, and  $C = |\mathbb{K}|/N$  denotes the fraction of selected clients. At round  $t$ , a selected client  $k \in \mathbb{K}$  executes  $T_{\text{gd}}$  local gradient descent iterations on the common model  $\mathbf{w}_{t-1}$  using its own training data  $D_k$  ( $D = \cup_{k \in \mathbb{K}} D_k$ ), and obtains the updated model  $\mathbf{w}_t^k$ , where the number of weights is denoted by  $n$  (i.e.,  $|\mathbf{w}_t^k| = |\Delta \mathbf{w}_t^k| = n$  for all  $k$  and  $t$ ). Each client  $k$  submits the update  $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$  to the server, which then updates the common model as follows:  $\mathbf{w}_t = \mathbf{w}_{t-1} + \sum_{k \in \mathbb{K}} \frac{|D_k|}{\sum_j |D_j|} \Delta \mathbf{w}_t^k$ , where  $|D_k|$  is known to the server for all  $k$  (a client's update is weighted with the size of its training data). The server stops training after a fixed number of rounds  $T_{\text{cl}}$ , or when the performance of the common model does not improve on a held-out data.

Note that each  $D_k$  may be generated from different distributions (i.e., Non-IID case), that is, any client's local dataset may not be representative of the population distribution [2]. This can happen, for example, when not all output classes are represented in every client's training data. The federated learning of neural networks is summarized in Alg. 1. In the sequel, each client is assumed to use the same model architecture.

---

**Algorithm 1:** FL-STD: Federated Learning

---

```

1 Server:
2   Initialize common model  $w_0$ 
3   for  $t = 1$  to  $T_{\text{cl}}$  do
4     Select  $\mathbb{K}$  clients uniformly at random
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\Delta \mathbf{w}_t^k = \text{Client}_k(\mathbf{w}_{t-1})$ 
7     end
8      $\mathbf{w}_t = \mathbf{w}_{t-1} + \sum_k \frac{|D_k|}{\sum_j |D_j|} \Delta \mathbf{w}_t^k$ 
9   end
10  Output: Global model  $\mathbf{w}_t$ 
11 Client $_k(\mathbf{w}_{t-1}^k)$ :
12   $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{\text{gd}})$ 
13  Output: Model update  $(\mathbf{w}_t^k - \mathbf{w}_{t-1}^k)$ 

```

---

The motivation of federated learning is three-fold: first, it aims to provide confidentiality of each participant's training data by sharing only model updates instead of potentially sensitive training data. Second, in order to decrease communication costs, clients can perform multiple local SGD iterations before sending their update back to the server. Third, in each round, only a few clients are required to perform local training of the common model,

---

**Algorithm 2:** Stochastic Gradient Descent

---

**Input:**  $D$  : training data,  $T_{\text{gd}}$  : local epochs,  $\mathbf{w}$  : weights

```
1 for  $t = 1$  to  $T_{\text{gd}}$  do
2   Select batch  $\mathbb{B}$  from  $D$  randomly
3    $\mathbf{w} = \mathbf{w} - \eta \nabla f(\mathbb{B}; \mathbf{w})$ 
4 end
Output: Model  $\mathbf{w}$ 
```

---

which further diminishes communication costs and makes the approach especially appealing with large number of clients.

However, several prior works have demonstrated that model updates do leak potentially sensitive information [3], [4]. Hence, simply not sharing training data *per se* is not enough to guarantee their confidentiality.

## 2.2. Differential Privacy

Differential privacy allows a party to privately release information about a dataset: a function of an input dataset is perturbed, so that any information which can differentiate a record from the rest of the dataset is bounded [6].

**Definition 1** (Privacy loss). *Let  $\mathcal{A}$  be a privacy mechanism which assigns a value  $\text{Range}(\mathcal{A})$  to a dataset  $D$ . The privacy loss of  $\mathcal{A}$  with datasets  $D$  and  $D'$  at output  $O \in \text{Range}(\mathcal{A})$  is a random variable  $\mathcal{P}(\mathcal{A}, D, D', O) = \log \frac{\Pr[\mathcal{A}(D)=O]}{\Pr[\mathcal{A}(D')=O]}$  where the probability is taken on the randomness of  $\mathcal{A}$ .*

**Definition 2** ( $(\epsilon, \delta)$ -Differential Privacy [6]). *A privacy mechanism  $\mathcal{A}$  guarantees  $(\epsilon, \delta)$ -differential privacy if for any database  $D$  and  $D'$ , differing on at most one record,  $\Pr_{O \sim \mathcal{A}(D)}[\mathcal{P}(\mathcal{A}, D, D', O) > \epsilon] \leq \delta$ .*

Intuitively, this guarantees that an adversary, provided with the output of  $\mathcal{A}$ , can draw almost the same conclusions (up to  $\epsilon$  with probability larger than  $1 - \delta$ ) about any record no matter if it is included in the input of  $\mathcal{A}$  or not [6]. That is, for any record owner, a privacy breach is unlikely to be due to its participation in the dataset.

*Moments Accountant.* Differential privacy maintains composition; the privacy guarantee of the  $k$ -fold adaptive composition of  $\mathcal{A}_{1:k} = \mathcal{A}_1, \dots, \mathcal{A}_k$  can be computed using the moments accountant method [15]. In particular, it follows from Markov's inequality that  $\Pr[\mathcal{P}(\mathcal{A}, D, D', O) \geq \epsilon] \leq \mathbb{E}[\exp(\lambda \mathcal{P}(\mathcal{A}, D, D', O))] / \exp(\lambda \epsilon)$  for any output  $O \in \text{Range}(\mathcal{A})$  and  $\lambda > 0$ . This implies that  $\mathcal{A}$  is  $(\epsilon, \delta)$ -DP with  $\delta = \min_{\lambda} \exp(\alpha_{\mathcal{A}}(\lambda) - \lambda \epsilon)$ , where  $\alpha_{\mathcal{A}}(\lambda) = \max_{D, D'} \log \mathbb{E}_{O \sim \mathcal{A}(D)}[\exp(\lambda \mathcal{P}(\mathcal{A}, D, D', O))]$  is the log of the moment generating function of the privacy loss. The privacy guarantee of the composite mechanism  $\mathcal{A}_{1:k}$  can be computed using that  $\alpha_{\mathcal{A}_{1:k}}(\lambda) \leq \sum_{i=1}^k \alpha_{\mathcal{A}_i}(\lambda)$  [15].

*Gaussian Mechanism.* There are a few ways to achieve DP, including the Gaussian mechanism [6]. A fundamental concept of all of them is the *global sensitivity* of a function [6].

**Definition 3** (Global  $L_p$ -sensitivity). *For any function  $f : \mathcal{D} \rightarrow \mathbb{R}^n$ , the  $L_p$ -sensitivity of  $f$  is  $\Delta_p f = \max_{D, D'} \|f(D) - f(D')\|_p$ , for all  $D, D'$  differing in at most one record, where  $\|\cdot\|_p$  denotes the  $L_p$ -norm.*

The Gaussian Mechanism [6] consists of adding Gaussian noise to the true output of a function. In particular, for any function  $f : \mathcal{D} \rightarrow \mathbb{R}^n$ , the Gaussian mechanism is defined as adding i.i.d Gaussian noise with variance  $(\Delta_2 f \cdot \sigma)^2$  and zero mean to each coordinate value of  $f(D)$ . Recall that the pdf of the Gaussian distribution with mean  $\mu$  and variance  $\xi^2$  is

$$\text{pdf}_{\mathcal{G}(\mu, \xi)}(x) = \frac{1}{\sqrt{2\pi}\xi} \exp\left(-\frac{(x - \mu)^2}{2\xi^2}\right) \quad (1)$$

In fact, the Gaussian mechanism draws vector values from a multivariate spherical (or isotropic) Gaussian distribution which is described by random variable  $\mathcal{G}(f(D), \Delta_2 f \cdot \sigma \mathbf{I}_n)$ , where  $n$  is omitted if its unambiguous in the given context.

## 2.3. Compressive Sensing

Compressive Sensing (CS) introduced in [11]–[13] aims to recover the original signal from significantly fewer samples (or measurements) than other traditional sampling techniques, which are based on the Nyquist-Shannon theorem, by exploiting the sparsity of the signal.

Consider a signal  $\mathbf{x} \in \mathbb{R}^n$  which admits a sparse representation  $\mathbf{s} \in \mathbb{R}^n$ , that is, there exists a *sparsity orthonormal basis* with matrix  $\Psi \in \mathbb{R}^{n \times n}$  such that:

$$\mathbf{x} = \Psi \mathbf{s} \quad (2)$$

Here,  $\mathbf{s}$  is  $U$ -sparse if  $\|\mathbf{s}\|_0 = U$ .  $\Psi$  can denote any linear transformation, such as Discrete Fourier/Cosine or Wavelet Transform, which render the original signal  $\mathbf{x}$  sparse. If  $\mathbf{x}$  is already sparse, then  $\Psi$  can be the identity matrix which corresponds to the canonical sparsity basis.

In CS,  $\mathbf{x}$  is reconstructed from some of its *linear measurements*. For  $m$  measurements, the signal is “sampled” in  $m$  values  $\mathbf{y}_j = \langle \phi_j, \mathbf{x} \rangle$  ( $1 \leq j \leq m$ ), where the vectors  $\phi_j \in \mathbb{R}^n$  constitute the sensing basis matrix  $\Phi = (\phi_1, \phi_2, \dots, \phi_m)^\top \in \mathbb{R}^{m \times n}$ . Here,  $m = r \times n$ , where  $r$  is the compression ratio. Therefore, the compression operator  $\mathcal{C}$  is defined as:

$$\mathcal{C}(\mathbf{x}, m) = \mathbf{y} = \Phi \mathbf{x} = \Phi \Psi \mathbf{s} = \Theta \mathbf{s} \quad (3)$$

where  $\Theta$  is the sparsity sensing matrix.

