

Federated Learning with Sparsified Model Perturbation: Improving Accuracy under Client-Level Differential Privacy

Rui Hu, *Member, IEEE*, Yuanxiong Guo, *Senior Member, IEEE* and Yanmin Gong, *Senior Member, IEEE*

Abstract—Federated learning (FL) that enables edge devices to collaboratively learn a shared model while keeping their training data locally has received great attention recently and can protect privacy in comparison with the traditional centralized learning paradigm. However, sensitive information about the training data can still be inferred from model parameters shared in FL. Differential privacy (DP) is the state-of-the-art technique to defend against those attacks. The key challenge to achieving DP in FL lies in the adverse impact of DP noise on model accuracy, particularly for deep learning models with large numbers of parameters. This paper develops a novel differentially-private FL scheme named Fed-SMP that provides a client-level DP guarantee while maintaining high model accuracy. To mitigate the impact of privacy protection on model accuracy, Fed-SMP leverages a new technique called Sparsified Model Perturbation (SMP) where local models are sparsified first before being perturbed by Gaussian noise. We provide a tight end-to-end privacy analysis for Fed-SMP using Rényi DP and prove the convergence of Fed-SMP with both unbiased and biased sparsifications. Extensive experiments on real-world datasets are conducted to demonstrate the effectiveness of Fed-SMP in improving model accuracy with the same DP guarantee and saving communication cost simultaneously.

Index Terms—Federated learning, edge computing, differential privacy, communication efficiency, sparsification.

1 INTRODUCTION

THE proliferation of edge devices such as smartphones and Internet-of-things (IoT) devices, each equipped with rich sensing, computation, and storage resources, leads to tremendous data being generated on a daily basis at the network edge. These data can be analyzed to build machine learning models that enable a wide range of intelligent services such as personal fitness tracking [1], traffic monitoring [2], and smart home security [3]. Traditional machine learning paradigm requires transferring all the raw data to the cloud before training a model, leading to high communication cost and severe privacy risk.

As an alternative machine learning paradigm, *Federated Learning (FL)* has attracted significant attention recently due to its benefits in communication efficiency and privacy [4]. In FL, multiple clients collaboratively learn a shared statistical model under the orchestration of the cloud without sharing their local data [5]. Although only model updates instead of raw data are shared by each client in FL, it is not sufficient to ensure privacy as the sensitive training data can still be inferred from the shared model parameters by using advanced inference attacks. For instance, given an input sample and a target model, the membership inference attack [6] can train an attack model to determine whether the sample was used for training the target model or not. Also, given a target model and class, the model inversion attack [7] can recover the typical representations of the target

class using an inversion model learned from the correlation between the inputs and outputs of the target model.

As a cryptography-inspired rigorous definition of privacy, differential privacy (DP) has become the de-facto standard for achieving data privacy and can give a strong privacy guarantee against an adversary with arbitrary auxiliary information [8]. For privacy protection in FL, client-level DP is often more relevant than record-level DP: client-level DP guarantee protects the participation of a client, while record-level DP guarantee protects only a single data sample of a client [9], [10]. While DP can be straightforwardly achieved using Gaussian or Laplacian mechanism [8], achieving client-level DP in the FL setting faces several major challenges in maintaining high model accuracy. Firstly, FL is an iterative learning process where model updates are exchanged in multiple rounds, leading to more privacy leakage compared to the one-shot inference. Secondly, the intensity of added DP noise is linearly proportional to the model size, which can be very large (e.g., millions of model parameters) for modern deep neural networks (DNNs), and will severely degrade the accuracy of the trained model. Thirdly, it is more challenging to achieve client-level DP than record-level DP because the entire dataset of a client rather than a single data sample of that client needs to be protected. Existing studies [11]–[13] on FL with client-level DP suffer from significant accuracy degradation due to the inherent challenge of large additive random noise required to achieve a certain level of client-level DP for DNNs.

In this paper, we propose a new differentially-private FL scheme called *Fed-SMP*, which guarantees client-level DP while preserving model accuracy and saving communication cost simultaneously. Fed-SMP utilizes *Sparsified Model*

- R. Hu is with the Department of Computer Science & Engineering, University of Nevada, Reno, Reno, NV, 89557 USA (E-mail: ruihu@unr.edu); Y. Guo is with the Department of Information Systems and Cyber Security, University of Texas at San Antonio, San Antonio, TX, 78249 USA (E-mail: yuanxiong.guo@utsa.edu); Y. Gong is with the Department of Electrical and Computer Engineering, University of Texas at San Antonio, San Antonio, TX, 78249 USA (E-mail: yanmin.gong@utsa.edu).

Perturbation (SMP) to improve the privacy-accuracy tradeoff of DP in FL. Specifically, SMP first sparsifies the local model update of a client at each round by selecting only a subset of coordinates to keep and then adds Gaussian noise to perturb the values at those selected coordinates. As we will show later in this paper, by using SMP in FL, the sparsification will have an amplification effect on the privacy guarantee offered by the added Gaussian noise, leading to a better privacy-accuracy tradeoff. Meanwhile, it can reduce the communication cost by compressing the shared model updates at each round. In summary, the main contributions of this paper are summarized as follows.

- We propose a new differentially private FL scheme called Fed-SMP, which achieves client-level DP guarantee for FL with a large number of clients. Fed-SMP only makes lightweight modifications to FedAvg [5], the most common learning method for FL, which enables easy integration into existing packages.
- Compared with DP-FedAvg [11], the state-of-art of differentially private FL schemes, Fed-SMP improves the privacy-utility tradeoff and communication efficiency of FL with client-level DP by using sparsification as a tool for amplifying privacy and reducing the communication cost at the same time.
- By integrating different sparsification operators with model perturbation, we design two algorithms of Fed-SMP: Fed-SMP with unbiased random sparsification and Fed-SMP with biased top- k sparsification. The resulting algorithms require less amount of added random noise to achieve the same level of DP and are compatible with secure aggregation, which is a crucial privacy-enhancing technique to achieve client-level DP in practical FL systems.
- To prove the (ϵ, δ) -DP guarantee of Fed-SMP, we use Rényi differential privacy (RDP) to tightly account the end-to-end privacy loss of Fed-SMP. We also theoretically analyze the impact of sparsified model perturbation on the convergence of Fed-SMP with both unbiased and biased sparsifications, filling the gap in the state-of-arts. The theoretical results indicate that under a certain DP guarantee, the optimal compression level of Fed-SMP needs to balance the increased compression error and reduced privacy error due to sparsification.
- We empirically evaluate the performances of Fed-SMP on both IID and non-IID datasets and compare the results with those of the state-of-art baselines. Experimental results demonstrate that Fed-SMP can achieve higher model accuracy than baseline approaches under the same level of DP and save the communication cost simultaneously.

The rest of the paper is organized as follows. Preliminaries on DP are described in Section 2. Section 3 introduces the problem formulation and presents the proposed Fed-SMP scheme. The privacy and convergence properties of Fed-SMP with random and top- k sparsification strategies are rigorously analyzed in Section 5. Section 6 shows the experimental results. Finally, Section 7 reviews the related work, and Section 8 concludes the paper.

2 PRELIMINARIES

DP has been proposed as a rigorous privacy notion for measuring privacy risk. The classic notion of DP, (ϵ, δ) -DP, is defined as follows:

Definition 1 ((ϵ, δ) -DP [8]). *Given privacy parameters $\epsilon > 0$ and $0 \leq \delta < 1$, a randomized mechanism \mathcal{M} satisfies (ϵ, δ) -DP if for any two adjacent datasets D, D' and any subset of outputs $O \subseteq \text{range}(\mathcal{M})$,*

$$\Pr[\mathcal{M}(D) \in O] \leq e^\epsilon \Pr[\mathcal{M}(D') \in O] + \delta. \quad (1)$$

When $\delta = 0$, we have ϵ -DP, or Pure DP.

To better calculate the privacy loss over multiple iterations in differentially private learning algorithms, Rényi differential privacy (RDP) has been proposed as follows:

Definition 2 ((α, ρ) -RDP [14]). *Given a real number $\alpha > 1$ and privacy parameter $\rho \geq 0$, a randomized mechanism \mathcal{M} satisfies (α, ρ) -RDP if for any two adjacent datasets D, D' , the Rényi α -divergence between $\mathcal{M}(D)$ and $\mathcal{M}(D')$ satisfies*

$$D_\alpha[\mathcal{M}(D) \parallel \mathcal{M}(D')] := \frac{1}{\alpha - 1} \log \mathbb{E} \left[\left(\frac{\mathcal{M}(D)}{\mathcal{M}(D')} \right)^\alpha \right] \leq \rho(\alpha).$$

RDP is a relaxed version of pure DP with a tighter composition bound. Thus, it is more suitable to analyze the end-to-end privacy loss of iterative algorithms. We can convert RDP to (ϵ, δ) -DP for any $\delta > 0$ using the following lemma:

Lemma 1 (From RDP to (ϵ, δ) -DP [15]). *If the randomized mechanism \mathcal{M} satisfies $(\alpha, \rho(\alpha))$ -RDP, then it also satisfies $(\rho(\alpha) + \frac{\log(1/\delta)}{\alpha-1}, \delta)$ -DP.*

In the following, we provide some useful definitions and lemmas about DP and RDP that will be used to derive our main results in the rest of the paper.

Definition 3 (ℓ_2 -sensitivity [8]). *Let $h : \mathcal{D} \rightarrow \mathbb{R}^d$ be a query function over a dataset. The ℓ_2 -sensitivity of h is defined as $\psi(h) := \sup_{D, D' \in \mathcal{D}, D \sim D'} \|h(D) - h(D')\|_2$ where $D \sim D'$ denotes that D and D' are two adjacent datasets.*

Lemma 2 (Gaussian Mechanism [14]). *Let $h : \mathcal{D} \rightarrow \mathbb{R}^d$ be a query function with ℓ_2 -sensitivity $\psi(h)$. The Gaussian mechanism $\mathcal{M} = h(D) + \mathcal{N}(0, \sigma^2 \psi(h)^2 \mathbf{I}_d)$ satisfies $(\alpha, \alpha/2\sigma^2)$ -RDP.*

Lemma 3 (RDP for Subsampling Mechanism [15], [16]). *For a Gaussian mechanism \mathcal{M} and any m -datapoints dataset D , define $\mathcal{M} \circ \text{SUBSAMPLE}$ as 1) subsample without replacement B datapoints from the dataset (denote $q = B/m$ as the sampling ratio); and 2) apply \mathcal{M} on the subsampled dataset as input. Then if \mathcal{M} satisfies $(\alpha, \rho(\alpha))$ -RDP with respect to the subsampled dataset for all integers $\alpha \geq 2$, then the new randomized mechanism $\mathcal{M} \circ \text{SUBSAMPLE}$ satisfies $(\alpha, \rho'(\alpha))$ -RDP w.r.t D , where*

$$\rho'(\alpha) \leq \frac{1}{\alpha - 1} \log \left(1 + q^2 \binom{\alpha}{2} \min\{4(e^{\rho(2)} - 1), 2e^{\rho(2)}\} + \sum_{j=3}^{\alpha} q^j \binom{\alpha}{j} 2e^{(j-1)\rho(j)} \right).$$

If $\sigma^2 \geq 0.7$ and $\alpha \leq (2/3)\sigma^2\psi^2(h) \log(1/q\alpha(1 + \sigma^2)) + 1$, $\mathcal{M} \circ \text{SUBSAMPLE}$ satisfies $(\alpha, 3.5q^2\alpha/\sigma^2)$ -RDP.

Lemma 4 (RDP Composition [14]). *For randomized mechanisms \mathcal{M}_1 and \mathcal{M}_2 applied on dataset D , if \mathcal{M}_1 satisfies (α, ρ_1) -RDP and \mathcal{M}_2 satisfies (α, ρ_2) -RDP, then their composition $\mathcal{M}_1 \circ \mathcal{M}_2$ satisfies $(\alpha, \rho_1 + \rho_2)$ -RDP.*

3 SYSTEM MODELING AND PROBLEM FORMULATION

In this section, we first present the problem formulation of FL and the attack model, then describe the classic method to solve the problem and its privacy-preserving variant.

3.1 Problem Formulation

A typical FL system consists of n clients (e.g., smartphones or IoT devices) and a central server (e.g., the cloud), as shown in Fig. 1. Each client $i \in [n]$ has a local dataset D_i , and those clients collaboratively train a global model $\theta \in \mathbb{R}^d$ based on their collective datasets under the orchestration of the central server. The goal of FL is to solve the following optimization problem:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) := \frac{1}{n} \sum_{i \in [n]} f_i(\theta), \quad (2)$$

where $f_i(\theta) = \mathbb{E}_{z \in D_i} [l(\theta; z)]$ is the local loss function of client i , z represents a datapoint sampled from D_i , and $l(\theta; z)$ denotes the loss of model θ on datapoint z . For $i \neq j$, the data distributions of D_i and D_j may be different.

3.2 Attack Model

The adversary can be the ‘‘honest-but-curious’’ server or clients in the system. The adversary will honestly follow the designed training protocol but is curious about a target client’s private data and wants to infer it from the shared messages. Furthermore, some clients can collude with the server or each other to infer private information about a specific victim client. Besides, the adversary could also be the passive outside attacker. These attackers can eavesdrop on all the shared messages during the execution of the training but will not actively inject false messages into or interrupt message transmissions.

3.3 Achieving Client-level DP in FL

As the classic and most widely-used algorithm in the FL setting, Federated Averaging (FedAvg) [5] solves (2) by running multiple iterations of stochastic gradient descent (SGD) in parallel on a subset of clients and then averaging the resulting local model updates at a central server periodically. Compared with distributed SGD, FedAvg is shown to achieve the same model accuracy with fewer communication rounds. Specifically, FedAvg involves T communication rounds, and each round can be divided into four stages as shown in Fig. 1. First, at the beginning of round $t \in \{0, \dots, T - 1\}$, the server randomly selects a subset \mathcal{W}^t of r clients and sends them the latest global model θ^t to perform local computations. Second, each client $i \in \mathcal{W}^t$ runs τ iterations of SGD on its local dataset to update its local model. Let $\theta_i^{t,s}$ denote client i ’s model after the s -th local iteration at the t -th communication round where

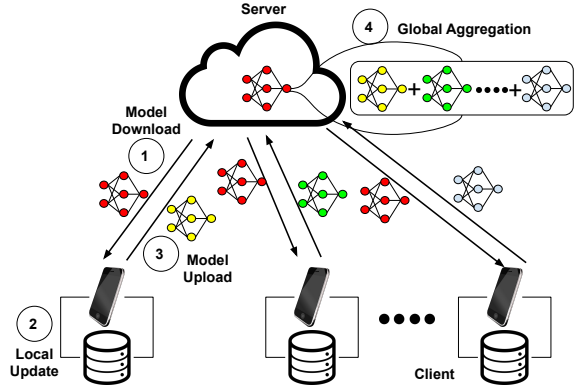


Fig. 1: An exemplary FL system.