There are several options to select the sensing matrix  $\Phi$ . When  $\Phi$  is a random matrix (e.g., each element of  $\Phi$  is an iid sample from  $\mathcal{G}(0, 1/m)$ ), then  $\Psi$  works well with an arbitrary sparsity basis [16]. On the other hand, the numerical reconstruction of  $\mathbf{x}$  in that case has a complexity of  $O(mn)$  which can be very large (recall that  $n$  is the model size in the order of  $10^6$ ). Another (faster) option for the sensing matrix  $\Phi$  is when it is composed of random  $m$  rows of the matrix of the (real) Discrete Fourier/Cosine Transform. Then, matrix multiplication can be executed with the Fast Fourier Transform (FFT) in  $O(n \ln n)$ , but such sensing matrix provides accurate reconstruction if  $\Psi$  is the identity matrix, i.e.  $\mathbf{x}$  is already sparse [16]. Fortunately, this usually holds for gradient vectors (or can be made as such by sparsification without significantly affecting convergence) and hence we will use this option in this paper.

In order to recover  $\mathbf{s}$  from  $\mathbf{y}$ , one has to solve a system of linear equations with  $m$  equations and  $n$  unknowns. Although this system seems underdetermined



because  $m < n$ , CS exploits the  $U$ -sparsity of  $\mathbf{s}$  for the reconstruction. It aims to reconstruct the sparse vector  $\mathbf{s}$  from  $\mathbf{y} = \Theta\mathbf{s}$  given the sparsity sensing basis  $\Theta$  by solving the following optimization problem:

$$\arg \min_{\mathbf{s}} \|\mathbf{s}\|_0 \quad \text{s.t.} \quad \mathbf{y} = \Theta\mathbf{s}$$

Since this optimization problem is NP-complete [16], [17], it is further relaxed into the following slightly different problem called Basis Pursuit (BP) [18]:

$$\arg \min_{\mathbf{s}} \|\mathbf{s}\|_1 \quad \text{s.t.} \quad \mathbf{y} = \Theta\mathbf{s}$$

Indeed, the convex  $L_1$ -norm usually approximates the non-convex  $L_0$ -norm well, and the relaxed optimization problem can be efficiently solved with any convex optimization technique [16] (e.g., with an LP solver).

When the measurements  $\mathbf{y}$  are noisy (i.e.,  $\mathbf{y} = \Theta\mathbf{s} + \mathbf{z}$ , where  $\mathbf{z} \in \mathbb{R}^m$  is the additional bounded iid noise, i.e.  $\|\mathbf{z}\|_2 \leq \kappa$ ), then the following convex quadratic variant of BP called Basis Pursuit Denoising (BPDN) is rather considered:

$$\mathcal{R}(\mathbf{y}, \kappa) = \arg \min_{\mathbf{s}} \|\mathbf{s}\|_1 \quad \text{s.t.} \quad \|\mathbf{y} - \Theta\mathbf{s}\|_2 \leq \kappa$$

and therefore

$$\mathcal{D}(\mathbf{y}, n) = \Psi \left( \arg \min_{\mathbf{s}} \frac{1}{2} \|\mathbf{y} - \Theta\mathbf{s}\|_2^2 + \lambda \|\mathbf{s}\|_1 \right), \quad (4)$$

Eq. (4) defines our decompression operator and is an instance of convex quadratic programming. In this paper, we use the Orthant-Wise Limited-memory Quasi-Newton (OWL-QN) algorithm [19], [20], an extension of Limited-memory BFGS, which is a numerical scalable optimization procedure that can efficiently solve Eq. (4).

When  $\lambda \rightarrow 0$ , the problem in Eq. (4) becomes BP because  $\lambda \|\mathbf{s}\|_1$  tends to 0. In the case of non-noisy sensing measurements, a BP decoder is more adapted to reconstruct the sparse signal  $\mathbf{s}$ . Otherwise, BPDN is more suited. This has particular importance in our case when the compressed vector (measurements) are noised to guarantee Differential Privacy, i.e.,  $\mathbf{y} = \Theta\mathbf{s} + \mathbf{z}$  where  $\mathbf{z} \sim \mathcal{N}(0, S\mathbf{I}\sigma)$  (see Section 3.2.2). Approximate signal reconstruction from noisy measurements have been theoretically justified in [21] from the Restricted Isometry Property of  $\Theta$ .

**Definition 4** (Restricted Isometry Property (RIP) [22]). *The  $U$ -restricted isometry constants  $0 \leq \delta_U < 1$  of a matrix  $\Theta \in \mathbb{R}^{m \times n}$  is defined as the smallest number such that:*

$$(1 - \delta_U) \|\mathbf{s}\|_2^2 \leq \|\Theta\mathbf{s}\|_2^2 \leq (1 + \delta_U) \|\mathbf{s}\|_2^2$$

for all  $U$ -sparse vector  $\mathbf{s} \in \mathbb{R}^n$  and we say that the matrix  $\Theta$  obeys the Restricted Isometry Property (or  $\text{RIP}(U, \delta_U)$ ) of order  $U < m$ .

**Theorem 1** (Reconstruction error of BPDN [21]). *If  $\Theta$  is  $\text{RIP}(2U, \delta_U)$  and  $\delta_U < \sqrt{2} - 1$ , then  $\|\mathbf{s} - \mathcal{R}(\mathbf{y}, \kappa)\|_2 \leq C\kappa + (D/\sqrt{K})\|\mathbf{s} - \mathbf{s}_K\|_1$ , where  $C$  and  $D$  are constants and  $\mathbf{s}_K$  is a vector with all but the  $K$ -largest entries of  $\mathbf{s}$  set to zero<sup>2</sup>.*

2. For instance, for  $\delta_U = 0.2$ ,  $C < 4.2$  and  $D < 8.5$

Finally, notice that the compression operator  $\mathcal{C}$  in Eq. (3) is *linear*, which means that:

$$\sum_i \mathcal{C}(\mathbf{x}_i, m) = \mathcal{C} \left( \sum_i \mathbf{x}_i, m \right)$$

and therefore

$$\mathcal{D} \left( \sum_i \mathcal{C}(\mathbf{x}_i, m) \right) \approx \sum_i \mathbf{x}_i$$

This linearity allows to combine secure aggregation and compressive sensing described in Section 3.2.2.

## 2.4. Error Propagation

Biased estimation of the gradients may prevent model convergence unless the approximation error introduced by lossy compression techniques, such as compressive sensing, sketching, or quantization, is accumulated and re-injected in every optimization round [23] as follows:

$$\begin{aligned} \mathbf{g}_t &= \nabla f(\mathbb{B}, \mathbf{w}_{t-1}) : \text{Computing gradients on batch } \mathbb{B} \\ \mathbf{p}_t &= \eta \mathbf{g}_t + \mathbf{e}_{t-1} : \text{Error feedback (correction)} \\ \Delta_t &= \mathcal{D}(\mathcal{C}(\mathbf{p}_t)) : \text{Reconstruction of } \mathbf{p}_t \\ \mathbf{w}_t &= \mathbf{w}_{t-1} - \Delta_t : \text{Updating model parameters (weights)} \\ \mathbf{e}_t &= \mathbf{p}_t - \Delta_t : \text{Error accumulation} \end{aligned}$$

The corrected direction  $\mathbf{p}_t$  is obtained by adding the error  $\mathbf{e}_{t-1}$  accumulated over all iterations to  $\mathbf{g}_t$  (see Alg. 2 in [23] for more details). Here, the error is calculated based on the biased estimation of the update given by  $\mathcal{D}(\mathcal{C}(\mathbf{p}_t))$ .

## 3. Federated Learning with Compressive Sensing

In the FL-STD scheme, presented in Section 2, each randomly selected client sends its complete model update to the server. Knowing that a model has on average millions of parameters (each is a floating point value represented on 32 bits), the network can suffer from large traffic.

To decrease large network traffic, we adapt compressive sensing to federated learning. The new algorithm is called FL-CS. Moreover, this scheme is also extended to a privacy-preserving version, called FL-CS-DP, which aims to protect the training data of every participant. We show that compression improves model performance with Differential Privacy by reducing the added noise compared to the uncompressed DP variant of federated learning. Hence, both FL-CS and FL-CS-DP improve bandwidth efficiency, and in addition, FL-CS-DP also boosts the accuracy of differentially private federated learning.

In what follows, we will first describe the non-private scheme FL-CS and then the privacy-preserving FL-CS-DP.

### 3.1. FL-CS: Federated Learning with Compressive Sensing

CS assumes the sparsity of the reconstructed signal in a specific basis domain  $\Psi$  as explained in Section 2.

We assume the model update (as a signal) to be already sparse in the time domain, that is,  $\Psi$  is canonical sparsity basis (i.e.,  $\Psi = \mathbf{I}$ ), and therefore, the compression operator is  $\mathcal{C}(\Delta\mathbf{w}, m) = \Phi\Delta\mathbf{w}$ , where  $\Phi$  is composed of the first  $m$  rows of the matrix of the Discrete Cosine Transform (DCT) [24], [25]. Indeed, due to the large energy compaction property of DCT, the first coefficients, which correspond to the low frequency components of  $\Delta\mathbf{w}$ , tend to have the largest magnitude and hence convey the most information about the model update [26]. In fact, for a canonical sparsity basis  $\Psi = \mathbf{I}$ ,  $\Theta = \Phi$  is RIP with overwhelming probability as soon as  $m = O(U \ln^4 n)$  if  $\Delta\mathbf{w}$  is  $U$ -sparse [27]. Therefore, reconstruction is possible according to Theorem 1.

The decompression operator  $\mathcal{D}$  is defined Eq. (4). Note that the compression operator can be computed in  $O(n \ln n)$  with FFT and the decompression (or reconstruction) operator is implemented with the OWL-QN algorithm [19] which makes our approach reasonably fast in practice.

FL-CS is described in Alg. 3. A client first computes its update  $\Delta\mathbf{w}_t^k$  with SGD, and then transfers the compressed update  $\mathcal{C}(\Delta\mathbf{w}_t^k, m)$ , which consists of the first  $m$  DCT coefficients of the update (Line 18). The server takes the average of the client’s updates (Line 8), updates the momentum (Line 9), and computes the error  $\mathbf{e}_t$  (Line 10-12) due to compression following the error propagation technique described in Section 2.4. This error is accumulated over all federated rounds and added to the model (Line 13) to compensate its negative effect on convergence. The server uses OWL-QN [19], [20] to reconstruct the error-compensated aggregated model update  $\mathbf{s}_t \in \mathbb{R}^n$ . Finally, the server updates the global model as  $\mathbf{w}_t = \mathbf{w}_{t-1} + \mathbf{s}_t$  before re-distributing the updated model to a new set of clients  $\mathbb{K}$ .

Notice that the error  $\mathbf{e}_t$ , the averaged model update  $\mathbf{y}_t$ , as well as the momentum are maintained in the compressed domain and have a size of  $m$  instead of  $n$ . This is possible due to the linearity of the compression scheme which is detailed in Section 2.3.

**Scalable reconstruction:** Although OWL-QN is reasonably fast in practice, its computational overhead may not be tolerated with very large models. A scalable reconstruction is proposed as follows. On the client side, the update vector  $\Delta\mathbf{w}_t$  is shuffled and then splitted into  $P$  equally-sized chunks. Then, the compression operator  $\mathcal{C}$  is applied on each individual chunk. Finally, the compressed chunks are transferred to the server. On the server side, each chunk is reconstructed independently using OWL-QN. The decompressed chunks are concatenated, and the resulted vector with size  $n$  is reshuffled to obtain  $\mathbf{s}_t$  by inverting the client-side shuffling.

Shuffling is performed by each client identically which guarantees that the compressed chunks can still be aggregated by the server. In practice, this can be implemented by sharing a common random seed among all participants to initialize the shuffler. As the server also knows this seed, it can invert this shuffling and reconstruct the aggregated model updates.

Notice that, instead of reconstructing the complete update vector at once, the server performs reconstruction on smaller chunks which makes decompression faster.

In addition, shuffling guarantees that the sparsity of the chunks is proportional to the sparsity of the whole update vector (i.e., if the update vector is  $U$ -sparse then all its chunks are  $U/P$ -sparse). Hence, the same compression operator  $\mathcal{C}(\cdot, m/P)$  can be applied on every chunk without increasing the compression ratio (i.e., the compressed update still has a size of  $m$ ).

Note that shuffling is also performed identically over all the rounds to maintain the error.

## 3.2. FL-CS-DP: Differentially Private Federated Learning with Compressive Sensing

**3.2.1. Privacy Model.** We consider an adversary, or a set of colluding adversaries, who can access any update vector sent by the server or any clients at each round of the protocol. A plausible adversary is a participating entity, i.e. a malicious client or server, that wants to infer the training data used by other participants. The adversary is *passive* (i.e., honest-but-curious), that is, it follows the learning protocol faithfully.

Different privacy requirements can be considered depending on what information the adversary aims to infer. In general, private information can be inferred about:

- any record (user) in any dataset of any client (*record-level privacy*),
- any client/party (*client-level privacy*).

To illustrate the above requirements, suppose that several banks build a common model to predict the credit-worthiness of their customers. A bank certainly does not want other banks to learn the financial status of any of their customers (record privacy) and perhaps not even the average income of all their customers (client privacy).

Record-level privacy is a standard requirement used in the privacy literature and is usually weaker than client-level privacy. Indeed, client-level privacy requires to hide any information which is unique to a client including perhaps all its training data.

We aim at developing a solution that provides *client-level privacy and is also bandwidth efficient*. For example, in the scenario of collaborating banks, we aim at protecting any information that is unique to each single bank’s training data. The adversary should not be able to learn from the received model or its updates whether any client’s data is involved in the federated run (up to  $\epsilon$  and  $\delta$ ). We believe that this adversarial model is reasonable in many practical applications when the confidential information spans over multiple samples in the training data of a single client (e.g., the presence of a group a samples, such as people from a certain race). Differential Privacy guarantees plausible deniability not only to any groups of samples of a client but also to any client in the federated run. Therefore, any negative privacy impact on a party (or its training samples) cannot be attributed to their involvement in the protocol run.

**3.2.2. Operation.** FL-CS-DP is described in Alg. 4. Client-level differential privacy requires each client to add Gaussian noise to the compressed model updates. In particular, each client first calculates  $\mathbf{c}_t^k = \mathcal{C}(\Delta\mathbf{w}_t^k, m)$  (in Line 19), which is then clipped (in Line 20) to obtain  $\hat{\mathbf{c}}_t^k$  with  $L_2$ -norm at most  $S$ . Then, random noise

$\mathbf{z}_k \sim \mathcal{G}(0, S\sigma\mathbf{I}/\sqrt{|\mathbb{K}|})$  is added to  $\hat{\mathbf{c}}_t^k$  such that  $\sum_{k \in \mathbb{K}} (\hat{\mathbf{c}}_t^k + \mathbf{z}_k) = \sum_{k \in \mathbb{K}} \hat{\mathbf{c}}_t^k + \mathcal{G}(0, S\sigma\mathbf{I})$  as the sum of Gaussian random variables also follows Gaussian distribution<sup>3</sup> and then differential privacy is satisfied where  $\varepsilon$  and  $\delta$  can be computed using the moments accountant described in Section 2.2.

However, as the noise is inversely proportional to  $\sqrt{|\mathbb{K}|}$ ,  $\mathbf{z}_k$  is likely to be small if  $|\mathbb{K}|$  is too large. Therefore, the adversary accessing an individual update  $\hat{\mathbf{c}}_t^k + \mathbf{z}_k$  can almost learn a non-noisy update since  $\mathbf{z}_k$  is small. Hence, each client uses secure aggregation to encrypt its individual update before sending it to the server. Upon reception, the server sums the encrypted updates as:

$$\begin{aligned} \sum_{k \in \mathbb{K}} \mathbf{y}_t^k &= \sum_{k \in \mathbb{K}} \text{Enc}_{K_k}(\hat{\mathbf{c}}_t^k + \mathbf{z}_k) \\ &= \sum_{k \in \mathbb{K}} \hat{\mathbf{c}}_t^k + \sum_{k \in \mathbb{K}} \mathbf{z}_k \\ &= \sum_{k \in \mathbb{K}} \hat{\mathbf{c}}_t^k + \mathcal{G}(0, S\sigma\mathbf{I}) \end{aligned} \quad (5)$$

where  $\text{Enc}_{K_k}(\hat{\mathbf{c}}_t^k + \mathbf{z}_k) = \hat{\mathbf{c}}_t^k + \mathbf{z}_k + \mathbf{K}_k \pmod{p}$  and  $\sum_k \mathbf{K}_k = 0$  (see [7], [28] for details). Here the modulo is taken element-wise and  $p = 2^{\lceil \log_2(\max_k \|\hat{\mathbf{c}}_t^k + \mathbf{z}_k\|_\infty |\mathbb{K}|) \rceil}$ . Let  $\gamma_t^k = 1/\max\left(1, \frac{\|\mathbf{c}_t^k\|_2}{S}\right)$ . Then,

$$\begin{aligned} \sum_{k \in \mathbb{K}} \hat{\mathbf{c}}_t^k &= \sum_{k \in \mathbb{K}} \gamma_t^k \mathbf{c}_t^k \\ &= \sum_{k \in \mathbb{K}} \gamma_t^k \mathcal{C}(\Delta \mathbf{w}_t^k, m) \\ &= \mathcal{C}\left(\sum_{k \in \mathbb{K}} \gamma_t^k \Delta \mathbf{w}_t^k, m\right) \end{aligned} \quad (6)$$

where the last equality comes from the linearity of the compression operation (see Section 2.3). Plugging Eq. (6) into Eq. (5), we get that

$$\sum_{k \in \mathbb{K}} \mathbf{y}_t^k = \mathcal{C}\left(\sum_{k \in \mathbb{K}} \gamma_t^k \Delta \mathbf{w}_t^k, m\right) + \mathcal{G}(0, S\sigma\mathbf{I})$$

This is an instance of BPDN (see Section 2.3), and therefore the direct reconstruction of  $\sum_{k \in \mathbb{K}} \mathbf{y}_t^k$  would be an approximation of  $\sum_{k \in \mathbb{K}} \gamma_t^k \Delta \mathbf{w}_t^k$ . However, analogously to FL-CS, the server applies error propagation and computes the (noisy) error  $\mathbf{e}_t$  from  $\mathbf{y}_t = (1/|\mathbb{K}|) \sum_{k \in \mathbb{K}} \mathbf{y}_t^k$  (in Line 10), and decompresses  $\mathbf{e}_t$  into  $\mathbf{s}_t$  by using OWL-QN. Recall that the reconstruction algorithm solves the BPDN problem, where a sparse vector  $\mathbf{s}$  is reconstructed from  $m$  noisy measurements of the form  $\Theta \mathbf{s} + \mathbf{z}$ , where the noise  $\mathbf{z} \in \mathbb{R}^m$  is assumed to be identically and independently distributed over its elements with a Gaussian distribution [16], [18]. Since  $\mathbf{z} \sim \mathcal{G}(0, S\mathbf{I}\sigma)$  in our case, the reconstruction algorithm is therefore optimized to reconstruct the differentially private compressed vectors (see Theorem 1).

**Privacy analysis:** The server can only access the noisy aggregate which is sufficiently perturbed to ensure differential privacy; any client-specific information that could be inferred from the noisy aggregate is tracked and quantified by the moments accountant, described in Section 2.2, as follows.

3. More precisely,  $\sum_i \mathcal{G}(\nu_i, \xi_i) = \mathcal{G}(\sum_i \nu_i, \sqrt{\sum_i \xi_i^2})$

Let  $\eta_0(x|\xi) = \text{pdf}_{\mathcal{G}(0,\xi)}(x)$  and  $\eta_1(x|\xi) = (1 - C)\text{pdf}_{\mathcal{G}(0,\xi)}(x) + C\text{pdf}_{\mathcal{G}(1,\xi)}(x)$  where  $C$  is the sampling probability of a single client in a single round. Let

$$\alpha(\lambda|C) = \log \max(E_1(\lambda, \xi, C), E_2(\lambda, \xi, C)) \quad (7)$$

where  $E_1(\lambda, \xi, C) = \int_{\mathbb{R}} \eta_0(x|\xi, C) \cdot \left(\frac{\eta_0(x|\xi, C)}{\eta_1(x|\xi, C)}\right)^\lambda dx$  and  $E_2(\lambda, \xi, C) = \int_{\mathbb{R}} \eta_1(x|\xi, C) \cdot \left(\frac{\eta_1(x|\xi, C)}{\eta_0(x|\xi, C)}\right)^\lambda dx$ .