$s \in [0, \tau - 1]$. By model initialization, we have $\theta_i^{t,0} = \theta^t$. Then the update rule at client i is represented as

$$\theta_i^{t,s+1} = \theta_i^{t,s} - \eta_t g_i^{t,s}, \forall s = 0, \dots, \tau - 1, \quad (3)$$

where η_t is the local learning rate and $g_i^{t,s} := (1/B) \sum_{z \in \xi_i^{t,s}} \nabla l(\theta_i^{t,s}, z)$ represents the stochastic gradient over a mini-batch $\xi_i^{t,s}$ of B datapoints sampled from D_i . Third, the client i uploads its model update $\Delta_i^t := \theta^t - \theta_i^{t,\tau}$ to the server. Fourth, after receiving all the local model updates, the server updates the global model by

$$\theta^{t+1} = \theta^t - \frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t. \quad (4)$$

The same procedure repeats at the next round $t+1$ thereafter until satisfying certain convergence criteria.

In FedAvg, clients repeatedly upload the local model updates of large dimension d (e.g., millions of model parameters for DNNs) to the server and download the newly-updated global model from the server many times in order to learn an accurate global model. Since the bandwidth between the server and clients could be rather limited, especially for uplink transmissions, the overall communication cost could be very high. Furthermore, although FedAvg can avoid direct information leakage by keeping the local dataset on the client, the model updates shared at each round can still leak private information about the local dataset, as demonstrated in recent advanced attacks such as model inversion attack [7] and membership inference attack [6].

As the commonly-used privacy protection technique for machine learning, DP mechanisms have been used to mitigate the above-mentioned privacy leakage in FedAvg. According to Section 2, the DP definitions apply to a range of different granularities, depending on how the adjacent datasets are defined. In this paper, we define the adjacent datasets by adding or removing the entire local data of a client and aim to protect whether one client participates in training or not in FL, resulting in the *client-level DP* [11]. In comparison, the commonly used privacy notion in standard non-federated learning DP is *record-level DP* [10], where the adjacent datasets are defined by adding or removing a single training example of a client, and only the privacy of one training example is protected. Therefore, client-level DP is stronger than record-level DP and has been shown to

be more relevant to the cross-device federated learning we considered, where there are a large number of participating devices, and each device can contribute multiple data records [11].

To provide client-level DP in FL without a fully trusted server, prior studies [9], [11] have proposed the model perturbation mechanism to prevent the privacy leakage from model updates in FedAvg, known as DP-FedAvg. As the pseudo-code of DP-FedAvg given in Algorithm 1, DP-FedAvg has the following changes to FedAvg at each FL round: 1) Each client’s model update is clipped to have a bounded ℓ_2 -norm (line 17); 2) Gaussian noise is added to the clipped model update at each client (line 18); 3) Final local model updates are encrypted following a secure aggregation protocol (e.g., [17], [18]) and sent to the server for aggregation (line 19).

Algorithm 1 DP-FedAvg

Require: number of selected clients per round r , number of training rounds T , local update period τ , local learning rate η_l , clipping threshold C , noise multiplier σ .

Server executes:

- 1: Initialize $\theta^0 \in \mathbb{R}^d$
- 2: **for** $t = 0$ to $T - 1$ **do**
- 3: Sample a set of r clients uniformly at random without replacement, denoted by $\mathcal{W}^t \subseteq [n]$
- 4: Broadcast θ^t to all clients in \mathcal{W}^t
- 5: **for** each clients $i \in \mathcal{W}^t$ **in parallel do**
- 6: $\mathbf{y}_i^t \leftarrow \text{ClientUpdate}(\theta^t)$
- 7: **end for**
- 8: $\theta^{t+1} \leftarrow \theta^t - (1/r) \sum_{i \in \mathcal{W}^t} \mathbf{y}_i^t$
- 9: **end for**
- 10: return θ^T

ClientUpdate(θ^t):

- 11: $\theta_i^{t,0} \leftarrow \theta^t$
 - 12: **for** $s = 0$ to $\tau - 1$ **do**
 - 13: Compute a mini-batch stochastic gradient $\mathbf{g}_i^{t,s}$
 - 14: $\theta_i^{t,s+1} \leftarrow \theta_i^{t,s} - \eta_l \mathbf{g}_i^{t,s}$
 - 15: **end for**
 - 16: $\hat{\Delta}_i^t \leftarrow \theta^t - \theta_i^{t,\tau}$
 - 17: $\bar{\Delta}_i^t \leftarrow \hat{\Delta}_i^t \times \min(1, C / \|\hat{\Delta}_i^t\|_2)$
 - 18: $\Delta_i^t \leftarrow \bar{\Delta}_i^t + \mathcal{N}(0, (C^2 \sigma^2 / r) \cdot \mathbf{I}_d)$
 - 19: Encrypt Δ_i^t and send it to the server via secure aggregation
-

Note that the use of secure aggregation is a common practice in the literature to achieve client-level DP in FL, and the design of a new secure aggregation protocol is out of the scope of this paper. Secure aggregation enables the server to learn just an aggregate function of the clients’ local model updates, typically the sum, and nothing else, so the Gaussian noise is added to prevent privacy leakage from the sum of local model updates. Specifically, in secure aggregation, clients generate randomly sampled zero-sum mask vectors locally by working in the space of integers modulo m and sampling the elements of the mask uniformly from \mathbb{Z}_m . When the server computes the modular sum of all the masked updates, the masks cancel out, and the server obtains the exact sum of local model updates. As with the existing works in FL [19]–[21], we ignore the finite precision

and modular summation arithmetic associated with secure aggregation in this paper, noting that one can follow the strategy in [22] to transform the real-valued vectors into integers for minimizing the approximation error of recovering the sum.

Although we can improve privacy and achieve client-level DP in FL by adding Gaussian noise locally and using secure aggregation, the resulting accuracy of the trained model is often low due to the significant intensity of added Gaussian noise. Moreover, communication cost has never been considered in those studies. This motivates us to develop a new differentially private FL scheme that can maintain high model accuracy while reducing the communication cost.

4 FED-SMP: FEDERATED LEARNING WITH SPARSIFIED MODEL PERTURBATION

In this section, we develop a new FL scheme called Fed-SMP to provide client-level DP with high model accuracy while improving communication efficiency at the same time. To ensure easy integration into existing packages/systems, Fed-SMP follows the same overall procedure of FedAvg as depicted in Fig. 1 and employs a novel integration of Gaussian mechanism and sparsification in the local update stage. This guarantees that each client’s shared local model update is both sparse and differentially private. The two specific sparsifiers considered in Fed-SMP are defined as follows:

Definition 4 (Random and Top- k Sparsifiers). *For a parameter $1 \leq k \leq d$ and vector $\mathbf{x} \in \mathbb{R}^d$, the random sparsifier $\text{rand}_k : \mathbb{R}^d \times \Omega_k \rightarrow \mathbb{R}^d$ and top- k sparsifier $\text{top}_k : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are defined as*

$$[\text{rand}_k(\mathbf{x}, \omega)]_j := \begin{cases} [\mathbf{x}]_j, & \text{if } j \in \omega \\ 0, & \text{otherwise} \end{cases}, \quad (5)$$

$$[\text{top}_k(\mathbf{x})]_j := \begin{cases} [\mathbf{x}]_{\pi(j)}, & \text{if } j \leq k \\ 0, & \text{otherwise} \end{cases}, \quad (6)$$

where $\Omega_k = \binom{[d]}{k}$ denotes the set of all k -element subsets of $[d]$, ω is chosen uniformly at random, i.e., $\omega \sim_{u.a.r} \Omega_k$, and π is a permutation of $[d]$ such that $|[\mathbf{x}]_{\pi(j)}| \geq |[\mathbf{x}]_{\pi(j+1)}|$ for $j \in [1, d - 1]$.

In Fed-SMP with both rand_k and top_k sparsifier, only k coordinates of the local model update will be sent out. Generally speaking, rand_k sparsifier may discard coordinates that are actually important, which inevitably degrades model accuracy when k is small. On the other hand, top_k sparsifier keeps the coordinates with the largest magnitude and can achieve higher model accuracy for small k . This seems to indicate that top_k is always the better choice. However, the two sparsifiers have different privacy implications. For rand_k sparsifier, since the set of selected coordinates ω is chosen uniformly at random, the coordinates themselves are not data-dependent and thus do not contain any private information of client data. On the other hand, the set of selected coordinates for top_k sparsifier depends on the values of model parameters and hence contains private information of client data and cannot be used like that

Algorithm 2 Fed-SMP: Federated Learning with Sparsified Model Perturbation

Require: number of selected clients per round r , number of training rounds T , local update period τ , local learning rate η_l , clipping threshold C , noise multiplier σ , compression parameter k , sparsifier spar (rand_k or top_k).

Server executes:

```

1: Initialize  $\theta^0 \in \mathbb{R}^d$ 
2: for  $t = 0$  to  $T - 1$  do
3:   Sample a set of  $r$  clients uniformly at random without replacement, denoted by  $\mathcal{W}^t \subseteq [n]$ 
4:   if  $\text{spar}$  is  $\text{rand}_k$  then
5:      $\mathbf{m}^t \leftarrow \text{SelectRandk}(\theta^t, k)$ 
6:   else if  $\text{spar}$  is  $\text{top}_k$  then
7:      $\mathbf{m}^t \leftarrow \text{SelectTopk}(\theta^t, k)$ 
8:   end if
9:   Broadcast  $\theta^t$  and  $\mathbf{m}^t$  to all clients in  $\mathcal{W}^t$ 
10:  for each clients  $i \in \mathcal{W}^t$  in parallel do
11:     $\mathbf{y}_i^t \leftarrow \text{ClientUpdate}(\theta^t, \mathbf{m}^t, \text{spar})$ 
12:  end for
13:   $\theta^{t+1} \leftarrow \theta^t - (1/r) \sum_{i \in \mathcal{W}^t} \mathbf{y}_i^t$ 
14: end for
15: return  $\theta^T$ 

```

SelectRandk(θ^t, k):

```

16: Select a random set of  $k$  coordinates of  $\theta^t$  and create a corresponding mask vector  $\mathbf{m}^t \in \{0, 1\}^d$ 
17: return  $\mathbf{m}^t$ 

```

SelectTopk(θ^t, k):

of rand_k sparsifier. It is worth noting that although rand_k sparsifier is a randomized mechanism, it does not provide any privacy guarantee by itself in terms of DP and needs to be combined with Gaussian mechanism for rigorous DP guarantee.

Using the above sparsifiers for client-level DP in FL has new challenges. As mentioned before, secure aggregation is a key privacy-enhancing technique to achieve client-level DP in practical FL systems. However, the naive application of rand_k or top_k sparsifier to the local model update of each client typically results in a different set of k coordinates for each client, preventing us from only encrypting the k selected coordinates of each client in the secure aggregation protocol. One can apply secure aggregation to all the d coordinates, but the communication efficiency benefit of sparsification will get lost. In the following, we design new rand_k and top_k sparsifiers in Fed-SMP that are compatible with secure aggregation. Specifically, we let the selected clients keep the same set of k active coordinates at each round; therefore, those clients can transmit the sparsified model update directly using the secure aggregation protocol and save the communication cost.

The pseudo-code for the proposed Fed-SMP is provided in Algorithm 2. At the beginning of round t , the server randomly selects a set \mathcal{W}^t of r clients and broadcasts to them the current global model θ^t and mask vector $\mathbf{m}^t \in \{0, 1\}^d$ (lines 3-9). The j -th coordinate of mask vector \mathbf{m}^t equals

Require: public dataset D_p , update period τ_p

```

18:  $\theta_p^0 \leftarrow \theta^t$ 
19: for  $s = 0$  to  $\tau - 1$  do
20:   Compute a mini-batch stochastic gradient  $\mathbf{g}_p^s$  over  $D_p$ 
21:    $\theta_p^{s+1} \leftarrow \theta_p^s - \eta_l \mathbf{g}_p^s$ 
22: end for
23:  $\Delta_p \leftarrow \theta^t - \theta_p^\tau$ 
24: Select the top  $k$  coordinates of  $|\Delta_p|$  and create a corresponding mask vector  $\mathbf{m}^t \in \{0, 1\}^d$ 
25: return  $\mathbf{m}^t$ 

```

ClientUpdate($\theta^t, \mathbf{m}^t, \text{spar}$):

```

26:  $\theta_i^{t,0} \leftarrow \theta^t$ 
27: for  $s = 0$  to  $\tau - 1$  do
28:   Compute a mini-batch stochastic gradient  $\mathbf{g}_i^{t,s}$ 
29:    $\theta_i^{t,s+1} \leftarrow \theta_i^{t,s} - \eta_l \mathbf{g}_i^{t,s}$ 
30: end for
31:  $\Delta_i^t \leftarrow \text{DP-spar}(\theta^t - \theta_i^{t,\tau}, \mathbf{m}^t, \text{spar})$ 
32: Encrypt  $\Delta_i^t$  and send it to the server via secure aggregation
    tion

```

DP-spar($\Delta, \mathbf{m}^t, \text{spar}$):

```

33: if  $\text{spar}$  is  $\text{rand}_k$  then
34:    $\Delta' \leftarrow \frac{d}{k} \times \Delta \odot \mathbf{m}^t$ 
35: else if  $\text{spar}$  is  $\text{top}_k$  then
36:    $\Delta' \leftarrow \Delta \odot \mathbf{m}^t$ 
37: end if
38:  $\hat{\Delta} \leftarrow \Delta' \times \min(1, C/\|\Delta'\|_2)$ 
39: return  $\hat{\Delta} + (\mathcal{N}(0, (C^2\sigma^2/r) \cdot \mathbf{I}_d) \odot \mathbf{m}^t)$ 

```

to 1 if that coordinate is selected by the sparsifier at round t and 0 otherwise. In Fed-SMP with rand_k sparsifier, the k coordinates are selected randomly from $[d]$ by the server (i.e., the procedure **SelectRandk**(\cdot)). In Fed-SMP with top_k sparsifier, we let the server choose a set of top k coordinates for all clients using a small public dataset D_p at each round, avoiding privacy leakage from the selected coordinates. The distribution of the public dataset D_p is assumed to be similar to the overall dataset distribution of clients. This assumption is common in the FL literature [23]–[26], where the server is assumed to have a small public dataset for model validation that mimics the overall dataset distribution. Specifically, at round t , the server first performs multiple iterations of SGD similar to (3) on the global model θ^t using the public dataset D_p and obtains the model difference Δ_p . Then, the server selects a set of k coordinates with the largest absolute values in Δ_p and generates a corresponding mask vector \mathbf{m}^t (i.e., the procedure **SelectTopk**(\cdot)).

After receiving the global model θ^t and mask vector \mathbf{m}^t from the server, each client $i \in \mathcal{W}^t$ initializes its local model to θ^t and runs τ iterations of SGD to update its local model in parallel (lines 26-30). Then, the client i sparsifies its local model update $\theta^t - \theta_i^{t,\tau}$ using the mask vector \mathbf{m}^t (lines 33-37). The operator \odot in Algorithm 2 represents the element-wise multiplication. Note that for rand_k sparsifier, the model update will be scaled by d/k to ensure an unbiased estimation on the sparsified model update (line 34). Since there is

no a priori bound on the size of the sparsified model update Δ' , each client will clip its sparsified model update in ℓ_2 -norm with clipping threshold C so that $\|\hat{\Delta}\|_2 \leq C$ (line 38). Next, each client perturbs its sparsified model update by adding independent Gaussian noise $\mathcal{N}(0, C^2\sigma^2/r)$ on the k selected coordinates (line 39), where σ represents the noise multiplier. The noisy sparsified model update Δ_i^t is then encrypted as an input into a secure aggregation protocol and sent to the server (line 32). Finally, the server computes the modular sum of all encrypted models to obtain the exact sum of local model updates (i.e., $\sum_{i \in \mathcal{W}^t} \mathbf{y}_i^t = \sum_{i \in \mathcal{W}^t} \Delta_i^t$) and updates the global model for the next round (line 13).

5 MAIN THEORETICAL RESULTS

In this section, we analyze the end-to-end privacy guarantee and convergence results of Fed-SMP. For better readability, we state the main theorems and only give the proof sketches in this section while leaving the complete proofs in the appendix. Before presenting our theoretical results, we make the following assumptions:

Assumption 1 (Smoothness). *Each local objective function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth, i.e., for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,*

$$\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L\|\mathbf{y} - \mathbf{x}\|, \forall i \in [n].$$

Assumption 2 (Unbiased Gradient and Bounded Variance). *The stochastic gradient at each client is an unbiased estimator of the local gradient: $\mathbb{E}[\mathbf{g}_i(\mathbf{x})] = \nabla f_i(\mathbf{x})$, and has bounded variance: $\mathbb{E}[\|\mathbf{g}_i(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2] \leq \zeta_i^2, \forall \mathbf{x} \in \mathbb{R}^d, i \in [n]$, where the expectation is over all the local mini-batches. We also denote $\bar{\zeta}^2 := (1/n) \sum_{i=1}^n \zeta_i^2$ for convenience.*

Assumption 3 (Bounded Dissimilarity). *There exist constants $\beta^2 \geq 1, \kappa^2 \geq 0$ such that $(1/n) \sum_{i=1}^n \|\nabla f_i(\mathbf{x})\|^2 \leq \beta^2 \|(1/n) \sum_{i=1}^n \nabla f_i(\mathbf{x})\|^2 + \kappa^2$. If local objective functions are identical to each other, then we have $\beta^2 = 1$ and $\kappa^2 = 0$.*

Assumptions 1 and 2 are standard in the analysis of SGD [27], and Assumption 3 is commonly used in the federated optimization literature [28], [29] to capture the dissimilarities of local objectives under non-IID data distribution.

5.1 Privacy Analysis

In this subsection, we provide the end-to-end privacy analysis of Fed-SMP based on RDP. Given the fact that the server only knows the sum of local model updates $\sum_{i \in \mathcal{W}^t} \Delta_i^t$ due to the use of secure aggregation, we need to compute the privacy loss incurred from releasing the sum of local model updates. Assume the client sets \mathcal{W} and \mathcal{W}' differ in one client index c such that $\mathcal{W}' := \mathcal{W} \cup \{c\}$. For any adjacent datasets $D := \{D_i\}_{i \in \mathcal{W}}$ and $D' := \{D_j\}_{j \in \mathcal{W}'} = \{D_i\}_{i \in \mathcal{W}} \cup D_c$, according to Definition 3, the ℓ_2 -sensitivity of the sum of local model updates is

$$\psi_{\Delta} := \sup_{D, D'} \left\| \sum_{i \in \mathcal{W}} \Delta_i^t(D_i) - \sum_{j \in \mathcal{W}'} \Delta_j^t(D_j) \right\|_2.$$

Due to the clipping, we have $\|\Delta_i^t(D_i)\|_2 \leq C, \forall i \in [n]$, and therefore $\psi_{\Delta} = \sup_{D, D'} \|\Delta_c^t(D_c)\|_2 \leq C$. As the sum of Gaussian random variables is still a Gaussian random variable, the variance of the Gaussian noise added to each

selected coordinate of the sum of local model updates is $C^2\sigma^2$. According to Lemmas 1–4, we compute the overall privacy guarantee of Fed-SMP as follows:

Theorem 1 (Privacy Guarantee of Fed-SMP). *Suppose the client is sampled without replacement with probability $q := r/n$ at each round. For any $\epsilon < 2 \log(1/\delta)$ and $\delta \in (0, 1)$, Fed-SMP satisfies (ϵ, δ) -DP after T communication rounds if*

$$\sigma^2 \geq \frac{7q^2T(\epsilon + 2 \log(1/\delta))}{\epsilon^2}.$$

Proof. After adding noise $\mathcal{N}(0, C^2\sigma^2)$ to the k selected coordinates, the sum of local model updates satisfies $(\alpha, \alpha/2\sigma^2)$ -RDP for each client in \mathcal{W}^t at round t by Lemma 2. As r clients are uniformly sampled from all clients at each round, the per-round privacy guarantee of Fed-SMP can be further amplified according to the subsampling amplification property of RDP in Lemma 3. The overall privacy guarantee of Fed-SMP follows by using the composition property in Lemma 4 to compute the RDP guarantee after T rounds and Lemma 1 to convert RDP to (ϵ, δ) -DP. The details are given in Appendix A. \square

5.2 Convergence Analysis

In this subsection, we provide the convergence result of Fed-SMP under the general non-convex setting in Theorem 2.

Theorem 2 (Convergence result of Fed-SMP). *Under Assumptions 1–3, assume the local learning rate satisfies $\eta \leq \min\{(1/24\tau L(\phi_k + 1)\beta^2, 1/4\tau L\sqrt{4\beta^2 + 2}, 1/12\tau L\}$ and $\|\hat{\Delta}\|_2 \leq C$, then the sequence of outputs θ^t generated by Algorithm 2 satisfies:*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\theta^t)\|^2 &\leq \frac{8e_0}{T\eta\tau} + 8\eta\tau L(3\kappa^2 + 2\bar{\zeta}^2) \\ &\quad + (8\eta\tau L(2\kappa^2 + \bar{\zeta}^2) + \zeta')\phi_k + \frac{4LkC^2\sigma^2}{\eta\tau r^2}, \end{aligned} \quad (7)$$

where $e_0 := f(\theta^0) - f^*$, $\phi_k := 1 - k/d$ and $\zeta' := 4\bar{\zeta}^2/\gamma$ for top_k sparsifier, $\phi_k := d/k - 1$ and $\zeta' := 0$ for rand_k sparsifier, and f^* represents the optimal objective value.

Proof. The proof is given in Appendix B-C. \square

The convergence bound (7) contains three parts. The first two terms $8e_0/T\eta\tau + 8\eta\tau L(3\kappa^2 + 2\bar{\zeta}^2)$ represent the optimization error bound in FedAvg. The third term $(8\eta\tau L(2\kappa^2 + \bar{\zeta}^2) + \zeta')\phi_k$ is the *compression error* resulted from applying sparsification on local model updates. The last term $4LkC^2\sigma^2/\eta\tau r^2$ is the *privacy error* resulted from adding DP noise to perturb local model updates. Both compression error and privacy error increase the error floor at convergence. When no sparsification is applied (i.e., $k = d$ and $\phi_k = 0$), the compression error is zero. When no DP noise is added (i.e., $\sigma = 0$), the privacy error is zero. The above result shows an explicit tradeoff between compression error and privacy error in Fed-SMP. As k decreases, the variance of sparsification ϕ_k gets larger which leads to a larger compression error, but the privacy error decreases. Therefore, there exists an optimal parameter k that can balance those two errors to minimize the convergence bound.

Compression ratio	Algorithm	Performance		
		Accuracy (%)	Cost (MB)	Privacy (ϵ)
$p = 0.001$	Fed-SMP-top $_k$	77.18 \pm 0.48	0.02	1.01
	Fed-SMP-rand $_k$	27.16 \pm 10.20	0.02	1.01
$p = 0.005$	Fed-SMP-top $_k$	80.76 \pm 0.27	0.10	1.01
	Fed-SMP-rand $_k$	56.04 \pm 5.22	0.10	1.01
$p = 0.01$	Fed-SMP-top $_k$	80.76 \pm 0.47	0.20	1.01
	Fed-SMP-rand $_k$	65.61 \pm 1.04	0.20	1.01
$p = 0.1$	Fed-SMP-top $_k$	80.49 \pm 0.23	2.00	1.01
	Fed-SMP-rand $_k$	77.59 \pm 0.53	2.00	1.01
$p = 0.2$	Fed-SMP-top $_k$	80.16 \pm 0.09	3.99	1.01
	Fed-SMP-rand $_k$	79.12 \pm 0.40	3.99	1.01
$p = 0.4$	Fed-SMP-top $_k$	80.26 \pm 0.25	7.98	1.01
	Fed-SMP-rand $_k$	79.88 \pm 0.37	7.98	1.01
$p = 0.8$	Fed-SMP-top $_k$	79.77 \pm 1.00	15.97	1.01
	Fed-SMP-rand $_k$	79.82 \pm 0.35	15.97	1.01
$p = 1.0$	DP-FedAvg	72.72 \pm 3.72	19.96	1.01
	FedAvg	86.98 \pm 0.12	19.96	-

TABLE 1: Summary of results on Fashion-MNIST dataset.

Compression ratio	Algorithm	Performance		
		Accuracy (%)	Cost (MB)	Privacy (ϵ)
$p = 0.005$	Fed-SMP-top $_k$	80.94 \pm 0.85	0.21	1.01
	Fed-SMP-rand $_k$	19.45 \pm 0.42	0.21	1.01
$p = 0.01$	Fed-SMP-top $_k$	81.12 \pm 1.08	0.41	1.01
	Fed-SMP-rand $_k$	19.75 \pm 0.12	0.41	1.01
$p = 0.05$	Fed-SMP-top $_k$	80.22 \pm 1.02	2.06	1.01
	Fed-SMP-rand $_k$	28.15 \pm 2.46	2.06	1.01
$p = 0.1$	Fed-SMP-top $_k$	79.64 \pm 0.29	4.12	1.01
	Fed-SMP-rand $_k$	64.50 \pm 2.95	4.12	1.01
$p = 0.2$	Fed-SMP-top $_k$	77.91 \pm 0.37	8.24	1.01
	Fed-SMP-rand $_k$	76.94 \pm 0.58	8.24	1.01
$p = 0.4$	Fed-SMP-top $_k$	75.58 \pm 0.40	16.47	1.01
	Fed-SMP-rand $_k$	77.07 \pm 0.48	16.47	1.01
$p = 0.8$	Fed-SMP-top $_k$	73.18 \pm 0.69	32.94	1.01
	Fed-SMP-rand $_k$	73.90 \pm 0.37	32.94	1.01
$p = 1.0$	DP-FedAvg	72.46 \pm 0.54	41.18	1.01
	FedAvg	88.47 \pm 0.17	41.18	-

TABLE 2: Summary of results on SVHN dataset.

6 PERFORMANCE EVALUATION

In this section, we evaluate the performance of Fed-SMP with rand $_k$ and top $_k$ sparsifiers (denoted by Fed-SMP-rand $_k$ and Fed-SMP-top $_k$, respectively) by comparing it with the following baselines:

- FedAvg: the classic FL algorithm served as the baseline without privacy consideration;
- FedAvg-top $_k$: a communication-efficient variant of FedAvg, where the local model update of each client is compressed using top $_k$ sparsifier before being uploaded;
- FedAvg-rand $_k$: another communication-efficient variant of FedAvg, where the local model update of each client is compressed using rand $_k$ sparsifier before being uploaded;
- DP-FedAvg [11]: the state-of-art differentially-private variant of FedAvg that achieves client-level DP, where the full-precision local model update from each client is clipped with clipping threshold C and then perturbed by adding Gaussian noise drawn from the distribution $\mathcal{N}(0, (C^2\sigma^2/r) \cdot \mathbf{I}_d)$, where σ is

the noise multiplier and r is the number of selected clients per round.

For fair comparisons, all the algorithms use the same secure aggregation protocol. When $p = 1.0$, Fed-SMP-top $_k$ and Fed-SMP-rand $_k$ reduce to DP-FedAvg, and FedAvg-top $_k$ and FedAvg-rand $_k$ reduce to FedAvg.

6.1 Dataset and Model

We use Fashion-MNIST [30] dataset, SVHN dataset [31] and Shakespeare dataset [5], three common benchmarks for differentially private machine/federated learning. The first two are image classification datasets, and the last one is a text dataset for next-character-prediction. Note that while Fashion-MNIST and SVHN datasets are considered as “solved” in the computer vision community, achieving high accuracy with strong privacy guarantee remains difficult on these datasets [32].

The Fashion-MNIST dataset consists of a training set of 60,000 examples and a test set of 10,000 examples. Each example is a 28×28 grayscale image associated with a label from 10 classes. For Fed-SMP-top $_k$, we randomly sample 1,000 examples from the training set as the public dataset D_p and evenly partition the remaining 59,000 examples across 6,000 clients. For other algorithms, we evenly split the training data across 6,000 clients. We use the convolutional neural network (CNN) model in [5] on Fashion-MNIST that consists of two 5×5 convolution layers (the first with 32 filters, the second with 64 filters, each followed with 2×2 max pooling and ReLu activation), a fully connected layer with 512 units and ReLu activation, and a final softmax output later, which results in about 1.6 million total parameters.

The SVHN dataset contains 73,257 training examples and 26,032 testing examples, and each of them is a 32×32 colored image of digits from 0 to 9. We extend the training set with another 531,131 additional examples. For all algorithms except Fed-SMP-top $_k$, we evenly partition the training set across 6,000 clients. For Fed-SMP-top $_k$, we randomly sample 2,000 examples from the training set as the public dataset D_p and evenly partition the remaining across 6,000 clients. We train a CNN model used in [33] on SVHN dataset, which stacks two 5×5 convolution layers (the first with 64 filters, the second with 128 filters, each followed with 2×2 max pooling and ReLu activation), two fully connected layers (the first with 384 units, the second with 192 units, each followed with ReLu activation), and a final softmax output later, resulting in about 3.4 million parameters in total.

The Shakespeare dataset is a natural federated dataset built from *The Complete Works of William Shakespeare*, where each of the total 715 clients corresponds to a speaking role with at least two lines. The training set and test set are obtained by splitting the lines from each client, and each client has at least one line for training or testing. We refer the reader to [34] for more details on processing the raw data. Note that for Fed-SMP-top $_k$, we randomly select 1000 samples in total from the processed training set of clients as the public dataset D_p . For Shakespeare, we use the recurrent neural network (RNN) model in [34] to predict the next word based on the preceding words in a line. The RNN model takes an 80-character sequence as input and consists

of a 80×8 embedding layers and 2 LSTM layers (each with 256 nodes) followed by a dense layer, resulting in about 0.8 million parameters in total.

Compression ratio	Algorithm	Performance		
		Accuracy (%)	Cost (MB)	Privacy (ϵ)
$p = 0.005$	Fed-SMP-top $_k$	55.41 \pm 0.20	0.23	6.99
	Fed-SMP-rand $_k$	36.41 \pm 0.15	0.23	6.99
$p = 0.01$	Fed-SMP-top $_k$	56.04 \pm 0.23	0.46	6.99
	Fed-SMP-rand $_k$	38.99 \pm 1.79	0.46	6.99
$p = 0.05$	Fed-SMP-top $_k$	56.76 \pm 0.13	2.30	6.99
	Fed-SMP-rand $_k$	50.79 \pm 0.27	2.30	6.99
$p = 0.1$	Fed-SMP-top $_k$	55.74 \pm 0.06	4.60	6.99
	Fed-SMP-rand $_k$	53.15 \pm 0.14	4.60	6.99
$p = 0.2$	Fed-SMP-top $_k$	55.89 \pm 0.14	9.20	6.99
	Fed-SMP-rand $_k$	54.22 \pm 0.15	9.20	6.99
$p = 0.4$	Fed-SMP-top $_k$	54.36 \pm 0.22	18.41	6.99
	Fed-SMP-rand $_k$	53.68 \pm 0.14	18.41	6.99
$p = 0.8$	Fed-SMP-top $_k$	52.14 \pm 0.31	36.82	6.99
	Fed-SMP-rand $_k$	51.85 \pm 0.29	36.82	6.99
$p = 1.0$	DP-FedAvg	51.03 \pm 0.28	46.02	6.99
	FedAvg	62.77 \pm 0.05	46.02	-

TABLE 3: Summary of results on Shakespeare dataset.