**Theorem 2** (Privacy of FL-CS-DP). *FL-CS-DP is  $(\min_\lambda(T_{\text{cl}} \cdot \alpha(\lambda|C) - \log \delta)/\lambda, \delta)$ -DP.*

Given a fixed value of  $\delta$ ,  $\varepsilon$  is computed numerically as in [15], [29].

The magnitude of the added Gaussian noise is proportional to the sensitivity  $S$ , which is in turn often proportional to the model size  $n$  [10]. Hence, when  $n$  becomes large, SGD often fails to converge due to the perturbation error caused by the added noise [10]. In our approach, the perturbation error is less since Gaussian noise is added to the compressed vector with size  $m < n$ . On the other hand, compression also induces some reconstruction error owing to its lossy nature. The total error is the sum of the reconstruction and the perturbation error and is quantified in Theorem 1. Finding the right trade-off between these two errors is the key to achieve good model quality.

---

#### Algorithm 3: FL-CS: Federated Learning

---

```

1 Server:
2   Initialize common model  $w_0$ ,  $\eta_G$ ,  $\rho$ ,  $\mathbf{u}_t = 0$ ,
    $\mathbf{e}_t = 0$ 
3   for  $t = 1$  to  $T_{\text{cl}}$  do
4     Select  $\mathbb{K}$  clients uniformly at random
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\mathbf{y}_t^k = \text{Client}_k(\mathbf{w}_{t-1})$ 
7     end
8      $\mathbf{y}_t = \sum_{k=1}^{|\mathbb{K}|} \frac{\mathbf{y}_t^k}{|\mathbb{K}|}$  : Averaging
9      $\mathbf{u}_t = \rho \mathbf{u}_{t-1} + \mathbf{y}_t$  : Momentum
10     $\mathbf{e}_t = \eta_G \mathbf{u}_t + \mathbf{e}_{t-1}$  : Error Feedback
11     $\mathbf{s}_t = \mathcal{D}(\mathbf{e}_t, n)$  : Reconstruction
12     $\mathbf{e}_t = \mathbf{e}_t - \mathcal{C}(\mathbf{s}_t, m)$  : Error accumulation
13     $\mathbf{w}_t = \mathbf{w}_{t-1} + \mathbf{s}_t$  : Update
14  end
Output: Global model  $\mathbf{w}_t$ 

15
16 Client $_k(\mathbf{w}_{t-1}^k)$ :
17    $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{\text{gd}})$ 
18    $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$ 
Output: Model update  $\mathcal{C}(\Delta \mathbf{w}_t^k, m)$ 

```

---

## 4. Experimental Results

The goal of this section is to evaluate the performance of our proposed schemes FL-CS and FL-CS-DP on a benchmark dataset and a realistic in-hospital mortality prediction scenario. We aim at evaluating their performance with different levels of compression and comparing them with the performance of the following learning protocols:

- FL-STD: It is described in Section 2.1 (see Alg. 1).
- FL-RND: This baseline follows the algorithm of FL-STD except that a random subset of the update vector with size  $m \leq n$  is sent to the server instead

---

**Algorithm 4:** FL-CS-DP: Private Compressive Sensing Federated Learning

---

```
1 Server:
2   Initialize common model  $w_0$ ,  $\eta_G$ ,  $\rho$ ,  $\mathbf{u}_t = 0$ ,
    $\mathbf{e}_t = 0$ 
3   for  $t = 1$  to  $T_{cl}$  do
4     Select  $\mathbb{K}$  clients uniformly at random
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\mathbf{y}_t^k = \text{Client}_k(\mathbf{w}_{t-1})$ 
7     end
8      $\mathbf{y}_t = \sum_{k=1}^{|\mathbb{K}|} \frac{\mathbf{y}_t^k}{|\mathbb{K}|}$  : Averaging
9      $\mathbf{u}_t = \rho \mathbf{u}_{t-1} + \mathbf{y}_t$  : Momentum
10     $\mathbf{e}_t = \eta_G \mathbf{u}_t + \mathbf{e}_{t-1}$  : Error Feedback
11     $\mathbf{s}_t = \mathcal{D}(\mathbf{e}_t, \mathbf{n})$  : Reconstruction
12     $\mathbf{e}_t = \mathbf{e}_t - \mathcal{C}(\mathbf{s}_t, m)$  : Error accumulation
13     $\mathbf{w}_t = \mathbf{w}_{t-1} + \mathbf{s}_t$  : Update
14  end
Output: Global model  $\mathbf{w}_t$ 

15
16 Client $_k(\mathbf{w}_{t-1}^k)$ :
17    $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{gd})$ 
18    $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$ 
19    $\mathbf{c}_t^k = \mathcal{C}(\Delta \mathbf{w}_t^k, m)$ 
20    $\hat{\mathbf{c}}_t^k = \mathbf{c}_t^k / \max\left(1, \frac{\|\mathbf{c}_t^k\|_2}{S}\right)$ 
Output:  $\text{Enc}_{K_k}(\mathcal{G}(\hat{\mathbf{c}}_t^k, S\mathbf{I}\sigma/\sqrt{|K|}))$ 
```

---

---

**Algorithm 5:** FL-STD-DP: Federated Learning with Client Privacy

---

```
1 Server:
2   Initialize common model  $w_0$ 
3   for  $t = 1$  to  $T_{cl}$  do
4     Select  $\mathbb{K}$  clients randomly
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\Delta \tilde{\mathbf{w}}_t^k = \text{Client}_k(\mathbf{w}_{t-1})$ 
7     end
8      $\mathbf{w}_t = \mathbf{w}_{t-1} + \frac{1}{|\mathbb{K}|} \sum_k \Delta \tilde{\mathbf{w}}_t^k$ 
9   end
10 Client $_k(\mathbf{w})$ :
11    $\mathbf{w}_{t-1}^k = \mathbf{w}$ 
12    $\Delta \mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{gd}) - \mathbf{w}_{t-1}^k$ 
13    $\Delta \hat{\mathbf{w}}_t^k = \Delta \mathbf{w}_t^k / \max\left(1, \frac{\|\Delta \mathbf{w}_t^k\|_2}{S}\right)$ 
Output:  $\text{Enc}_{K_k}(\mathcal{G}(\Delta \hat{\mathbf{w}}_t^k, S\mathbf{I}\sigma/\sqrt{|K|}))$ 
```

---

of the complete update of size  $n$ . Each client selects the same random subset of coordinates from the update vector, but a different subset in every round. The server then averages the received updates before updating only the corresponding  $m$  weights. Note that if  $m = n$ , FL-RND is equivalent to FL-STD (see Alg. 6).

- FL-FREQ: In this baseline, a client transforms the model update to the frequency domain by using DCT [24], [25], and then the first  $m$  coefficients (low frequency components) are extracted and sent to the server as in FL-CS. However, as opposed to FL-CS, the server reconstructs the aggregated update vector by applying the inverse DCT on the aggregated compressed vectors where the last  $n - m$  coefficients are zeroed out (see Alg. 7). This baseline corresponds to a low-pass filter applied on the update vector.  $\Phi$  in Alg. 7 is composed of

the first  $m$  rows of the matrix of the DCT.

## 4.1. Medical Dataset

### 4.1.1. The In-hospital Mortality Prediction Scenario.

The ability to accurately predict the risks in the patient’s perspectives of evolution is a crucial prerequisite in order to adapt the care that certain patients receive [30].

We consider the scenario where several hospitals are collaborating to train models for in-hospital mortality prediction using our Federated Learning schemes. This well-studied real-world problem consists in trying to precisely identify the patients who are at risk of dying from complications during their hospital stay [30]–[32]. As commonly found in the literature [30], for such predictions, we focus on hospital admissions of adults hospitalized for at least 3 days, excluding elective admissions.

**4.1.2. The Premier Healthcare Database.** We used EHR data from the Premier healthcare database<sup>4</sup> which is one of the largest clinical databases in the United States, collecting information from millions of patients over a period of 12 months from 415 hospitals in the USA [30]. These hospitals are supposedly representative of the United States hospital experience [30]. Each hospital in the database provides discharge files that are dated records of all billable items (including therapeutic and diagnostic procedures, medication, and laboratory usage) which are all linked to a given patient’s admission [30], [33].

The initial snapshot of the database used in our work (before pre-processing step) comprises the EHR data of 1,271,733 hospital admissions. Electronic Health Record (EHR) is a digital version of a patient’s paper chart readily available in hospitals. For developing supervised learning and specifically deep learning models, we focus on a specific set of features from EHR data. The features of interest that capture the patients information are summarized in Table 1. There is a total of 24,428 features per patient, mainly due to the variety of drugs possibly served. As in [31], we also removed all the features which appear on less than 100 patients’ records, hence, the number of features was reduced to 7,280 features.

The Medication regimen complexity index (MRCI) [34] is an aggregate score computed from a total of 65 items, whose purpose is to indicate the complexity of the patient’s situation. The minimum MRCI score for a patient is 1.5, which represents a single tablet or capsule taken once a day as needed (single medication). However the maximum is not defined since the number of medications increases the score [34]. In our case, after statistical analysis of our dataset, we consider the MRCI score as ranging from 2 to 60.

Most real datasets like ours are generally imbalanced with a skewed distribution between the classes. In our case, the positive cases (patients who die during their hospital stay) represent only 3% of all patients. Table 2 gives more details about this distribution after the pre-processing step which is discussed in A.1. To deal with this well-known problem, we have decided to use down-

4. <https://www.premierinc.com/newsroom/education/premier-healthcare-database-whitepaper>



TABLE 1: Descriptions of features

Features	Descriptions
Age	Value in the range of 15 and 89
Gender	Male, Female or Unknown
Admission type	Emergency, Urgent, Trauma Center: visits to a trauma center/hospital or Unknown
MRCI	Medication regimen complexity index score (ranging from 2 to 60)
Drugs	Drugs given to the patient on the 1 <sup>st</sup> day of hospitalization. There is a total of 24,419 possible drugs that can be served.

sampling technique [35], [36], a standard solution used for this purpose<sup>5</sup>.