6.2 Experimental Settings

For Fashion-MNIST and SVHN datasets, the server randomly samples $r = 100$ clients to participate in the training at each round, and for Shakespeare, $r = 10$. For the local optimizer on clients, we use momentum SGD for both datasets. Specifically, for Fashion-MNIST, we set the local momentum coefficient as 0.5, local learning rate as $\eta_l = 0.125$ decaying at a rate of 0.99 at each communication round, batch size $B = 10$ and local update period to be 10 epochs; for SVHN, we set the local momentum coefficient as 0.8, local learning rate as $\eta_l = 0.05$ decaying at a rate of 0.99 at each communication round, batch size $B = 50$, local update period to be 5 epochs and number of communication rounds $T = 180$; for Shakespeare, we set the local momentum coefficient as 0.9, local learning rate as $\eta_l = 1$ decaying at a rate of 0.99 at every 50 round, batch size $B = 4$ and local update period to be 1 epoch and number of rounds $T = 1000$. It is worth noting that the momentum at each client is initialized at the beginning of every communication round, since local momentum will be stale due to the partial participation of clients. Specially, for Fed-SMP-top $_k$, the server uses the same optimizer as the clients with the same local update period to compute the global model update Δ_p .

For the privacy-preserving algorithms, we compute the end-to-end privacy loss using the API provided in [35]. For all datasets, we follow [11] to set $\delta = 1/n^{1.1}$ such that δ is less than $1/n$. In particular, as it is challenging to obtain a small ϵ on Shakespeare dataset due to the small n and the minimal r sufficient for convergence, we follow [36] and [11] to report privacy with a hypothetical population size $n = 10^9$ for Shakespeare dataset. Moreover, the noise multiplier is set as $\sigma = 1.4$ for Fashion-MNIST and SVHN datasets and $\sigma = 0.3$ for Shakespeare dataset by default. We tune the clipping threshold over the grid $C = \{0.1, 0.2, 0.4, 0.6, 0.8, 1.0, 2.0\}$ and obtain the optimal clipping threshold, i.e., $C = 1.0$ for Fashion-MNIST dataset, $C = 0.4$ for SVHN and Shakespeare datasets.

In addition, we define the *compression ratio* for rand $_k$ and top $_k$ sparsifiers as $p := k/d$. When p is smaller, more coordinates of the local model updates are set to zero, leading to a higher level of communication compression.

6.3 Experiment Results

Table 1, Table 2 and Table 3 summarize the performance of FedAvg, DP-FedAvg and Fed-SMP after T rounds on Fashion-MNIST, SVHN and Shakespeare, respectively. We run each experiment with 5 random seeds and report the average and standard deviation of testing *accuracy*. Note that we always report the best testing accuracy across all rounds in each experiment. We also report the uplink communication *cost* that denotes the size of the data sent from a client to the server as the input into the secure aggregation protocol, which equals to $32k \times T \times (r/n)$ bits for Fed-SMP and $32d \times T \times (r/n)$ bits for DP-FedAvg and FedAvg. The *privacy* for DP-FedAvg and Fed-SMP is the accumulated privacy loss ϵ at the end of the training.

From Table 1, we can see that the non-private FedAvg achieves a testing accuracy of 86.61% on Fashion-MNIST dataset. When adding Gaussian noise to ensure $(1.01, n^{-1.1})$ -DP, the accuracy drops to 71.40% for DP-FedAvg. By choosing a proper value of p , Fed-SMP can reach a higher accuracy than DP-FedAvg under the same level of DP guarantee. For example, the highest testing accuracy of Fed-SMP-top $_k$ is 80.76% when $p = 0.005$, outperforming DP-FedAvg by 11%; the highest testing accuracy of Fed-SMP-rand $_k$ is 79.88% when $p = 0.4$, outperforming DP-FedAvg by 10%.

From Table 2, we can see that the accuracy of non-private FedAvg on SVHN dataset is 88.47%, and then it decreases to 72.46% to achieve $(1.01, n^{-1.1})$ -DP guarantee for DP-FedAvg. For Fed-SMP with the same level of DP guarantee, Fed-SMP-top $_k$ can achieve a testing accuracy of 81.12% on SVHN dataset when $p = 0.01$, outperforming DP-FedAvg by 12%, and Fed-SMP-rand $_k$ can achieve a testing accuracy of 77.07% on SVHN dataset when $p = 0.4$, outperforming DP-FedAvg by 6%.

From the results of Shakespeare dataset in Table 3, we can see that the accuracy of non-private FedAvg is 62.77% and then decreases to 51.03% to achieve $(6.99, n^{-1.1})$ -DP guarantee for DP-FedAvg. For Fed-SMP with the same level of DP guarantee, Fed-SMP-top $_k$ achieves a testing accuracy of 56.76% on Shakespeare when $p = 0.05$, outperforming DP-FedAvg by 11%, and Fed-SMP-rand $_k$ achieves a testing accuracy of 54.22% on Shakespeare when $p = 0.2$, outperforming DP-FedAvg by 6%. In the following, we further evaluate Fed-SMP with more experimental settings.

6.3.1 Communication efficiency of Fed-SMP

For the communication cost, we can observe from Table 1-3 that the uplink communication cost for Fed-SMP on both datasets is relatively lower than that of DP-FedAvg and FedAvg under the same privacy guarantee, since our scheme integrates well with the secure aggregation protocol and takes advantage of model update compression. To further demonstrate the communication efficiency of Fed-SMP, we show the convergence speed and communication cost of our Fed-SMP algorithms. Specifically, we select the Fed-SMP algorithms with the optimal compression ratio and

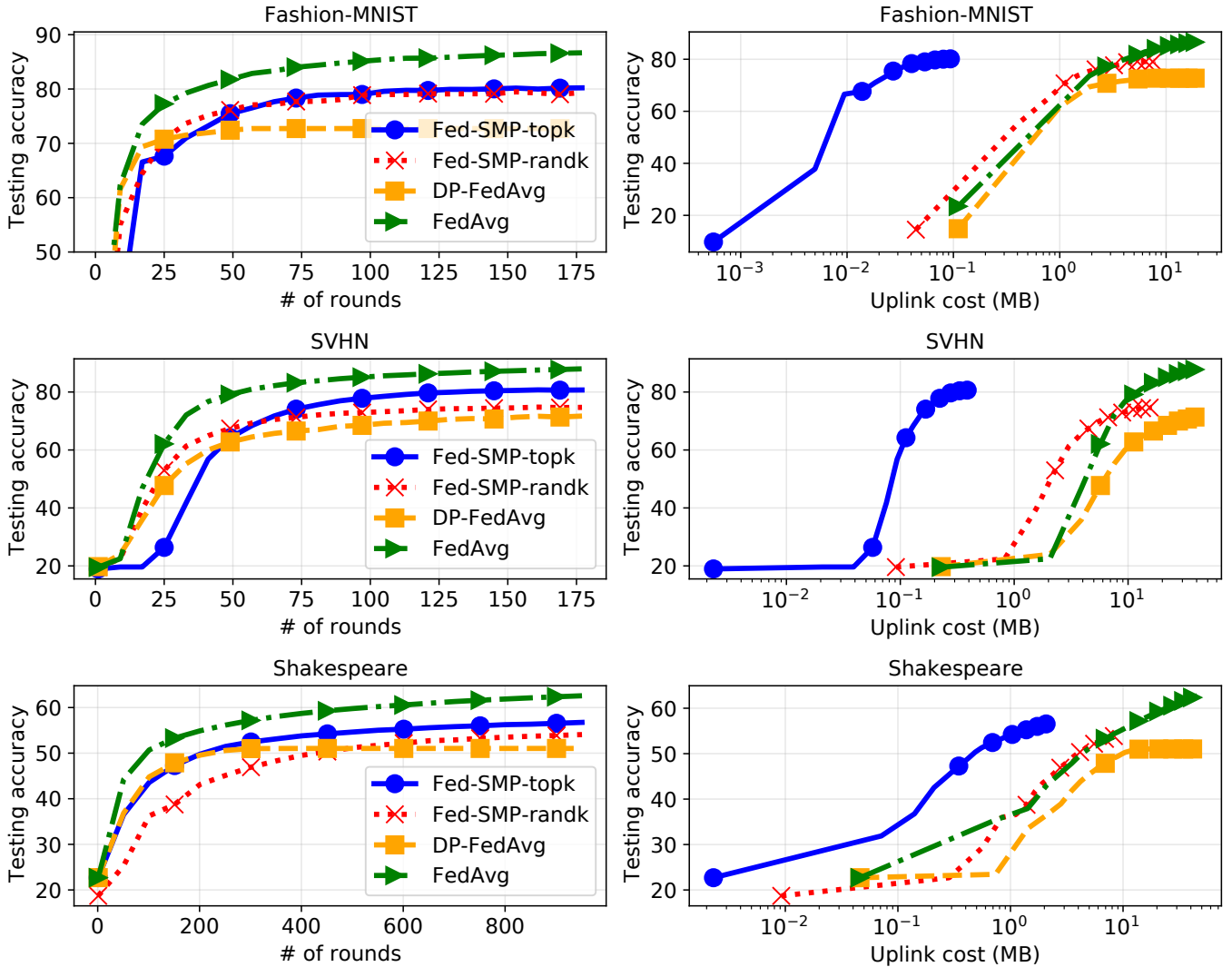


Fig. 2: Testing accuracy of Fed-SMP and DP-FedAvg w.r.t. communication round (left) and uplink communication cost (right).

show its testing accuracy with respect to the communication round and uplink communication cost in Fig. 2, compared with DP-FedAvg and FedAvg: For Fashion-MNIST dataset, we draw the results of Fed-SMP-top_k with $p = 0.005$ and Fed-SMP-rand_k with $p = 0.4$, while achieving $(1.01, n^{-1.1})$ -DP; for SVHN dataset, we draw the results of Fed-SMP-top_k with $p = 0.01$ and Fed-SMP-rand_k with $p = 0.4$, while achieving $(1.01, n^{-1.1})$ -DP; for Shakespeare dataset, we draw the results of Fed-SMP-top_k with $p = 0.05$ and Fed-SMP-rand_k with $p = 0.2$, while achieving $(6.99, n^{-1.1})$ -DP.

From Fig. 2, we observe that DP-FedAvg converges to a lower accuracy than FedAvg due to the added DP noise. Fed-SMP converges slower than DP-FedAvg due to the use of sparsification. However, the final accuracy of Fed-SMP is higher than DP-FedAvg because of the advantage of sparsification in improving privacy and model accuracy. Moreover, Fig. 2 also demonstrates that Fed-SMP is more communication-efficient than DP-FedAvg, and Fed-SMP-top_k is more communication-efficient than Fed-SMP-rand_k. For instance, to achieve a target accuracy 72% on Fashion-

MNIST dataset, Fed-SMP-top_k and Fed-SMP-rand_k save 99% and 70% of uplink communication cost compared with DP-FedAvg, respectively; to achieve a target accuracy 72% on SVHN dataset, Fed-SMP-top_k and Fed-SMP-rand_k save 99% and 82% of uplink communication cost compared with DP-FedAvg, respectively; to achieve a target accuracy 50% on Shakespeare dataset, Fed-SMP-top_k and Fed-SMP-rand_k save 95% and 60% of uplink communication cost compared with DP-FedAvg, respectively.

6.3.2 Tradeoff between privacy and compression errors

From Table 1-Table 3, we have observed that as compression ratio p decreases from 1.0 to a small value (e.g., 0.005), the testing accuracy of Fed-SMP first increases and then decreases. This is due to the change of privacy error and compression error in the convergence bound (7), which is controlled by k . Since the compression ratio $p = k/d$, the smaller p is, the lower the privacy error is, but the compression error could increase. In the following, we conduct additional experiments to further demonstrate this tradeoff under various compression ratios. As it is hard to directly

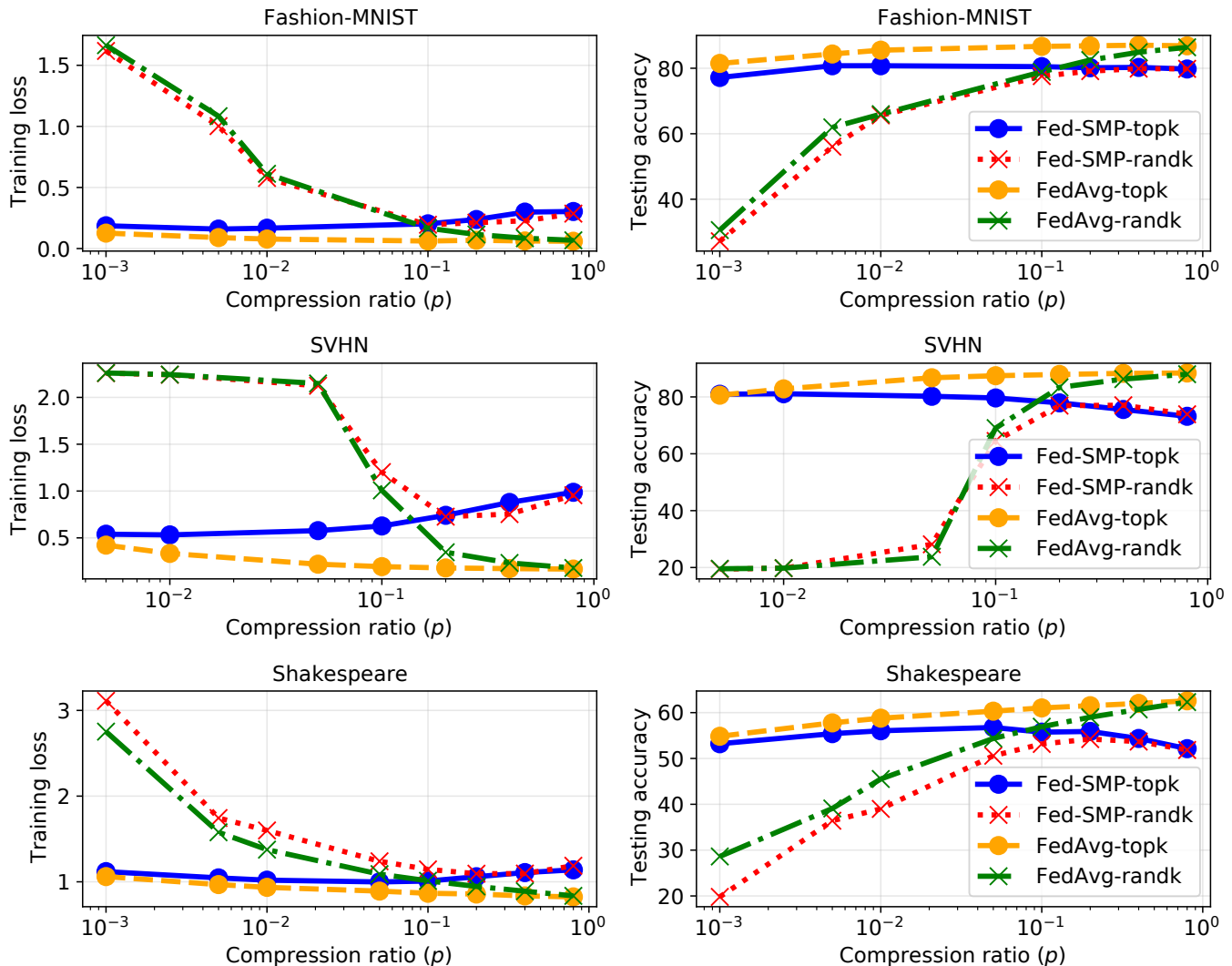


Fig. 3: Training loss and testing accuracy of Fed-SMP w.r.t compression ratio, compared with FedAvg-rand_k/top_k.

show the tradeoff between privacy error and compression error, we study the difference between the performances of Fed-SMP and FedAvg-rand_k/top_k, which differ only in DP noise addition.