**4.1.3. Model architecture.** As in [31], we use a fully connected neural network model with the following architecture: two hidden layers of 200 units, which use a Relu activation function followed by an output layer of 1 unit with sigmoid activation function and a binary cross entropy loss function. A dropout layer with a rate set to 0.5 is used between each hidden layer and between the last hidden layer and the output layer. This results in 1,496,601 parameters in total. We tune  $\eta$  from 0.01 to 0.5 with an increment value of 0.005. As in [37], we fix the momentum parameter  $\rho$  to 0.9 and we tuned the global learning rate  $\eta_G$  from 0.05 to 2.0 with an increment value of 0.05. The number of chunks is  $P = 200$ . The hyperparameters used by each of the considered schemes are summarized in Table 6.

The sensitivity  $S$  is selected during an initialization round for each scheme by taking the median value over  $N$   $L_2$ -norm values. We also noticed that the sensitivity of FL-CS-DP, FL-RND and FL-FREQ are nearly equivalent for the same level of compression. For this reason, the same sensitivity value is used for each compressed scheme and for the same compression ratio. Table 5 and Table 6 show the selected clipping threshold (i.e., sensitivity  $S$ ) for each dataset and according to each compression ratio.

TABLE 2: Number of instances for our case study. The Medical dataset contains in total 1,271,733 records.

Data	Positive cases	Negative cases	Ratio	Total
Train	32,106	985,280	3.16%	1,017,386
Test	7,882	246,465	3.10%	254,347

## 4.2. Fashion-MNIST

**4.2.1. Data Description.** Fashion-MNIST database of fashion articles consists of 60,000 28x28 grayscale images of 10 fashion categories, along with a test set of 10,000 images [38] [39].

### 4.2.2. Data pre-processing & experimental setup.

**Preprocessing:** The pixel of each image is an unsigned integer in the range between 0 and 255. We rescale them to the range [0,1] instead.

**Model architecture:** For Fashion-MNIST, we use a model [2] with the following architecture: a convolutional neural network (CNN) with two 5x5 convolution layers (the first with 32 filters, the second with 64, each followed with

2x2 max pooling), a fully connected layer with 512 units and ReLu activation, and a final softmax output layer. This results in 1,663,370 parameters in total. We tune  $\eta$  from 0.01 to 0.5 with an increment value of 0.005. As in [37], we fix the momentum parameter  $\rho$  to 0.9 and we tuned the global learning rate  $\eta_G$  from 0.05 to 2.0 with an increment value of 0.05. The number of chunks used is  $P = 200$ . The hyperparameters used by each of the considered schemes are summarized in Table 5.

## 4.3. Computational environment

Our experiments were performed on a server running Ubuntu 18.04 LTS equipped with a Intel(R) Xeon(R) Silver 4114 CPU @ 2.20GHz, 192GB RAM, and two NVIDIA Quadro P5000 GPU card of 16 Go each. We use Keras 2.2.0 [40] with a TensorFlow backend 1.12.0 [41] and Numpy 1.14.3 [42] to implement our models and experiments. We use Python 3.6.5 and our code runs on a Docker container to simplify reproducibility.

## 4.4. Results

Table 3 represents the best accuracy over 200 rounds for each scheme on the Fashion-Mnist dataset. *Round* corresponds to the round when the best accuracy is reached and *Cost* is the average bandwidth consumption calculated as:  $r \times n \times 32 \times Round \times C$ , where 32 is the number of bits necessary to represent a float value,  $n$  is the uncompressed model size,  $r = \frac{m}{n}$ ,  $m$  is the compressed model size,  $C$  is the sampling probability of a client, and *Round* is the round when we get the the best accuracy.

Table 4 represents the best balanced accuracy over 100 rounds for each scheme on the Medical dataset. AUROC (area under the receiver operating characteristic curve [43]) corresponds to the AUROC value when the best balanced accuracy is reached, round is also the round when we get the best balanced accuracy, and finally, Cost is the average bandwidth consumption calculated as for the Fashion-MNIST dataset described above.

Without DP, notice that our FL-CS scheme outperforms FL-RND and FL-FREQ whatever the considered compression ratio or the dataset are. Also, compared to FL-STD, our scheme started to reach the same accuracy from a compression ratio  $r$  being equal or greater than 0.1 for both datasets, although the differences between FL-CS and FL-STD for a compression ratio of 0.05 are only of 6%<sup>6</sup> and 1%<sup>7</sup> for the Fashion-Mnist and the medical datasets, respectively. However, FL-STD consumes much more bandwidth than FL-CS. Indeed, FL-CS reduces the bandwidth cost by 95% compared to FL-STD with a compression ratio of 0.05 for both datasets, while the

5. We have also tested weighted loss function and oversampling techniques. But, we noticed experimentally that downsampling technique outperforms the others whatever the considered scheme.

6. Based on the accuracy

7. Based on the balanced accuracy [44], [45]

bandwidth cost is reduced to 80% and 85% with a compression ratio of 0.2 for Fashion-MNIST and the medical data, respectively.

Surprisingly, for the smallest compression ratio 5%, FL-CS-DP performs as well as FL-RND-DP or FL-FREQ in terms of accuracy and much better in terms of bandwidth consumption. Indeed, FL-CS-DP with a compression ratio of 5% reached 0.78 of accuracy on Fashion-MNIST, however, our baseline FL-RND needs a compression ratio of 10% to reach a similar result (0.77) and 20% to have slightly better result (0.80). The same holds for the medical dataset, where FL-CS-DP reached 0.69 and 0.76 of balanced accuracy and Auroc, respectively. However, our other baseline FL-FREQ needs a compression ratio of 20% to reach the same performance. FL-CS-DP performs better for a small compression ratio. Indeed, FL-CS-DP reaches 0.78, 0.73 and 0.66 for 5%, 10% and 20% of accuracy, respectively, on Fashion-MNIST. The accuracy degradation with DP can be explained by the fact increasing the compression ratio  $r$  also increases the sensitivity  $S$  which has a direct impact on the additive Gaussian noise as explained in Section 3.2.2. Indeed, the standard deviation of the normal distribution is  $\sigma \times S$ .

On both datasets, FL-STD-DP suffers from the noise due to the large sensitivity which is the largest one in Table 5 and 6 for a compression ratio of 1.0 (uncompressed model). Even for FL-RND and FL-FREQ, the gap between the non-private and the private version is larger when the compression ratio increases for both datasets. As the noise proportional to  $S$  and is added to every coordinate, the norm of the added noise increases with the model size  $n$ . This has negative impact on model convergence for a large  $n$  as discussed in [10]. By decreasing  $n$ , compression helps reach better utility.

On Fashion-MNIST, FL-CS-DP with a compression ratio of 0.05 outperforms FL-STD-DP on both utility and bandwidth preservation. However, and even though FL-CS-DP reduces the bandwidth cost by 95% which is not negligible, they have both comparable accuracy on the medical data. It can be explained by the reduction of the noise due to the reduction of  $S$  and  $\sigma$  (see Table 6 and Table 5) needed to reach an  $\epsilon$  value of at most 1 after  $T_{cl}$  rounds.

There is a possible tradeoff between the privacy, communication cost, and utility. Indeed, having a small  $\epsilon$  (better privacy) results in a reduction of the communication costs while it decreases accuracy. FL-STD-DP, for example, converges to the best accuracy (61%) after only 25 rounds with early stopping, which results on high privacy ( $\epsilon=0.69$ ) and low communication cost (only 22.18 Megabyte). However, the accuracy degradation is more important (about 30% which is the worst accuracy degradation indicated in Table 3). Indeed, the large amount of added noise impacts the convergence of the model which can not achieve an accuracy larger than 61%.

Finally, we highlight a trade-off for FL-CS-DP. As mentioned above, FL-CS-DP performs better when the smallest compression ratio  $r$  is used, as the sensitivity for this level of compression is the smallest one. On the other hand, the compression ratio cannot be decreased arbitrarily as it will result in large reconstruction error. Therefore, one has to find the smallest compression ratio that is small enough to reduce the perturbation error but

large enough to induce small reconstruction error.

## 5. Related work

**Privacy of Federated Learning:** There exist a few inference attacks specifically designed against federated learning schemes. In [4], the adversary’s goal is to infer whether records with a specific property are included in the training dataset of the other participants (called batch property inference). The authors demonstrate the attack by inferring whether black people are included in any of the training datasets, where the common model is trained for gender classification (i.e., the inferred property is independent of the learning objective). The adversary is supposed to have access to the aggregated model update of honest participants. In [3], the proposed attack infers if a specific person is included in the training dataset of the participants (aka, Membership inference). The adversary extracts the following features from every snapshot of the common model, which is a neural network: output value, hidden layers, loss values, and the gradient of the loss with respect to the parameters of each layer. These features are used to train a membership inference model, which is a convolutional neural network.

The concept of Client-based Differential Privacy has been introduced in [46] and [47], where the goal is to hide any information that is specific to a single client’s training data. These algorithms bound and noise the contribution of a single client’s instead of a single record in the client’s dataset. The noise is added by the server, hence, unlike our solution, these works assume that the server is trusted. Also, the noise is drawn from continuous distributions.

**Bandwidth Optimization in Federated Learning:** Different quantization methods have been proposed to save the bandwidth and reduce the communication costs in federated learning. They can be divided into two main groups: unbiased and biased methods. The unbiased approximation techniques use probabilistic quantization schemes to compress the stochastic gradient and attempt to approximate the true gradient value as much as possible [48] [49] [50] [51]. However, biased approximations of the stochastic gradient can still guarantee convergence both in theory and practice [52]–[54]. In signSGD [52], all the clients calculate the stochastic gradient based on a single mini-batch and then send the sign vector of this gradient to the server. The server calculates the aggregated sign vector by taking the median (majority vote) and sends the signs of the aggregated signs back to each client.

A different line of works exploit the sparsity of model updates to compress model updates. Our work belongs to this line. The authors in [55] use CS for low-complexity energy-efficient ECG compression. Although compressed sensing was primarily designed for compression [11], [13], it was extended for denoising as in [56], [57] where compressive sensing is used for the purpose of denoising. In [58], compressed sensing based denoising and certain artificial intelligence are combined to improve the prediction performance.

CS was also used with DP in [59]. The authors show that the amount of noise is reduced from  $O(\sqrt{n})$  to  $O(\log(n))$ , when the noise is added on the sampled coefficients instead of the original database.

Compression ratio ( $r$ )	Algorithms	Performance			
		Accuracy	Round	Cost (Megabyte)	$\epsilon$
0.05	FL-RND	0.73	192	8.52	N/A
	FL-FREQ	0.73	189	8.38	N/A
	FL-CS	0.82	200	8.87	N/A
	FL-RND-DP	0.73	196	8.69	0.99
	FL-FREQ-DP	0.72	200	8.87	1
	FL-CS-DP	0.78	197	8.74	1
0.1	FL-RND	0.78	200	17.74	N/A
	FL-FREQ	0.78	197	17.48	N/A
	FL-CS	0.85	199	17.65	N/A
	FL-RND-DP	0.77	199	17.65	1
	FL-FREQ-DP	0.76	200	17.74	1
	FL-CS-DP	0.73	101	8.96	0.84
0.2	FL-RND	0.82	200	35.49	N/A
	FL-FREQ	0.82	195	34.60	N/A
	FL-CS	0.87	193	34.24	N/A
	FL-RND-DP	0.80	199	35.31	1
	FL-FREQ-DP	0.79	200	35.49	1
	FL-CS-DP	0.66	150	26.61	0.92
1.0	FL-STD	0.87	191	169.44	N/A
	FL-STD-DP	0.61	25	22.18	0.69

TABLE 3: Summary of results on Fashion-MNIST dataset.

Compression ratio ( $r$ )	Algorithms	Performance				
		Bal_Acc	AUROC	Round	Cost(Megabyte)	$\epsilon$
0.05	FL-RND	0.60	0.69	99	4.73	N/A
	FL-FREQ	0.69	0.76	100	4.78	N/A
	FL-CS	0.73	0.80	100	4.78	N/A
	FL-RND-DP	0.60	0.69	100	4.78	1
	FL-FREQ-DP	0.65	0.72	100	4.78	1
	FL-CS-DP	0.69	0.76	100	4.78	1
0.1	FL-RND	0.66	0.73	100	9.56	N/A
	FL-FREQ	0.71	0.78	100	9.56	N/A
	FL-CS	0.73	0.81	87	8.31	N/A
	FL-RND-DP	0.65	0.72	100	9.56	1
	FL-FREQ-DP	0.67	0.74	100	9.56	1
	FL-CS-DP	0.69	0.76	99	9.46	1
0.2	FL-RND	0.69	0.76	100	19.11	N/A
	FL-FREQ	0.72	0.80	100	19.11	N/A
	FL-CS	0.73	0.81	74	14.14	N/A
	FL-RND-DP	0.67	0.74	99	18.92	1
	FL-FREQ-DP	0.69	0.76	100	19.11	1
	FL-CS-DP	0.68	0.74	64	12.23	0.92
1.0	FL-STD	0.74	0.82	99	94.62	N/A
	FL-STD-DP	0.70	0.77	93	88.88	0.99

TABLE 4: Summary of results on Medical dataset.

Existing works [60], [61] proposed to use a compressive sensing for federated learning in order to compress model updates without privacy guarantees. However, they assume that all clients participate in each round (as they maintain an error accumulation vector at each client due to the compression scheme), but as discussed in [62] this assumption is not always realistic. Recently in [63] another compressive sensing algorithm was proposed for federated learning for the denoising purpose (instead of the compression), where the added noise is due to the network transmission.

Sketching was adapted to federated learning for the purpose of compressing model updates in [37]. The authors proposed to use Count-Sketch [64] to retrieve the largest weights in the update vector on the server side. After that, the server uses two additional communication rounds to inform the clients about what gradient values they need to send back to the server. The server then takes the average of the received gradients and zeros-out the others before updating the model. The error due to the compression is maintained at each client, and the participation of all clients are required in each round

which, as per [62] and as discussed above, is not practical to federated learning. In [65], the aforementioned scheme is improved further by directly retrieving the most updated gradient values without asking for their positions in the update vector. This makes the scheme more efficient as it needs fewer communication rounds. Similarly to our approach, the error vector is also maintained on the server side instead of the client side, which is clearly a better fit for federated learning.

## 6. Conclusion

In this paper, we propose to extend Federated Learning with compressive sensing. Specifically, we propose two schemes: the first one (FL-CS) uses compressive sensing in order to reduce communication bandwidth. The second one (FL-CS-DP) combines compressive sensing and differential privacy in order to protect participants' information.

We present some experimental results that are based on the Fashion-MNIST dataset as well as on a medical dataset of 1.2 millions of US hospital patients. Results indicate

that using compressive sensing in Federated Learning allows to reduce the communication costs by up to 95% for a moderate loss of accuracy.

Results with the privacy-preserving extension FL-CS-DP indicate that compression happens to be especially useful and interesting in this context as it improves accuracy. This is due to the sensitivity reduction which is proportional to the added noise needed to guarantee differential privacy, and to the considered optimization problem (BPDN) which reconstructs data from noisy measurements.

We believe that the proposed privacy-preserving extension (FL-CS-DP) is an interesting alternative to differentially private federated learning (FL-STD-DP), as it improves both accuracy and bandwidth cost.

## References

- [1] Reza Shokri and Vitaly Shmatikov, "Privacy-preserving deep learning," in *ACM SIGSAC Conference on Computer and Communications Security, 2015*, 2015, pp. 1310–1321.
- [2] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *AISTATS*, 2016.
- [3] 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 *IEEE Symposium on Security and Privacy, 2019*, 2019, pp. 739–753.
- [4] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov, "Inference attacks against collaborative learning," *CoRR*, vol. abs/1805.04049, 2018.
- [5] Ligeng Zhu, Zhijian Liu, and Song Han, "Deep leakage from gradients," in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, Eds., 2019, pp. 14747–14756.
- [6] Cynthia Dwork and Aaron Roth, "The Algorithmic Foundations of Differential Privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3–4, 2014.
- [7] Keith Bonawitz et al., "Practical secure aggregation for federated learning on user-held data," *CoRR*, vol. abs/1611.04482, 2016.
- [8] Úlfar Erlingsson, Vasily Pihur, and Aleksandra Korolova, "RAP-POR: randomized aggregatable privacy-preserving ordinal response," in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, Scottsdale, AZ, USA, November 3-7, 2014*, Gail-Joon Ahn, Moti Yung, and Ninghui Li, Eds. 2014, pp. 1054–1067, ACM.
- [9] Stacey Truex and al., "A hybrid approach to privacy-preserving federated learning," *CoRR*, vol. abs/1812.03224, 2018.
- [10] Yuqing Zhu, Xiang Yu, Yi-Hsuan Tsai, Francesco Pittaluga, Masoud Faraki, Manmohan Chandraker, and Yu-Xiang Wang, "Voting-based approaches for differentially private federated learning," 2020.
- [11] David L Donoho, "Compressed sensing," *IEEE Transactions on information theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [12] Emmanuel J Candes and Terence Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?," *IEEE transactions on information theory*, vol. 52, no. 12, pp. 5406–5425, 2006.
- [13] Emmanuel J Candès, Justin Romberg, and Terence Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Transactions on information theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [14] François Chollet et al., "Keras early stopping," [https://keras.io/api/callbacks/early\\_stopping/](https://keras.io/api/callbacks/early_stopping/), 2015.
- [15] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang, "Deep learning with differential privacy," in *ACM CCS, New York, NY, USA, 2016*, pp. 308–318, ACM.
- [16] Laurent Jacques and Pierre Vandergheynst, "Compressed sensing: When sparsity meets sampling," Tech. Rep., Wiley-Blackwell, 2010.
- [17] Balas Kausik Natarajan, "Sparse approximate solutions to linear systems," *SIAM journal on computing*, vol. 24, no. 2, pp. 227–234, 1995.
- [18] Scott Shaobing Chen, David L Donoho, and Michael A Saunders, "Atomic decomposition by basis pursuit," *SIAM review*, vol. 43, no. 1, pp. 129–159, 2001.
- [19] Galen Andrew and Jianfeng Gao, "Scalable training of  $l_1$ -regularized log-linear models," in *Proceedings of the 24th international conference on Machine learning, 2007*, pp. 33–40.
- [20] Robert Taylor, "Orphant-wise limited-memory quasi-newton (owl-qn) algorithm implementation," <https://bitbucket.org/rtaylor/pylbfgs/src/master/>, 2020.
- [21] Emmanuel Candès, "The restricted isometry property and its implications for compressed sensing," *Compte Rendus de l'Academie des Sciences*, vol. 346, pp. 589–592, 05 2008.
- [22] Emmanuel J Candes and Terence Tao, "Decoding by linear programming," *IEEE transactions on information theory*, vol. 51, no. 12, pp. 4203–4215, 2005.
- [23] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U. Stich, and Martin Jaggi, "Error feedback fixes signsgd and other gradient compression schemes," *CoRR*, vol. abs/1901.09847, 2019.
- [24] Nasir Ahmed, T. Natarajan, and Kamisetty R Rao, "Discrete cosine transform," *IEEE transactions on Computers*, vol. 100, no. 1, pp. 90–93, 1974.
- [25] Nasir Ahmed, "How I came up with the discrete cosine transform," *Digit. Signal Process.*, vol. 1, no. 1, pp. 4–5, 1991.
- [26] K Ramamohan Rao and Ping Yip, *Discrete cosine transform: Algorithms, advantages, applications*, Academic press, 2014.
- [27] E. J. Candes and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?," *IEEE Transactions on Information Theory*, vol. 52, no. 12, pp. 5406–5425, 2006.
- [28] Gergely Ács and Claude Castelluccia, "I have a dream! (differentially private smart metering)," in *Information Hiding - 13th International Conference, IH 2011, 2011*, pp. 118–132.
- [29] Ilya Mironov, Kunal Talwar, and Li Zhang, "Rényi differential privacy of the sampled gaussian mechanism," *CoRR*, vol. abs/1908.10530, 2019.
- [30] A. Fejza, P. Genevès, N. Layaïda, and J. Bosson, "Scalable and interpretable predictive models for electronic health records," in *2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA)*, Oct 2018, pp. 341–350.
- [31] Anand Avati, Kenneth Jung, Stephanie Harman, Lance Downing, Andrew Ng, and Nigam H. Shah, "Improving palliative care with deep learning," *BMC Medical Informatics and Decision Making*, vol. 18, no. 4, pp. 122, Dec 2018.
- [32] Alvin Rajkomar and al., "Scalable and accurate deep learning with electronic health records," *npj Digital Medicine*, vol. 1, no. 1, pp. 18, 2018, url, An earlier version appeared in eprint arXiv:1801.07860.
- [33] Rupa Makadia and Patrick B. Ryan, "Transforming the premier perspective® hospital database into the observational medical outcomes partnership (omop) common data model," in *EGEMS*, 2014.
- [34] Margaret McDonald, Timothy Peng, Sridevi Sridharan, Janice Foust, Polina Kogan, Liliana Pezzin, and Penny Feldman, "Automating the medication regimen complexity index," *Journal of the American Medical Informatics Association : JAMIA*, vol. 20, 12 2012.
- [35] Ajinkya More, "Survey of resampling techniques for improving classification performance in unbalanced datasets," *arXiv preprint arXiv:1608.06048*, 2016.