Specifically, we show the training losses and testing accuracies of Fed-SMP, FedAvg-rand_k and FedAvg-top_k with respect to different compression ratios in Fig. 3. Note that here Fed-SMP achieves $(1.01, n^{-1.1})$ -DP for Fashion-MNIST and SVHN datasets and $(6.99, n^{-1.1})$ -DP for Shakespeare dataset. From the results, we can see that as the compression ratio p increases, the training losses of FedAvg-rand_k and FedAvg-top_k on all datasets monotonically decrease, and the corresponding testing accuracies monotonically increase, due to the decreasing compression error incurred from sparsification. However, the training losses of Fed-SMP-rand_k and Fed-SMP-top_k decrease at first and then increase, and the corresponding testing accuracies increase at first and then decrease. This matches our observations from Theorem 2, i.e., the performance of Fed-SMP depends on the tradeoff between the compression error and privacy error. When p is too small (e.g., $p < 0.1$ for Fed-SMP-rand_k on Fashion-MNIST datasets), despite the privacy amplifi-

cation effect of compression that reduces the amount of added Gaussian noise, the compression error is significant and dominates the total convergence error, leading to a higher training loss and lower testing accuracy. As p increases, the compression error decreases, leading to a lower training loss. However, if p becomes too large, the privacy amplification effect of compression becomes negligible, and the privacy error starts to dominate the convergence error, leading to a higher training loss. For example, the testing accuracies of Fed-SMP-rand_k start to decrease when $p > 0.4$ for Fashion-MNIST and SVHN datasets and $p > 0.2$ for Shakespeare dataset, and the testing accuracy of Fed-SMP-top_k starts to decrease when $p > 0.005$ on Fashion-MNIST dataset, $p > 0.01$ on SVHN dataset, and $p > 0.05$ for Shakespeare dataset. Therefore, to achieve the best performance of Fed-SMP, the choice of p needs to balance the privacy error and compression error. Furthermore, we find that Fed-SMP-top_k tends to achieve its best performance when p is small (i.e., 0.005, 0.01 and 0.1 for Fashion-MNIST, SVHN and Shakespeare datasets, respectively), and Fed-SMP-rand_k tends to achieve its best performance when p is large (i.e., 0.4, 0.4, 0.2 for Fashion-MNIST, SVHN and

Shakespeare datasets, respectively).

6.3.3 Tradeoff between privacy and accuracy

In this part, we study the tradeoff between privacy and model accuracy of Fed-SMP. In Table 4-7, we show the testing accuracy of Fed-SMP with respect to the privacy loss, under different compression ratios. Note that here we vary the noise multiplier σ to achieve different privacy losses. We can see that under the same compression ratio p , the testing accuracy of Fed-SMP always decreases when the privacy loss ϵ decreases. For example, the testing accuracy of Fed-SMP- rand_k with $p = 0.005$ decreases from 60.02% to 51.72% on Fashion-MNIST dataset when the privacy loss ϵ decreases from 2.02 to 0.44. The reason is that as ϵ decreases, the noise multiplier σ increases (according to Theorem 1), and hence, the convergence error increases (according to Theorem 2). Furthermore, we also observe that under the same privacy loss, the testing accuracy of Fed-SMP always increases at first and then decreases, as the compression ratio p increases. For example, when the privacy loss $\epsilon = 2.02$, the testing accuracy of Fed-SMP- rand_k on SVHN dataset increases from 28.18% to 81.35% and then decreases to 80.39%, as the compression ratio p increases from 0.05 to 0.8. This matches with our observations from Fig. 3, i.e., there exists an optimal compression ratio that achieves the highest testing accuracy under the same privacy loss (as highlighted in the table) due to the tradeoff between privacy and compression errors.

Besides, the optimal compression ratio decreases as the privacy loss decreases. For example, from Table 9, the optimal p for Fed-SMP- top_k on Shakespeare dataset is 0.1 when $\epsilon = 11.51$, and as ϵ decreases to 3.39, the optimal p decreases to 0.01. We observe similar trends on Table 4-8. It is because that the optimal p always balances the tradeoff between privacy and compression errors to minimize the overall convergence error. As ϵ decreases, the noise multiplier σ increases, and thus, the privacy error increases. In this case, the optimal p should decrease to reduce the privacy error.

7 RELATED WORK

FL, or distributed learning in general, with DP has attracted increasing attention recently. Different from centralized learning, FL involves multiple communication rounds between clients and the server; therefore, DP needs to be preserved for all communication rounds. The main challenge of achieving DP in FL lies in maintaining high model accuracy under a reasonable DP guarantee. Prior related works rely on specialized techniques such as shuffling [12], [13], [37], sparsification [38]–[40], and aggregation directly over wireless channels [41] to boost DP guarantee with the same amount of noise, which are only applicable to some specific settings. For instance, the model shuffling method in [12] uses a shuffler to permute the local model updates before aggregation, improving the DP guarantee for protecting the local model update of each client. In [38], the top_k sparsification is integrated with record-level DP to protect the local model update of a client, and in [40] random sparsification is used with gradient perturbation to achieve record-level DP. Our approach is orthogonal to theirs as we aim to achieve client-level DP that protects the

Privacy Loss	Testing Accuracy (%)				
	$p = 0.005$	$p = 0.01$	$p = 0.1$	$p = 0.4$	$p = 0.8$
$\epsilon = 2.02$	60.02	67.77	77.76	81.50	81.89
$\epsilon = 1.01$	56.04	65.61	77.59	79.88	79.82
$\epsilon = 0.58$	54.13	61.90	75.80	69.62	53.39
$\epsilon = 0.44$	51.72	59.08	74.54	52.69	29.26

TABLE 4: Privacy-accuracy tradeoff of Fed-SMP- rand_k with different compression ratios on Fashion-MNIST dataset.

Privacy Loss	Testing Accuracy (%)				
	$p = 0.001$	$p = 0.005$	$p = 0.01$	$p = 0.1$	$p = 0.8$
$\epsilon = 2.02$	78.11	81.10	81.62	82.04	81.71
$\epsilon = 1.01$	77.18	80.76	80.76	80.49	80.16
$\epsilon = 0.58$	76.42	78.76	78.44	75.04	56.02
$\epsilon = 0.44$	74.55	76.90	75.77	67.79	28.55

TABLE 5: Privacy-accuracy tradeoff of Fed-SMP- top_k with different compression ratios on Fashion-MNIST dataset.

local dataset of a client and can be combined to achieve even better privacy-accuracy tradeoffs. Moreover, our work theoretically analyzes the convergence of the proposed methods to demonstrate the benefit of integrating biased/unbiased sparsification with DP in the classic FedAvg algorithm in a rigorous manner, filling the gap in the state-of-arts. Some works [42], [43] also develop differentially private versions of distributed algorithms that require fewer iterations than SGD-based algorithms to converge such as alternating direction method of multipliers, but they are only applicable to convex settings and do not allow partial client participation at each communication round, which is a key feature in FL.

Client-level DP provides more realistic protections than record-level DP against information leakage in the FL setting. With client-level DP, the participation of a client rather than a single data record needs to be protected, therefore requiring the addition of a larger amount of noise than record-level DP to achieve the same level of DP [11], [44]. Our work follows this line of research and proposes to integrate SMP with FedAvg to achieve higher model accuracy than prior schemes, under the same client-level DP guarantee in

Privacy Loss	Testing Accuracy (%)				
	$p = 0.05$	$p = 0.1$	$p = 0.2$	$p = 0.4$	$p = 0.8$
$\epsilon = 2.02$	28.18	65.85	79.76	81.35	80.39
$\epsilon = 1.01$	28.15	64.50	76.94	77.07	73.90
$\epsilon = 0.58$	26.30	60.21	69.90	65.75	61.39
$\epsilon = 0.44$	26.07	55.37	61.69	54.81	24.99

TABLE 6: Privacy-accuracy tradeoff of Fed-SMP- rand_k with different compression ratios on SVHN dataset.

Privacy Loss	Testing Accuracy (%)				
	$p = 0.05$	$p = 0.1$	$p = 0.2$	$p = 0.4$	$p = 0.8$
$\epsilon = 2.02$	81.24	81.86	83.44	81.50	80.08
$\epsilon = 1.01$	80.94	81.12	80.22	75.58	73.18
$\epsilon = 0.58$	78.42	78.85	75.07	62.53	59.50
$\epsilon = 0.44$	78.76	76.45	67.11	32.32	22.92

TABLE 7: Privacy-accuracy tradeoff of Fed-SMP- top_k with different compression ratios on SVHN dataset.

Privacy Loss	Testing Accuracy (%)				
	$p = 0.05$	$p = 0.1$	$p = 0.2$	$p = 0.4$	$p = 0.8$
$\epsilon = 11.51$	50.87	53.16	54.78	55.34	54.30
$\epsilon = 6.99$	50.79	53.15	54.22	53.68	51.85
$\epsilon = 4.70$	50.59	52.87	53.26	51.73	49.36
$\epsilon = 3.39$	50.35	52.39	51.93	49.74	46.18

TABLE 8: Privacy-accuracy tradeoff of Fed-SMP- rand_k with different compression ratios on Shakespeare dataset.

Privacy Loss	Testing Accuracy (%)				
	$p = 0.005$	$p = 0.01$	$p = 0.05$	$p = 0.1$	$p = 0.2$
$\epsilon = 11.51$	55.43	56.10	57.05	57.30	56.97
$\epsilon = 6.99$	55.41	56.04	56.76	55.74	55.89
$\epsilon = 4.70$	55.34	56.02	56.24	54.85	54.61
$\epsilon = 3.39$	55.00	55.87	54.97	54.34	53.00

TABLE 9: Privacy-accuracy tradeoff of Fed-SMP- top_k with different compression ratios on Shakespeare dataset.

FL. Note there are other approaches to preserve privacy in distributed learning, such as secure multi-party computation [45] or homomorphic encryption [46]. However, those approaches have different privacy goals and cannot protect against information leakage incurred from observing the final trained models at the server.

Besides privacy, another bottleneck in FL is the high communication cost of transmitting model parameters. Typical techniques to address this issue include quantization and sparsification [47]–[51]. For example, in [50], sparsification and quantization are integrated with secure aggregation to mitigate the communication cost and privacy concern of FL simultaneously; in [51] the top- k coordinates of a client’s gradient are selected cyclically for all clients to do gradient sparsification to save the communication cost in distributed learning. However, these works do not consider the rigorous privacy protection for clients. Some recent works [52]–[54] start to jointly consider communication efficiency and DP in distributed learning. Agarwal et al. [52] propose cpSGD by making modifications to distributed SGD to make the method both private and communication-efficient. However, the authors treat these as separate issues and develop different approaches to address each within cpSGD. As shown in [55], by separately reducing communication and enforcing privacy, errors in cpSGD are compounded and higher than considering privacy only. In comparison, our proposed Fed-SMP scheme considers those two issues jointly and provides substantially higher accuracy by leveraging sparsification to amplify privacy protection. Zhang et al. [53] also combine sparsification with Gaussian mechanism to achieve communication efficiency and DP in decentralized SGD, but their scheme adds noise before using sparsification and cannot leverage sparsification to improve privacy protection. Kerkouche et al. [54] propose to use compressive sensing before adding noise to model updates in FL but relies on the strong assumption that model updates only have a few non-zero coordinates, which seldom holds in practice.

8 CONCLUSIONS

This paper has proposed Fed-SMP, a new FL scheme based on sparsified model perturbation. Fed-SMP achieves client-level DP while maintaining high model accuracy and is also communication-efficient. We have rigorously analyzed the convergence and end-to-end DP guarantee of the proposed scheme and extensively evaluated its performance on two common benchmark datasets. Experimental results have shown that Fed-SMP with both rand_k and top_k sparsification strategies can improve the privacy-accuracy tradeoff and communication efficiency simultaneously compared with the existing methods. In the future, we will consider other compression methods such as low-rank approximation and quantization, and extend the scheme to alternative FL settings such as shuffled model.

REFERENCES

- [1] A. Pothitos, “IoT and wearables: Fitness tracking,” 2017, <http://www.mobileindustryreview.com/2017/03/iot-wearables-fitness-tracking.html>.
- [2] P. Goldstein, “Smart cities gain efficiencies from IoT traffic sensors and data,” 2018, <https://statetechmagazine.com/article/2018/12/smart-cities-gain-efficiencies-iot-traffic-sensors-and-data-perfcon>.
- [3] A. Weinreic, “The future of the smart home: Smart homes and IoT: A century in the making,” 2018, <https://statetechmagazine.com/article/2018/12/smart-cities-gain-efficiencies-iot-traffic-sensors-and-data-perfcon>.
- [4] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings et al., “Advances and open problems in federated learning,” *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial Intelligence and Statistics*, 2017, pp. 1273–1282.
- [6] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, “Membership inference attacks against machine learning models,” in *IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 3–18.
- [7] M. Fredrikson, S. Jha, and T. Ristenpart, “Model inversion attacks that exploit confidence information and basic countermeasures,” in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2015, pp. 1322–1333.
- [8] C. Dwork, A. Roth et al., “The algorithmic foundations of differential privacy,” *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- [9] R. C. Geyer, T. Klein, and M. Nabi, “Differentially private federated learning: A client level perspective,” *arXiv preprint arXiv:1712.07557*, 2018.
- [10] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security*, 2016, pp. 308–318.
- [11] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, “Learning differentially private recurrent language models,” in *International Conference on Learning Representations*, 2018.
- [12] R. Liu, Y. Cao, H. Chen, R. Guo, and M. Yoshikawa, “FLAME: Differentially private federated learning in the shuffle model,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [13] Ú. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta, “Amplification by shuffling: From local to central differential privacy via anonymity,” in *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2019, pp. 2468–2479.
- [14] I. Mironov, “Rényi differential privacy,” in *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. IEEE, 2017, pp. 263–275.
- [15] Y. Wang, B. Balle, and S. P. Kasiviswanathan, “Subsampled rényi differential privacy and analytical moments accountant,” in *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019, pp. 1226–1235.