Algos	Parameters
FL-STD & FL-STD-DP ( $r=1.0$ )	$S = 2.15; C = 1/60; N = 6000; T_{cl} = 200; T_{gd} = 5;  \mathbb{B}  = 10;  D_k  = 10; n = 1, 663, 370; \delta = 10^{-5}; SGD(\eta = 0.215); \sigma = 1.54$
FL-CS,FL-RND,FL-FREQ and their private extensions ( $r=0.2$ )	$S = 0.98; C = 1/60; N = 6000; T_{cl} = 200; T_{gd} = 5;  \mathbb{B}  = 10;  D_k  = 10; n = 1, 663, 370; \delta = 10^{-5}; SGD(\eta = 0.215); \eta_G = 0.35; \rho = 0.9; P = 200; \sigma = 1.54$
FL-CS,FL-RND,FL-FREQ and their private extensions ( $r=0.1$ )	$S = 0.69; C = 1/60; N = 6000; T_{cl} = 200; T_{gd} = 5;  \mathbb{B}  = 10;  D_k  = 10; n = 1, 663, 370; \delta = 10^{-5}; SGD(\eta = 0.215); \eta_G = 0.35; \rho = 0.9; P = 200; \sigma = 1.54$
FL-CS,FL-RND,FL-FREQ and their private extensions ( $r=0.05$ )	$S = 0.47; C = 1/60; N = 6000; T_{cl} = 200; T_{gd} = 5;  \mathbb{B}  = 10;  D_k  = 10; n = 1, 663, 370; \delta = 10^{-5}; SGD(\eta = 0.215); \eta_G = 0.35; \rho = 0.9; P = 200; \sigma = 1.54$

TABLE 5: Common environment between the schemes on Fashion-MNIST.  $\rho$ ,  $\eta_G$  and  $P$  are only used with FL-CS and FL-CS-DP.

Algos	Parameters
FL-STD & FL-STD-DP ( $r=1.0$ )	$S = 0.31; C = 100/5011; N = 5011; T_{cl} = 100; T_{gd} = 5; n = 1, 496, 601; \delta = 10^{-5}; SGD(\eta = 0.1); \sigma = 1.49$
FL-CS,FL-RND,FL-FREQ and their private extensions ( $r=0.2$ )	$S = 0.14; C = 100/5011; N = 5011; T_{cl} = 100; T_{gd} = 5; n = 1, 496, 601; \delta = 10^{-5}; SGD(\eta = 0.1); \eta_G = 1.0; \rho = 0.9; P = 200; \sigma = 1.49$
FL-CS,FL-RND,FL-FREQ and their private extensions ( $r=0.1$ )	$S = 0.1; C = 100/5011; N = 5011; T_{cl} = 100; T_{gd} = 5; n = 1, 496, 601; \delta = 10^{-5}; SGD(\eta = 0.1); \eta_G = 1.0; \rho = 0.9; P = 200; \sigma = 1.49$
FL-CS,FL-RND,FL-FREQ and their private extensions ( $r=0.05$ )	$S = 0.07; C = 100/5011; N = 5011; T_{cl} = 100; T_{gd} = 5; n = 1, 496, 601; \delta = 10^{-5}; SGD(\eta = 0.1); \eta_G = 1.0; \rho = 0.9; P = 200; \sigma = 1.49$

TABLE 6: Common environment between the schemes on the Medical dataset.  $\rho$ ,  $\eta_G$  and  $P$  are only used with FL-CS and FL-CS-DP.

- [36] Haibo He and Eduardo A Garcia, "Learning from imbalanced data," *IEEE Transactions on knowledge and data engineering*, vol. 21, no. 9, pp. 1263–1284, 2009.
- [37] Nikita Ivkin, Daniel Rothchild, Enayat Ullah, Ion Stoica, Raman Arora, et al., "Communication-efficient distributed sgd with sketching," in *Advances in Neural Information Processing Systems*, 2019, pp. 13144–13154.
- [38] Han Xiao, Kashif Rasul, and Roland Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *CoRR*, vol. abs/1708.07747, 2017.
- [39] François Chollet et al., "Keras datasets," <https://keras.io/datasets/>, 2015.
- [40] François Chollet et al., "Keras," <https://keras.io>, 2015.
- [41] Martín Abadi, et al., "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, Software available from tensorflow.org.
- [42] Travis E Oliphant, *A guide to NumPy*, vol. 1, Trelgol Publishing USA, 2006.
- [43] Sarang Narkhede, "Understanding auc - roc curve," <https://towardsdatascience.com/understanding-auc-roc-curve-68b2303cc9c5>, 2018.
- [44] Kay Henning Brodersen, Cheng Soon Ong, Klaas Enno Stephan, and Joachim M Buhmann, "The balanced accuracy and its posterior distribution," in *2010 20th International Conference on Pattern Recognition*. IEEE, 2010, pp. 3121–3124.
- [45] Mohamed Bekkar, Hassiba Djema, and T.A. Alitouche, "Evaluation measures for models assessment over imbalanced data sets," *Journal of Information Engineering and Applications*, vol. 3, pp. 27–38, 01 2013.
- [46] H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang, "Learning differentially private recurrent language models," in *International Conference on Learning Representations*, 2018.
- [47] Robin C. Geyer, Tassilo Klein, and Moin Nabi, "Differentially private federated learning: A client level perspective," *CoRR*, vol. abs/1712.07557, 2017.
- [48] Dan Alistarh, Jerry Li, Ryota Tomioka, and Milan Vojnovic, "QSGD: randomized quantization for communication-optimal stochastic gradient descent," *CoRR*, vol. abs/1610.02132, 2016.
- [49] Wei Wen and al., "Terngrad: Ternary gradients to reduce communication in distributed deep learning," *CoRR*, vol. abs/1705.07878, 2017.
- [50] Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary B. Charles, Dimitris S. Papailiopoulos, and Stephen Wright, "Atomo: Communication-efficient learning via atomic sparsification," in *NeurIPS*, 2018.
- [51] Jakub Konečný, H. Brendan McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon, "Federated learning: Strategies for improving communication efficiency," *CoRR*, vol. abs/1610.05492, 2016.
- [52] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Anima Anandkumar, "signsgd: compressed optimisation for non-convex problems," *CoRR*, vol. abs/1802.04434, 2018.
- [53] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and Bill Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," in *International Conference on Learning Representations, ICLR 2018*, 2018.
- [54] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu, "1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns," in *INTERSPEECH 2014*, 2014, pp. 1058–1062.
- [55] Hossein Mamaghanian, Nadia Khaled, David Atienza, and Pierre Vanderghenst, "Compressed sensing for real-time energy-efficient ecg compression on wireless body sensor nodes," *IEEE Transactions on Biomedical Engineering*, vol. 58, no. 9, pp. 2456–2466, 2011.
- [56] Christopher A Metzler, Arian Maleki, and Richard G Baraniuk, "From denoising to compressed sensing," *IEEE Transactions on Information Theory*, vol. 62, no. 9, pp. 5117–5144, 2016.
- [57] Amin Tavakoli and Ali Pourmohammad, "Image denoising based on compressed sensing," *International Journal of Computer Theory and Engineering*, vol. 4, no. 2, pp. 266, 2012.
- [58] Lean Yu, Yang Zhao, and Ling Tang, "A compressed sensing based ai learning paradigm for crude oil price forecasting," *Energy Economics*, vol. 46, pp. 236–245, 2014.
- [59] Yang D. Li, Zhenjie Zhang, Marianne Winslett, and Yin Yang, "Compressive mechanism: Utilizing sparse representation in differential privacy," in *Proceedings of the 10th Annual ACM Workshop on Privacy in the Electronic Society*, New York, NY, USA, 2011, WPES '11, p. 177–182, Association for Computing Machinery.
- [60] Mohammad Mohammadi Amiri and Deniz Gündüz, "Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air," *CoRR*, vol. abs/1901.00844, 2019.
- [61] Mohammad Mohammadi Amiri and Deniz Gündüz, "Federated learning over wireless fading channels," *CoRR*, vol. abs/1907.09769, 2019.

- [62] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al., “Advances and open problems in federated learning,” *arXiv preprint arXiv:1912.04977*, 2019.
- [63] Yo-Seb Jeon, Mohammad Mohammadi Amiri, Jun Li, and H. Vincent Poor, “A compressive sensing approach for federated learning over massive mimo communication systems,” 2020.
- [64] Moses Charikar, Kevin Chen, and Martin Farach-Colton, “Finding frequent items in data streams,” in *International Colloquium on Automata, Languages, and Programming*. Springer, 2002, pp. 693–703.
- [65] Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman, Joseph Gonzalez, and Raman Arora, “Fetchsgd: Communication-efficient federated learning with sketching,” 2020.
- [66] Josephine Akosa, “Predictive accuracy: a misleading performance measure for highly imbalanced data,” in *Proceedings of the SAS Global Forum*, 2017, pp. 2–5.

## Appendix A.

### Medical data: Data pre-processing & experimental setup details

This section describes the experimental setting which is used to evaluate the accuracy and the privacy of our proposals.

#### A.1. Preprocessing

- 1) **Features normalization:** we extract from the dataset the values of each feature represented in Table 1. For gender, we use one-hot encoding: Male, Female and Unknown. Similarly, for admission type we use 4 features: Emergency, Urgent, Trauma Center, and Unknown<sup>8</sup>. For drugs, we extract 24,419 features which correspond to the different drugs (name and dosage). A given patient receives only a few of the possible drugs served, resulting in a very sparse patient’s record. We use a MinMax normalization for age and MRCI in order to rescale the values of these features between 0 and 1 (using MinMaxScaler class of scikit-learn<sup>9</sup>). The labels that we consider are boolean: true means that the patient died during his hospital stay while false means she survived.
- 2) **Patients filtering:** We consider patient and drug information of the first day at the hospital so that we can make predictions 24 hours after admission (as commonly found in the literature [30], [32]). We filter out the pregnant and new-born patients because the medication types and admission services are not the same for these two categories of patients. Our model prediction is built without patients’ historical medical data. This has the advantage to require minimum patient’s information and to work for new patients.
- 3) **Hospitals filtering:** The dataset contains 415 hospitals for a total size of 1,271,733 records. We

8. <https://www.resdac.org/cms-data/variables/claim-inpatient-admission-type-code-ffs>

9. <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>

split randomly the dataset into disjoint training and testing data (80% and 20% respectively). The final dataset for testing contains 254,347 patients, with 7,882 deceased patients and 246,465 non-deceased patients (see Table 2).

Using Client-Level differential privacy requires to add more noise than Record-Level differential privacy, because the privacy purposes are not the same as detailed in Section 2. To reduce the noise (when  $\epsilon$  is fixed) and then improve the utility, we have to reduce the number of iterations or to reduce the sampling probability which are the parameters used to compute  $\epsilon$ . We therefore have two options to reduce the sampling probability:

- Reducing the number of clients selected at each round  $|\mathbb{K}|$ . However this option also decreases the amount of data, and hence have a negative impact on the utility. We therefore preferred to use the next option.
- Increasing the total number of clients  $N$ : we created more hospitals by splitting randomly the training data over 5011 “virtual” hospitals. We also, took care to have at least one in-hospital dead patient per hospital. Each hospital contains 203 patients except one which has 356 patients. We created 5011 hospitals in order to have approximately the same number of patients per hospital, each of them with some in-hospital dead patients. In practise, Client-Level differential privacy is more adapted to an environment with a large set of clients as explained in [46], [47].

#### A.2. Imbalanced data

The dataset of each hospital is imbalanced because the proportion of patients that leave the hospital alive is, fortunately, much larger than in-hospital dead patients. To deal with this well-known problem, we have decided to use downsampling technique [35], [36], a standard solution used for this purpose.<sup>10</sup>

#### A.3. Performance Metrics

We use the following metrics:

- *Balanced accuracy* [44], [45] is computed as  $1/2 \cdot (\frac{TP}{P} + \frac{TN}{N}) = \frac{TPR+TNR}{2}$  and is mainly used with imbalanced data. *True Positive Rate (TPR)* and *True Negative Rate (TNR)*:  $TPR = \frac{TP}{P}$  and  $TNR = \frac{TN}{N}$ , where  $P$  and  $N$  are the number of positive and negative instances, respectively, and  $TP$  and  $TN$  are the number of true positive and true negative instances. We note that traditional (“non-balanced”) accuracy metrics such as  $\frac{TP+TN}{P+N}$  can be misleading for very imbalanced data [66]:

10. We have also tested weighted loss function and oversampling techniques. But, we noticed experimentally that downsampling technique outperforms the other techniques for all the schemes.

in our dataset, the minority class has only 3% of all the training samples (see Table 2), which means that a biased (and totally useless) model always predicting the majority class would have a (non-balanced) accuracy of 97%.

- The *area under the ROC curve (AUROC)* is also a frequently used accuracy metric. The ROC curve is calculated by varying the prediction threshold from 1 to 0, when *TPR* and *FPR* are calculated at each threshold. The area under this curve is then used to measure the quality of the predictions. A random guess has an *AUROC* value of 0.5, whereas a perfect prediction has the largest *AUROC* value of 1.

#### A.4. Evaluation Method.

First, we split randomly the dataset of each hospital into disjoint training and testing data (80% and 20% respectively). An entire federated run is executed with this split, and all the metrics are evaluated in every round on the union of all clients' testing data. All metric values of the round with the best balanced metric are recorded.

**Algorithm 6:** FL-RND

---

```

1 Server:
2   Initialize common model  $w_0$ 
3   for  $t = 1$  to  $T_{cl}$  do
4     Generate a random seed  $\zeta$  Select  $\mathbb{K}$  clients
       uniformly at random
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\mathbf{y}_t^k = \text{Client}_k(\mathbf{w}_{t-1}, \zeta)$ 
7     end
8      $\mathbf{y}_t = \sum_k \frac{|D_k|}{\sum_j |D_j|} \Delta \mathbf{w}_t^k$ 
9      $j = 0$ 
10    for each element  $i$  in  $\mathbf{G}$  do
11       $\mathbf{w}_t[i] = \mathbf{w}_{t-1}[i] + \mathbf{y}_t[j]$ 
12     $j = j + 1$ 
13    end
14  end
Output: Global model  $\mathbf{w}_t$ 
15
16 Client $_k(\mathbf{w}_{t-1}^k, \zeta)$ :
17  $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{gd})$ 
18  $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$ 
19 Generates a random set  $\mathbf{G} = \{x \in \{1, \dots, n\}\}$  of
     $m$  random integer values such that  $m \leq n$  based
    on the seed  $\zeta$ 
20  $\hat{\Delta w}_t^k = \text{Sample } m \text{ elements from } \Delta \mathbf{w}_t^k \text{ by taking}$ 
    each element of  $\mathbf{G}$  as a coordinate
Output: The sampled Model update  $\hat{\Delta w}_t^k$ 

```

---

**Algorithm 7:** FL-FREQ

---

```

1 Server:
2   Initialize common model  $w_0$ 
3   for  $t = 1$  to  $T_{cl}$  do
4     Select  $\mathbb{K}$  clients uniformly at random
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\Delta \mathbf{y}_t^k = \text{Client}_k(\mathbf{w}_{t-1})$ 
7     end
8      $\mathbf{y}_t = \sum_k \frac{|D_k|}{\sum_j |D_j|} \Delta \mathbf{w}_t^k$ 
9      $\hat{\mathbf{y}}_t = \Phi^{-1} \mathbf{y}_t$  : Transform to time domain
10     $\mathbf{w}_t = \mathbf{w}_{t-1} + \hat{\mathbf{y}}_t$ 
11  end
Output: Global model  $\mathbf{w}_t$ 
12
13 Client $_k(\mathbf{w}_{t-1}^k)$ :
14  $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{gd})$ 
15  $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$ 
Output: The sampled Model update  $\Phi \Delta \mathbf{w}_t^k$ 

```

---

**Algorithm 8:** FL-RND-DP

---

```

1 Server:
2   Initialize common model  $w_0$ 
3   for  $t = 1$  to  $T_{cl}$  do
4     Generate a random seed  $\zeta$  Select  $\mathbb{K}$  clients
       uniformly at random
5     for each client  $k$  in  $\mathbb{K}$  do
6        $\mathbf{y}_t^k = \text{Client}_k(\mathbf{w}_{t-1}, \zeta)$ 
7     end
8      $\mathbf{y}_t = \sum_k \frac{|D_k|}{\sum_j |D_j|} \Delta \mathbf{w}_t^k$ 
9      $j = 0$ 
10    for each element  $i$  in  $\mathbf{G}$  do
11       $\mathbf{w}_t[i] = \mathbf{w}_{t-1}[i] + \mathbf{y}_t[j]$ 
12     $j = j + 1$ 
13    end
14  end
Output: Global model  $\mathbf{w}_t$ 
15
16 Client $_k(\mathbf{w}_{t-1}^k, \zeta)$ :
17  $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{gd})$ 
18  $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$ 
19 Generates a random set  $\mathbf{G} = \{x \in \{1, \dots, n\}\}$  of
     $m$  random integer values such that  $m \leq n$  based
    on the seed  $\zeta$ 
20  $\hat{\Delta w}_t^k = \text{Sample } m \text{ elements from } \Delta \mathbf{w}_t^k \text{ by taking}$ 
    each element of  $\mathbf{G}$  as a coordinate
21  $\hat{\Delta w}_t^{k'} = \hat{\Delta w}_t^k / \max\left(1, \frac{\|\hat{\Delta w}_t^k\|_2}{S}\right)$ 
Output:  $\text{Enc}_{K_k}(\mathcal{G}(\hat{\Delta w}_t^{k'}, S, I\sigma / \sqrt{|K|}))$ 

```

---

---

**Algorithm 9: FL-FREQ-DP**

---

1 **Server:**  
2 Initialize common model  $w_0$   
3 **for**  $t = 1$  **to**  $T_{\text{cl}}$  **do**  
4     Select  $\mathbb{K}$  clients uniformly at random  
5     **for each client**  $k$  **in**  $\mathbb{K}$  **do**  
6          $\Delta \mathbf{y}_t^k = \text{Client}_k(\mathbf{w}_{t-1})$   
7     **end**  
8      $\mathbf{y}_t = \sum_k \frac{|D_k|}{\sum_j |D_j|} \Delta \mathbf{w}_t^k$   
9      $\hat{\mathbf{y}}_t = \Phi^{-1} \mathbf{y}_t$  : Transform to time domain  
10      $\mathbf{w}_t = \mathbf{w}_{t-1} + \hat{\mathbf{y}}_t$   
11 **end**  
12 **Output:** Global model  $\mathbf{w}_t$   
13 **Client** $_k(\mathbf{w}_{t-1}^k)$ :  
14      $\mathbf{w}_t^k = \text{SGD}(D_k, \mathbf{w}_{t-1}^k, T_{\text{gd}})$   
15      $\Delta \mathbf{w}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$   
16      $\hat{\Delta w}_t^k = \Phi \Delta \mathbf{w}_t^k / \max\left(1, \frac{\|\mathcal{C}(\hat{\Delta \mathbf{w}}_t^k, m)\|_2}{S}\right)$   
17 **Output:**  $\text{Enc}_{K_k}(\mathcal{G}(\hat{\Delta w}_t^k, S\mathbf{I}\sigma/\sqrt{|K|}))$

---