- [16] L. Wang, B. Jayaraman, D. Evans, and Q. Gu, "Efficient privacy-preserving nonconvex optimization," *arXiv preprint arXiv:1910.13659*, 2020.
- [17] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 1175–1191.
- [18] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, B. McMahan, T. Van Overveldt, D. Petrou, D. Ramage, and J. Roselander, "Towards federated learning at scale: System design," in *Proceedings of Machine Learning and Systems*, vol. 1, 2019, pp. 374–388.
- [19] S. Goryczka, L. Xiong, and V. Sunderam, "Secure multiparty aggregation with differential privacy: A comparative study," in *Proceedings of the Joint EDBT/ICDT 2013 Workshops*, 2013, pp. 155–163.
- [20] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou, "A hybrid approach to privacy-preserving federated learning," in *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, 2019, pp. 1–11.
- [21] F. Valovich and F. Alda, "Computational differential privacy from lattice-based cryptography," in *International Conference on Number-Theoretic Methods in Cryptology*. Springer, 2017, pp. 121–141.
- [22] K. Bonawitz, F. Salehi, J. Konečný, B. McMahan, and M. Gruteser, "Federated learning with autotuned communication-efficient secure aggregation," in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2019, pp. 1222–1226.
- [23] N. Papernot, S. Song, I. Mironov, A. Raghunathan, K. Talwar, and Ú. Erlingsson, "Scalable private learning with PATE," in *International Conference on Learning Representations*, 2018.
- [24] N. Alon, R. Bassily, and S. Moran, "Limits of private learning with access to public data," *Advances in neural information processing systems*, vol. 32, 2019.
- [25] J. Wang and Z.-H. Zhou, "Differentially private learning with small public data," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 6219–6226.
- [26] R. Sharma, A. Ramakrishna, A. MacLaughlin, A. Rumshisky, J. Majmudar, C. Chung, S. Avestimehr, and R. Gupta, "Federated learning with noisy user feedback," *arXiv preprint arXiv:2205.03092*, 2022.
- [27] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *SIAM Review*, vol. 60, no. 2, pp. 223–311, 2018.
- [28] R. Ward, X. Wu, and L. Bottou, "AdaGrad stepsizes: Sharp convergence over nonconvex landscapes," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6677–6686.
- [29] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of FedAvg on Non-IID data," in *International Conference on Learning Representations*, 2020.
- [30] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [31] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NeurIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [32] N. Papernot, A. Thakurta, S. Song, S. Chien, and U. Erlingsson, "Tempered sigmoid activations for deep learning with differential privacy," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 9312–9321.
- [33] Z. Liang, B. Wang, Q. Gu, S. Osher, and Y. Yao, "Differentially private federated learning with laplacian smoothing," *arXiv preprint arXiv:2005.00218*, 2020.
- [34] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, "Adaptive federated optimization," *arXiv preprint arXiv:2003.00295*, 2020.
- [35] "Opacus PyTorch library," Available from <https://opacus.ai>.
- [36] G. Andrew, O. Thakkar, B. McMahan, and S. Ramaswamy, "Differentially private learning with adaptive clipping," *Advances in Neural Information Processing Systems*, vol. 34, pp. 17 455–17 466, 2021.
- [37] A. Girgis, D. Data, S. Diggavi, P. Kairouz, and A. T. Suresh, "Shuffled model of differential privacy in federated learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 2521–2529.
- [38] R. Liu, Y. Cao, M. Yoshikawa, and H. Chen, "FedSel: Federated SGD under local differential privacy with top-k dimension selection," in *International Conference on Database Systems for Advanced Applications*. Springer, 2020, pp. 485–501.
- [39] L. Lyu, "Dp-signsgd: When efficiency meets privacy and robustness," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 3070–3074.
- [40] R. Hu, Y. Gong, and Y. Guo, "Federated learning with sparsification-amplified privacy and adaptive optimization," *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- [41] M. Seif, R. Tandon, and M. Li, "Wireless federated learning with local differential privacy," in *2020 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2020, pp. 2604–2609.
- [42] Y. Guo and Y. Gong, "Practical collaborative learning for crowdsensing in the internet of things with differential privacy," in *2018 IEEE Conference on Communications and Network Security (CNS)*. IEEE, 2018, pp. 1–9.
- [43] Z. Huang, R. Hu, Y. Guo, E. Chan-Tin, and Y. Gong, "DP-ADMM: ADMM-based distributed learning with differential privacy," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 1002–1012, 2019.
- [44] K. Wei, J. Li, M. Ding, C. Ma, H. Su, B. Zhang, and H. V. Poor, "User-level privacy-preserving federated learning: Analysis and performance optimization," *IEEE Transactions on Mobile Computing*, 2021.
- [45] R. Kanagavelu, Z. Li, J. Samsudin, Y. Yang, F. Yang, R. S. M. Goh, M. Cheah, P. Wiwatphonthana, K. Akkrajitsakul, and S. Wang, "Two-phase multi-party computation enabled privacy-preserving federated learning," in *20th IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing (CCGRID)*. IEEE, 2020, pp. 410–419.
- [46] Y. Aono, T. Hayashi, L. Wang, S. Moriai *et al.*, "Privacy-preserving deep learning via additively homomorphic encryption," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 5, pp. 1333–1345, 2017.
- [47] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, "QSGD: Communication-efficient SGD via gradient quantization and encoding," in *Advances in Neural Information Processing Systems*, 2017, pp. 1709–1720.
- [48] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGD: Compressed optimisation for non-convex problems," in *International Conference on Machine Learning*. PMLR, 2018, pp. 560–569.
- [49] H. Wang, S. Sievert, S. Liu, Z. Charles, D. Papailiopoulos, and S. Wright, "Atomo: Communication-efficient learning via atomic sparsification," in *Advances in Neural Information Processing Systems*, 2018, pp. 9850–9861.
- [50] C. Beguier, M. Andreux, and E. W. Tramel, "Efficient sparse secure aggregation for federated learning," *arXiv preprint arXiv:2007.14861*, 2020.
- [51] C.-Y. Chen, J. Ni, S. Lu, X. Cui, P.-Y. Chen, X. Sun, N. Wang, S. Venkataramani, V. V. Srinivasan, W. Zhang *et al.*, "Scalecom: Scalable sparsified gradient compression for communication-efficient distributed training," *Advances in Neural Information Processing Systems*, vol. 33, pp. 13 551–13 563, 2020.
- [52] N. Agarwal, A. T. Suresh, F. X. X. Yu, S. Kumar, and B. McMahan, "cpSGD: Communication-efficient and differentially-private distributed SGD," in *Advances in Neural Information Processing Systems*, 2018, pp. 7564–7575.
- [53] X. Zhang, M. Fang, J. Liu, and Z. Zhu, "Private and communication-efficient edge learning: A sparse differential gaussian-masking distributed SGD approach," in *Proceedings of the Twenty-First International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, 2020, pp. 261–270.
- [54] R. Kerkouche, G. Ács, C. Castelluccia, and P. Genevès, "Compression boosts differentially private federated learning," in *6th IEEE European Symposium on Security and Privacy*, 2021.
- [55] T. Li, Z. Liu, V. Sekar, and V. Smith, "Privacy for free: Communication-efficient learning with differential privacy using sketches," *arXiv preprint arXiv:1911.00972*, 2019.

APPENDIX A

PROOF OF THEOREM 1

Suppose the client is sampled without replacement with probability $q = r/n$ at each round. By Lemma 2 and Lemma 3, the t -th round of Fed-SMP satisfies $(\alpha, \rho_t(\alpha))$ -RDP, where

$$\rho_t(\alpha) = \frac{3.5q^2\alpha}{\sigma^2}, \quad (8)$$

if $\sigma^2 \geq 0.7$ and $\alpha \leq 1 + (2/3)C^2\sigma^2 \log(1/q\alpha(1 + \sigma^2))$. Then by Lemma 4, Fed-SMP satisfies $(\alpha, T\rho_t(\alpha))$ -RDP after T rounds of training. Next, in order to guarantee (ϵ, δ) -DP according to Lemma 1, we need

$$\frac{3.5q^2T\alpha}{\sigma^2} + \frac{\log(1/\delta)}{\alpha - 1} \leq \epsilon. \quad (9)$$

Suppose α and σ are chosen such that the conditions for (8) are satisfied. Choose $\alpha = 1 + 2\log(1/\delta)/\epsilon$ and rearrange the inequality in (9), we need

$$\sigma^2 \geq \frac{7q^2T(\epsilon + 2\log(1/\delta))}{\epsilon^2}. \quad (10)$$

Then using the constraint on ϵ concludes the proof.

APPENDIX B

PROOF OF THEOREM 2 WITH rand_k SPARSIFIER

B.1 Notations

For ease of expression, let $\theta_i^{t,s}$ denote client i 's local model at local iteration s of round t , and let \mathbf{b}_i^t denote the Gaussian noise added to the sparsified model update where $[\mathbf{b}_i^t]_j \sim \mathcal{N}(0, C^2\sigma^2/r)$ if $[\mathbf{m}^t]_j = 0$ and $[\mathbf{b}_i^t]_j = 0$ otherwise. Let spar denote the sparsifier. In Fed-SMP, the update rule of the global model can be summarized as:

$$\theta^{t+1} = \theta^t - \frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t, \text{ where } \Delta_i^t = \text{spar}(\eta_l \sum_{s=0}^{\tau-1} \mathbf{g}_i^{t,s}) + \mathbf{b}_i^t \quad (11)$$

where Δ_i^t represents the sparsified noisy model update of client i . Assume the clients are sampled uniformly at random without replacement, then we can directly validate that the client sampling is unbiased:

$$\mathbb{E}_{\mathcal{W}^t} \left[\frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right] = \frac{1}{r} \sum_{\substack{\mathcal{W} \in [n], \\ |\mathcal{W}|=r}} \mathbb{P}(\mathcal{W}^t = \mathcal{W}) \sum_{i \in \mathcal{W}^t} \Delta_i^t = \frac{1}{n} \sum_{i=1}^n \Delta_i^t. \quad (12)$$

Moreover, let $\nabla f_i(\theta_i^{t,s})$ represent the local gradient so that $\mathbb{E}_{\xi_i^{t,s}}[\mathbf{g}_i^{t,s}] = \nabla f_i(\theta_i^{t,s})$. For ease of expression, we let $\mathbf{d}_i^t = (1/\tau) \sum_{s=0}^{\tau-1} \mathbf{g}_i^{t,s}$ and $\mathbf{h}_i^t = (1/\tau) \sum_{s=0}^{\tau-1} \nabla f_i(\theta_i^{t,s})$, so we have $\Delta_i^t = \eta_l \tau \text{spar}(\mathbf{d}_i^t) + \mathbf{b}_i^t$ and $\mathbb{E}_t[\mathbf{d}_i^t] = \mathbf{h}_i^t$.

B.2 Useful Inequalities

Lemma 5 (Jensen's inequality). *For arbitrary set of n vectors $\{\mathbf{a}_i\}_{i=1}^n$, $\mathbf{a}_i \in \mathbb{R}^d$ and positive weights $\{w_i\}_{i \in [n]}$,*

$$\left\| \sum_{i=1}^n w_i \mathbf{a}_i \right\|^2 \leq \frac{\sum_{i=1}^n w_i \|\mathbf{a}_i\|^2}{\sum_{i=1}^n w_i}. \quad (13)$$

Lemma 6. *For arbitrary set of n vectors $\{\mathbf{a}_i\}_{i=1}^n$, $\mathbf{a}_i \in \mathbb{R}^d$,*

$$\left\| \sum_{i=1}^n \mathbf{a}_i \right\|^2 \leq n \sum_{i=1}^n \|\mathbf{a}_i\|^2. \quad (14)$$

Lemma 7. *For given two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$,*

$$\|\mathbf{a} + \mathbf{b}\|^2 \leq (1 + \alpha)\|\mathbf{a}\|^2 + (1 + \alpha^{-1})\|\mathbf{b}\|^2, \forall \alpha > 0. \quad (15)$$

Lemma 8. *For given two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$,*

$$2\langle \mathbf{a}, \mathbf{b} \rangle \leq \gamma\|\mathbf{a}\|^2 + \gamma^{-1}\|\mathbf{b}\|^2, \forall \gamma > 0. \quad (16)$$

B.3 Detailed Proof

Let $\text{spar} := (d/k) \text{rand}_k$ denote the unbiased random sparsifier. By the L -smoothness of function f , we have

$$\begin{aligned}
\mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] &\leq \mathbb{E}_t \langle \nabla f(\boldsymbol{\theta}^t), \boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t \rangle + \frac{L}{2} \mathbb{E}_t \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t\|^2 \\
&= -\mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_{\mathcal{W}^t} \left[\frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right] \right\rangle + \frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right\|^2 \\
&= -\mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \frac{1}{n} \sum_{i=1}^n \Delta_i^t \right\rangle + \frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right\|^2 \\
&= \underbrace{-\mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \frac{1}{n} \sum_{i=1}^n (\eta_l \tau \text{spar}(\mathbf{d}_i^t) + \mathbf{b}_i^t) \right\rangle}_{T_1} + \underbrace{\frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} (\eta_l \tau \text{spar}(\mathbf{d}_i^t) + \mathbf{b}_i^t) \right\|^2}_{T_2} \tag{17}
\end{aligned}$$

where the expectation $\mathbb{E}_t[\cdot]$ is taken over the sampled clients \mathcal{W}^t and mini-batches $\xi_i^s, \forall i \in [n], s \in \{0, \dots, \tau - 1\}$ at round t and the sparsifier spar . As $\text{spar} = (d/k) \text{rand}_k$, due to the unbiasedness of the stochastic gradient and Gaussian noise, we have

$$\begin{aligned}
T_1 &= -\left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \eta_l \tau \mathbf{d}_i^t \right] \right\rangle - \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^t \right] \right\rangle \\
&= -\eta_l \tau \mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\rangle \\
&= -\frac{\eta_l \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 - \frac{\eta_l \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 + \frac{\eta_l \tau}{2} \mathbb{E}_t \left\| \nabla f(\boldsymbol{\theta}^t) - \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 \\
&\leq -\frac{\eta_l \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta_l \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \frac{1}{\tau} \sum_{s=0}^{\tau-1} (\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}_i^{t,s})) \right\|^2 \\
&\leq -\frac{\eta_l \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta_l \tau}{2n} \sum_{i=1}^n \frac{1}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}_i^{t,s})\|^2 \\
&\leq -\frac{\eta_l \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta_l L^2}{2n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2, \tag{18}
\end{aligned}$$

where the first inequality results from the fact that $\|(1/n) \sum_{i=1}^n \mathbf{h}_i^t\|^2 \geq 0$, the second inequality uses Lemma 6, and the third inequality uses the L -smoothness of function f_i .

For T_2 , let $\phi_k := d/k - 1$, we have

$$\begin{aligned}
T_2 &\leq \frac{L\eta_l^2 \tau^2}{2} \mathbb{E}_t \left[\frac{1}{r} \sum_{i \in \mathcal{W}^t} \|\text{spar}(\mathbf{d}_i^t)\|^2 \right] + \frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{b}_i^t \right\|^2 \\
&= \frac{L\eta_l^2 \tau^2}{2n} \sum_{i=1}^n \mathbb{E}_t \|\text{spar}(\mathbf{d}_i^t) - \mathbf{d}_i^t + \mathbf{d}_i^t\|^2 + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq \frac{L\eta_l^2 \tau^2}{2n} \sum_{i=1}^n \mathbb{E}_t \left[2\|\text{spar}(\mathbf{d}_i^t) - \mathbf{d}_i^t\|^2 + 2\|\mathbf{d}_i^t\|^2 \right] + \frac{LC^2k\sigma^2}{2r^2} \\
&\leq \frac{L\eta_l^2 \tau^2 (\phi_k + 1)}{n} \sum_{i=1}^n \mathbb{E}_t \|\mathbf{d}_i^t\|^2 + \frac{LC^2k\sigma^2}{2r^2}, \tag{19}
\end{aligned}$$

where the first inequality uses the independence of Gaussian noise and Lemma 5, the second inequality also uses Lemma 6, and the last inequality results from Lemma 9. By Lemma 11, T_2 is bounded as follows:

$$T_2 \leq 2L\eta_l^2 \tau^2 (\phi_k + 1) \beta^2 \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{2L^3 \eta_l^2 \tau (\phi_k + 1)}{n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + L\eta_l^2 \tau^2 (\phi_k + 1) (\bar{\zeta}^2 + 2\kappa^2) + \frac{LC^2k\sigma^2}{2r^2} \tag{20}$$

Combining (17), (18) and (20), one yields

$$\begin{aligned}
\mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] &\leq \left(-\frac{\eta\tau}{2} + 2L\eta_l^2\tau^2(\phi_k + 1)\beta^2\right) \|\nabla f(\boldsymbol{\theta}^t)\|^2 \\
&\quad + \frac{\eta L^2 + 4L^3\eta_l^2\tau(\phi_k + 1)}{2n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + L\eta_l^2\tau^2(\phi_k + 1)(\bar{\zeta}^2 + 2\kappa^2) + \frac{LC^2k\sigma^2}{2r^2} \\
&\leq -\frac{\eta\tau}{4} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta L^2(2\beta^2 + 1)}{4n\beta^2} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + L\eta_l^2\tau^2(\phi_k + 1)(\bar{\zeta}^2 + 2\kappa^2) + \frac{LkC^2\sigma^2}{2r^2},
\end{aligned} \tag{21}$$

if we choose $\eta_l \leq 1/(8\tau L(\phi_k + 1)\beta^2)$. Next, by Lemma 10, we obtain that

$$\begin{aligned}
\mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] &\leq -\frac{\eta\tau}{4} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta L^2(2\beta^2 + 1)}{4\beta^2} \sum_{s=0}^{\tau-1} \left[16\eta_l^2\tau^2\beta^2 \|\nabla f(\boldsymbol{\theta}^t)\|^2 + 16\eta_l^2\tau^2\kappa^2 + 4\tau\eta_l^2\bar{\zeta}^2\right] \\
&\quad + L\eta_l^2\tau^2(\phi_k + 1)(\bar{\zeta}^2 + 2\kappa^2) + \frac{LC^2k\sigma^2}{2r^2} \\
&= \left(-\frac{\eta\tau}{4} + 4\eta L^2(2\beta^2 + 1)\eta_l^2\tau^3\right) \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta_l^3\tau^2L^2(2\beta^2 + 1)(4\tau\kappa^2 + \bar{\zeta}^2)}{\beta^2} \\
&\quad + L\eta_l^2\tau^2(\phi_k + 1)(\bar{\zeta}^2 + 2\kappa^2) + \frac{LC^2k\sigma^2}{2r^2} \\
&\leq -\frac{\eta\tau}{8} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \eta_l^2\tau^2L(12\tau\eta_lL + 2)\kappa^2 + \eta_l^2\tau^2L(3\eta_lL + 1)\bar{\zeta}^2 \\
&\quad + L\eta_l^2\tau^2\phi_k(\bar{\zeta}^2 + 2\kappa^2) + \frac{LC^2k\sigma^2}{2r^2} \\
&\leq -\frac{\eta\tau}{8} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \eta_l^2\tau^2L(3\kappa^2 + 2\bar{\zeta}^2) + \eta_l^2\tau^2L(2\kappa^2 + \bar{\zeta}^2)\phi_k + \frac{LC^2k\sigma^2}{2r^2}
\end{aligned} \tag{22}$$

if $\eta_l \leq \min\{1/(4\tau L\sqrt{4\beta^2 + 2}), 1/(12\tau L)\}$. Rearranging the above inequality in (22) and summing it from $t = 0$ to $T - 1$, we get

$$\begin{aligned}
\sum_{t=0}^{T-1} \|\nabla f(\boldsymbol{\theta}^t)\|^2 &\leq \frac{8}{\eta_l\tau} \mathbb{E}\left[\sum_{t=0}^{T-1} f(\boldsymbol{\theta}^t) - f(\boldsymbol{\theta}^{t+1})\right] + 8T\eta_l\tau L(3\kappa^2 + 2\bar{\zeta}^2) + 8T\eta_l\tau L(2\kappa^2 + \bar{\zeta}^2)\phi_k + \frac{4TLC^2k\sigma^2}{\eta_l\tau r^2} \\
&\leq \frac{8(f(\boldsymbol{\theta}^0) - f(\boldsymbol{\theta}^*))}{\eta_l\tau} + 8T\eta_l\tau L(3\kappa^2 + 2\bar{\zeta}^2) + 8T\eta_l\tau L(2\kappa^2 + \bar{\zeta}^2)\phi_k + \frac{4TLC^2k\sigma^2}{\eta_l\tau r^2},
\end{aligned} \tag{23}$$

where the expectation is taken over all rounds $t \in [0, T - 1]$. Dividing both sides of (23) by T , one yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\boldsymbol{\theta}^t)\|^2 \leq \frac{8(f(\boldsymbol{\theta}^0) - f^*)}{T\eta_l\tau} + 8\eta_l\tau L(3\kappa^2 + 2\bar{\zeta}^2) + 8\eta_l\tau L(2\kappa^2 + \bar{\zeta}^2)\phi_k + \frac{4LC^2k\sigma^2}{\eta_l\tau r^2},$$

if the learning rates η_l satisfy $\eta_l \leq \min\{1/(8\tau L(\phi_k + 1)\beta^2), 1/(4\tau L\sqrt{4\beta^2 + 2}), 1/(12\tau L)\}$. Here, we use the fact that $f^* \leq f(\boldsymbol{\theta}^T)$.

APPENDIX C

PROOF OF THEOREM 2 WITH top_k SPARSIFIER

Here, spar represents the operation of top_k sparsification in Fed-SMP, and we assume that the distribution of the public distribution is similar to the overall distribution $\{\mathcal{D}_i\}_i \in [n]$. By the L -smoothness of function f , we have

$$\begin{aligned}
\mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] &\leq \mathbb{E}_t \langle \nabla f(\boldsymbol{\theta}^t), \boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t \rangle + \frac{L}{2} \mathbb{E}_t \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t\|^2 \\
&= -\mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_{\mathcal{W}^t} \left[\frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right] \right\rangle + \frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right\|^2 \\
&= -\mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \frac{1}{n} \sum_{i=1}^n \Delta_i^t \right\rangle + \frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \Delta_i^t \right\|^2 \\
&= \underbrace{-\mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \frac{1}{n} \sum_{i=1}^n (\eta_l\tau \text{spar}(\mathbf{d}_i^t) + \mathbf{b}_i^t) \right\rangle}_{T_1} + \underbrace{\frac{L}{2} \mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} (\eta_l\tau \text{spar}(\mathbf{d}_i^t) + \mathbf{b}_i^t) \right\|^2}_{T_2}
\end{aligned} \tag{24}$$

where the expectation $\mathbb{E}_t[\cdot]$ is taken over the sampled clients \mathcal{W}^t and mini-batches $\xi_i^s, \forall i \in [n], s \in \{0, \dots, \tau - 1\}$ at round t . Due to the unbiasedness of the stochastic gradient and Gaussian noise, we have

$$\begin{aligned} T_1 &= - \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \eta \tau \mathbf{d}_i^t \right] \right\rangle + \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \eta \tau \mathbf{d}_i^t - \frac{1}{n} \sum_{i=1}^n \eta \tau \text{spar}(\mathbf{d}_i^t) \right] \right\rangle - \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \mathbf{b}_i^t \right] \right\rangle \\ &= \underbrace{-\eta \tau \mathbb{E}_t \left\langle \nabla f(\boldsymbol{\theta}^t), \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\rangle}_{A_1} + \underbrace{\eta \tau \left\langle \nabla f(\boldsymbol{\theta}^t), \mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t - \text{spar} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t \right) \right] \right\rangle}_{A_2} \end{aligned} \quad (25)$$

since the clients at each round use the same top_k sparsifier. Using the fact that $2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2$, we have

$$\begin{aligned} A_1 &= -\frac{\eta \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 - \frac{\eta \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 + \frac{\eta \tau}{2} \mathbb{E}_t \left\| \nabla f(\boldsymbol{\theta}^t) - \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 \\ &= -\frac{\eta \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \frac{1}{\tau} \sum_{s=0}^{\tau-1} (\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}_i^{t,s})) \right\|^2 - \frac{\eta \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 \\ &\leq -\frac{\eta \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau}{2n} \sum_{i=1}^n \frac{1}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}_i^{t,s})\|^2 - \frac{\eta \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 \\ &\leq -\frac{\eta \tau}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta L^2}{2n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 - \frac{\eta \tau}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2, \end{aligned} \quad (26)$$

where the first inequality uses Lemma 6, and the second inequality uses the L -smoothness of function f_i . Next, let $\phi_k := 1 - k/d$, we get

$$\begin{aligned} A_2 &\leq \frac{\eta \tau}{2} \left(\gamma \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \gamma^{-1} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t - \text{spar} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t \right) \right\|^2 \right) \\ &\leq \frac{\eta \tau \gamma}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau \phi_k}{2\gamma} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t \right\|^2 \\ &= \frac{\eta \tau \gamma}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau \phi_k}{2\gamma} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t - \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t + \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 \\ &= \frac{\eta \tau \gamma}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau \phi_k}{2\gamma} \mathbb{E}_t \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t - \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 + 2 \left\langle \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t, \frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t - \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\rangle \right] \\ &= \frac{\eta \tau \gamma}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau \phi_k}{2\gamma} \mathbb{E}_t \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{d}_i^t - \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 \right] \end{aligned} \quad (27)$$

$$\leq \frac{\eta \tau \gamma}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau \phi_k}{2n\gamma} \sum_{i=1}^n \mathbb{E}_t \|\mathbf{d}_i^t - \mathbf{h}_i^t\|^2 + \frac{\eta \tau \phi_k}{2\gamma} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2, \quad (28)$$

where the first inequality uses Lemma 8, the second inequality uses Lemma 9, and the third inequality uses Lemma 6. Given that

$$\mathbb{E}_t \|\mathbf{d}_i^t - \mathbf{h}_i^t\|^2 = \mathbb{E}_t \left\| \frac{1}{\tau} \sum_{s=0}^{\tau-1} (\mathbf{g}_i^{t,s} - \nabla f_i(\boldsymbol{\theta}_i^{t,s})) \right\|^2 \leq \frac{1}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\mathbf{g}_i^{t,s} - \nabla f_i(\boldsymbol{\theta}_i^{t,s})\|^2 \leq \frac{1}{\tau} \sum_{s=0}^{\tau-1} \zeta_i^2 = \zeta_i^2, \quad (29)$$

by Lemma 6 and Assumption 2, we have

$$A_2 \leq \frac{\eta \tau \gamma}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau \phi_k}{2n\gamma} \sum_{i=1}^n \zeta_i^2 + \frac{\eta \tau \phi_k}{2\gamma} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2. \quad (30)$$

Combining (25), (26) and (27) and let $\bar{\zeta}^2 := (1/n) \sum_{i=1}^n \zeta_i^2$, we get

$$\begin{aligned} T_1 &\leq -\frac{\eta \tau (1-\gamma)}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta \tau L^2}{2n} \sum_{i=1}^n \frac{1}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 - \frac{\eta \tau (1-\phi_k/\gamma)}{2} \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{h}_i^t \right\|^2 + \frac{\eta \tau \phi_k \bar{\zeta}^2}{2\gamma} \\ &\leq -\frac{\eta \tau (1-\gamma)}{2} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta L^2}{2n} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + \frac{\eta \tau \phi_k \bar{\zeta}^2}{2\gamma}, \end{aligned} \quad (31)$$

if $\phi_k \leq \gamma$.

To bound T_2 , we utilize the independence of Gaussian noise and obtain that

$$\begin{aligned}
T_2 &= \frac{L\eta_l^2\tau^2}{2}\mathbb{E}_t \left\| \text{spar} \left(\frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right) \right\|^2 + \frac{L}{2}\mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{b}_i^t \right\|^2 \\
&= \frac{L\eta_l^2\tau^2}{2}\mathbb{E}_t \left\| \text{spar} \left(\frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right) - \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t + \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right\|^2 + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq L\eta_l^2\tau^2\mathbb{E}_t \left\| \text{spar} \left(\frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right) - \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right\|^2 + L\eta_l^2\tau^2\mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right\|^2 + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq L\eta_l^2\tau^2(\phi_k + 1)\mathbb{E}_t \left\| \frac{1}{r} \sum_{i \in \mathcal{W}^t} \mathbf{d}_i^t \right\|^2 + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq L\eta_l^2\tau^2(\phi_k + 1)\mathbb{E}_t \left[\frac{1}{r} \sum_{i \in \mathcal{W}^t} \|\mathbf{d}_i^t\|^2 \right] + \frac{LkC^2\sigma^2}{2r^2} \\
&= L\eta_l^2\tau^2(\phi_k + 1)\mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \|\mathbf{d}_i^t\|^2 \right] + \frac{LkC^2\sigma^2}{2r^2}, \tag{32}
\end{aligned}$$

where the first inequality uses the fact that $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$, the second inequality uses Lemma 9, and the last inequality results from Lemma 5. Then, by Lemma 11, we get

$$T_2 \leq L\eta_l^2\tau^2(\phi_k + 1) \left(\frac{2L^2}{n\tau} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + 2(\beta^2 \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \kappa^2) + \bar{\zeta}^2 \right) + \frac{LkC^2\sigma^2}{2r^2} \tag{33}$$

Combining (24), (31) and (33), one yields

$$\begin{aligned}
\mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] &\leq \left(-\frac{\eta\tau(1-\gamma)}{2} + 2L\eta_l^2\tau^2(\phi_k + 1)\beta^2 \right) \|\nabla f(\boldsymbol{\theta}^t)\|^2 \\
&\quad + \left(\frac{\eta L^2}{2n} + \frac{2L\eta_l^2\tau^2(\phi_k + 1)L^2}{n\tau} \right) \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + \frac{\eta\tau\phi_k\bar{\zeta}^2}{2\gamma} \\
&\quad + L\eta_l^2\tau^2(\phi_k + 1)(2\kappa^2 + \bar{\zeta}^2) + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq -\frac{\eta\tau}{4} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta L^2(2\beta^2 + 1)}{4n\beta^2} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + \frac{\eta\tau\phi_k\bar{\zeta}^2}{2\gamma} \\
&\quad + L\eta_l^2\tau^2(\phi_k + 1)(2\kappa^2 + \bar{\zeta}^2) + \frac{LkC^2\sigma^2}{2r^2} \tag{34}
\end{aligned}$$

if $\eta_l \leq (1 - 2\gamma)/(8\tau L(\phi_k + 1)\beta^2)$ and $\gamma < 1/2$. Then, by Lemma 10, we have

$$\begin{aligned}
\mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] &\leq -\frac{\eta\tau}{4} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta L^2(2\beta^2 + 1)}{4\beta^2} \sum_{s=0}^{\tau-1} \left[16\eta_l^2\tau^2\beta^2 \|\nabla f(\boldsymbol{\theta}^t)\|^2 + 16\eta_l^2\tau^2\kappa^2 + 4\tau\eta_l^2\bar{\zeta}^2 \right] + \frac{\eta\tau\phi_k\bar{\zeta}^2}{2\gamma} \\
&\quad + L\eta_l^2\tau^2(\phi_k + 1)(2\kappa^2 + \bar{\zeta}^2) + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq \left(-\frac{\eta\tau}{4} + 4\eta_l^3\tau^3L^2(2\beta^2 + 1) \right) \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \frac{\eta_l^3\tau^2L^2(2\beta^2 + 1)(4\tau\kappa^2 + \bar{\zeta}^2)}{\beta^2} + \frac{\eta\tau\phi_k\bar{\zeta}^2}{2\gamma} \\
&\quad + L\eta_l^2\tau^2(\phi_k + 1)(2\kappa^2 + \bar{\zeta}^2) + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq -\frac{\eta\tau}{8} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \eta_l^2\tau^2L(12\eta_l\tau L + 2)\kappa^2 + \eta_l^2\tau^2L(3\eta_l L + 1)\bar{\zeta}^2 \\
&\quad + \left(L\eta_l^2\tau^2(2\kappa^2 + \bar{\zeta}^2) + \frac{\eta\tau\bar{\zeta}^2}{2\gamma} \right) \phi_k + \frac{LkC^2\sigma^2}{2r^2} \\
&\leq -\frac{\eta\tau}{8} \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \eta_l^2\tau^2L(3\kappa^2 + 2\bar{\zeta}^2) + \left(L\eta_l^2\tau^2(2\kappa^2 + \bar{\zeta}^2) + \frac{\eta\tau\bar{\zeta}^2}{2\gamma} \right) \phi_k + \frac{LkC^2\sigma^2}{2r^2}, \tag{35}
\end{aligned}$$

if $\eta_l \leq \min\{1/(4\tau L\sqrt{4\beta^2 + 2}), 1/(12\tau L)\}$. Rearranging the above inequality in (35) and summing it from $t = 0$ to $T - 1$, we get

$$\sum_{t=0}^{T-1} \|\nabla f(\boldsymbol{\theta}^t)\|^2 \leq \frac{8}{\eta_l \tau} \sum_{t=0}^{T-1} \mathbb{E}_t[f(\boldsymbol{\theta}^{t+1}) - f(\boldsymbol{\theta}^t)] + 8T\eta_l \tau L(3\kappa^2 + 2\bar{\zeta}^2) + 8T \left(\eta_l \tau L(2\kappa^2 + \bar{\zeta}^2) + \frac{\bar{\zeta}^2}{2\gamma} \right) \phi_k + \frac{4TLkC^2\sigma^2}{\eta_l \tau r^2}, \quad (36)$$

where the expectation is taken over all rounds $t \in [0, T - 1]$. Dividing both sides of (36) by T , one yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\boldsymbol{\theta}^t)\|^2 \leq \frac{8(f(\boldsymbol{\theta}^0) - f^*)}{T\eta_l \tau} + 8\eta_l \tau L(3\kappa^2 + 2\bar{\zeta}^2) + 8 \left(\eta_l \tau L(2\kappa^2 + \bar{\zeta}^2) + \frac{\bar{\zeta}^2}{2\gamma} \right) \phi_k + \frac{4LkC^2\sigma^2}{\eta_l \tau r^2},$$

Here, we use the fact that $f^* \leq f(\boldsymbol{\theta}^T)$. We summarize the convergence results of Fed-SMP with rand_k sparsifier and top_k sparsifier in Theorem 2 by selecting a reasonable constant $\gamma = 1/3$ (which satisfies $\gamma < 1/2$).

APPENDIX D INTERMEDIATE RESULTS

Lemma 9 (Bounded Sparsification). *Given a vector $\mathbf{x} \in \mathbb{R}^d$, parameter $k \in [d]$. The sparsifier $\{\text{rand}_k(\mathbf{x}), \text{top}_k(\mathbf{x})\}$ holds that*

$$\mathbb{E} \|\text{top}_k(\mathbf{x}) - \mathbf{x}\|^2 \leq \left(1 - \frac{k}{d}\right) \|\mathbf{x}\|^2; \quad \mathbb{E} \left\| \frac{d}{k} \text{rand}_k(\mathbf{x}) - \mathbf{x} \right\|^2 \leq \left(\frac{d}{k} - 1\right) \|\mathbf{x}\|^2$$

Proof. For the rand_k sparsifier, by applying the expectation over the active set ω , we have

$$\mathbb{E}_\omega \left[\frac{d}{k} \text{rand}_k(\mathbf{x}) \right] = \frac{d}{k} \left[\frac{k}{d} [\mathbf{x}]_1, \dots, \frac{k}{d} [\mathbf{x}]_d \right] = \mathbf{x},$$

$$\begin{aligned} \mathbb{E}_\omega \|\text{rand}_k(\mathbf{x}) - \mathbf{x}\|^2 &= \sum_{j=1}^d \left(\frac{k}{d} ([\mathbf{x}]_j - [\mathbf{x}]_j)^2 + \left(1 - \frac{k}{d}\right) [\mathbf{x}]_j^2 \right) \\ &= \left(1 - \frac{k}{d}\right) \|\mathbf{x}\|^2, \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}_\omega \left\| \frac{d}{k} \text{rand}_k(\mathbf{x}) - \mathbf{x} \right\|^2 &= \sum_{j=1}^d \left(\frac{d}{k} \left(\frac{k}{d} [\mathbf{x}]_j - [\mathbf{x}]_j \right)^2 + \left(1 - \frac{k}{d}\right) [\mathbf{x}]_j^2 \right) \\ &= \left(\frac{k}{d} \left(\frac{d}{k} - 1 \right)^2 + \left(1 - \frac{k}{d}\right) \right) \sum_{j=1}^d [\mathbf{x}]_j^2 \\ &= \left(\frac{d}{k} - 1 \right) \|\mathbf{x}\|^2, \end{aligned}$$

For the top_k sparsifier, as $\mathbb{E} \|\text{top}_k(\mathbf{x}) - \mathbf{x}\|^2 \leq \mathbb{E} \|\text{rand}_k(\mathbf{x}) - \mathbf{x}\|^2$, we have

$$\mathbb{E} \|\text{top}_k(\mathbf{x}) - \mathbf{x}\|^2 \leq \left(1 - \frac{k}{d}\right) \|\mathbf{x}\|^2. \quad \square$$

Lemma 10 (Bounded Local Divergence). *The local model difference at round t is bounded as follows:*

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s} \right\|^2 \leq 32\eta_l^2 \tau^2 \beta^2 \|\nabla f(\boldsymbol{\theta}^t)\|^2 + 32\eta_l^2 \tau^2 \kappa^2 + 4\tau\eta_l^2 \bar{\zeta}^2. \quad (37)$$

where $\bar{\zeta}^2 := (1/n) \sum_{i=1}^n \zeta_i^2$.

Proof. Plugging into the local update rule, we have

$$\begin{aligned}
\mathbb{E}_t \left\| \boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s} \right\|^2 &= \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t - \eta_l \mathbf{g}_i^{t,s-1} \right\|^2 \\
&= \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t - \eta_l \mathbf{g}_i^{t,s-1} + \eta_l \nabla f_i(\boldsymbol{\theta}_i^{t,s-1}) - \eta_l \nabla f_i(\boldsymbol{\theta}_i^{t,s-1}) + \eta_l \nabla f_i(\boldsymbol{\theta}^t) - \eta_l \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 \\
&= \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t - \eta_l \nabla f_i(\boldsymbol{\theta}_i^{t,s-1}) + \eta_l \nabla f_i(\boldsymbol{\theta}^t) - \eta_l \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 + \eta_l^2 \mathbb{E}_t \left\| \mathbf{g}_i^{t,s-1} - \nabla f_i(\boldsymbol{\theta}_i^{t,s-1}) \right\|^2 \\
&\leq \left(1 + \frac{1}{2\tau-1} \right) \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + 2\eta_l^2 \tau \mathbb{E}_t \left\| \nabla f_i(\boldsymbol{\theta}_i^{t,s-1}) + \nabla f_i(\boldsymbol{\theta}^t) - \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 + \eta_l^2 \zeta_i^2 \\
&\leq \left(1 + \frac{1}{2\tau-1} \right) \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + 4\eta_l^2 \tau \mathbb{E}_t \left\| \nabla f_i(\boldsymbol{\theta}_i^{t,s-1}) - \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 + 4\eta_l^2 \tau \left\| \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 + \eta_l^2 \zeta_i^2 \\
&\leq \left(1 + \frac{1}{2\tau-1} \right) \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + 4\eta_l^2 L^2 \tau \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + 4\eta_l^2 \tau \left\| \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 + \eta_l^2 \zeta_i^2,
\end{aligned}$$

by using Lemma 7 and Assumption 2, and hence,

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s} \right\|^2 &\leq \left(1 + \frac{1}{2\tau-1} + 4\eta_l^2 L^2 \tau \right) \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + \frac{4\eta_l^2 \tau}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 + \frac{\eta_l^2}{n} \sum_{i=1}^n \zeta_i^2 \\
&\leq \left(1 + \frac{1}{2\tau-1} + 4\eta_l^2 L^2 \tau \right) \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + 4\eta_l^2 \tau \beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + 4\eta_l^2 \tau \kappa^2 + \frac{\eta_l^2}{n} \sum_{i=1}^n \zeta_i^2 \\
&\leq \left(1 + \frac{1}{\tau-1} \right) \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \boldsymbol{\theta}_i^{t,s-1} - \boldsymbol{\theta}^t \right\|^2 + 4\eta_l^2 \tau \beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + 4\eta_l^2 \tau \kappa^2 + \frac{\eta_l^2}{n} \sum_{i=1}^n \zeta_i^2, \tag{38}
\end{aligned}$$

when $\eta_l \leq 1/3\tau L$. Unrolling the recursion, we obtain the following:

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s} \right\|^2 &\leq \sum_{h=0}^{s-1} \left(1 + \frac{1}{\tau-1} \right)^h \left[4\eta_l^2 \tau \beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + 4\eta_l^2 \tau \kappa^2 + \frac{\eta_l^2}{n} \sum_{i=1}^n \zeta_i^2 \right] \\
&\leq (\tau-1) \left[\left(1 + \frac{1}{\tau-1} \right)^\tau - 1 \right] \times \left[4\eta_l^2 \tau \beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + 4\eta_l^2 \tau \kappa^2 + \frac{\eta_l^2}{n} \sum_{i=1}^n \zeta_i^2 \right] \\
&\leq 16\eta_l^2 \tau^2 \beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + 16\eta_l^2 \tau^2 \kappa^2 + \frac{4\tau\eta_l^2}{n} \sum_{i=1}^n \zeta_i^2, \tag{39}
\end{aligned}$$

where the last inequality results from the fact that $\left(1 + \frac{1}{\tau-1} \right)^\tau \leq 5$ when $\tau > 1$. \square

Lemma 11 (Bounded Local Model Update). *The local model update \mathbf{d}_i^t at round t is bounded as follows:*

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \mathbf{d}_i^t \right\|^2 \leq \frac{2L^2}{n\tau} \sum_{i=1}^n \sum_{s=0}^{\tau-1} \mathbb{E}_t \left\| \boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s} \right\|^2 + 2(\beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + \kappa^2) + \bar{\zeta}^2. \tag{40}$$

where $\bar{\zeta}^2 := (1/n) \sum_{i=1}^n \zeta_i^2$.

Proof. Since $\mathbb{E}_t[\mathbf{d}_i^t] = \mathbf{h}_i^t$, we have

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \mathbf{d}_i^t \right\|^2 &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \mathbf{d}_i^t - \mathbf{h}_i^t + \mathbf{h}_i^t \right\|^2 \\
&= \frac{1}{n} \sum_{i=1}^n \left[\mathbb{E}_t \left\| \mathbf{d}_i^t - \mathbf{h}_i^t \right\|^2 + \mathbb{E}_t \left\| \mathbf{h}_i^t \right\|^2 + \mathbb{E}_t \langle \mathbf{d}_i^t - \mathbf{h}_i^t, \mathbf{h}_i^t \rangle \right] \\
&= \frac{1}{n} \sum_{i=1}^n \left[\mathbb{E}_t \left\| \mathbf{d}_i^t - \mathbf{h}_i^t \right\|^2 + \mathbb{E}_t \left\| \mathbf{h}_i^t \right\|^2 \right] \\
&= \frac{1}{n} \sum_{i=1}^n \left[\mathbb{E}_t \left\| \mathbf{d}_i^t - \mathbf{h}_i^t \right\|^2 + \mathbb{E}_t \left\| \mathbf{h}_i^t - \nabla f_i(\boldsymbol{\theta}^t) + \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 \right] \\
&\leq \frac{1}{n} \sum_{i=1}^n \left[\mathbb{E}_t \left\| \mathbf{d}_i^t - \mathbf{h}_i^t \right\|^2 + 2\mathbb{E}_t \left\| \mathbf{h}_i^t - \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 \right] + \frac{2}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 \\
&\leq \frac{1}{n} \sum_{i=1}^n \left[\mathbb{E}_t \left\| \mathbf{d}_i^t - \mathbf{h}_i^t \right\|^2 + 2\mathbb{E}_t \left\| \mathbf{h}_i^t - \nabla f_i(\boldsymbol{\theta}^t) \right\|^2 \right] + 2(\beta^2 \left\| \nabla f(\boldsymbol{\theta}^t) \right\|^2 + \kappa^2) \tag{41}
\end{aligned}$$

where the first inequality uses Lemma 6, and the last inequality uses Assumption 3. Given that

$$\mathbb{E}_t \|\mathbf{d}_i^t - \mathbf{h}_i^t\|^2 = \mathbb{E}_t \left\| \frac{1}{\tau} \sum_{s=0}^{\tau-1} (\mathbf{g}_i^{t,s} - \nabla f_i(\boldsymbol{\theta}_i^{t,s})) \right\|^2 \leq \frac{1}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\mathbf{g}_i^{t,s} - \nabla f_i(\boldsymbol{\theta}_i^{t,s})\|^2 \leq \frac{1}{\tau} \sum_{s=0}^{\tau-1} \zeta_i^2 = \zeta_i^2, \quad (42)$$

by Lemma 6 and Assumption 2, and

$$\mathbb{E}_t \|\mathbf{h}_i^t - \nabla f_i(\boldsymbol{\theta}^t)\|^2 = \mathbb{E}_t \left\| \frac{1}{\tau} \sum_{s=0}^{\tau-1} (\nabla f_i(\boldsymbol{\theta}_i^{t,s}) - \nabla f_i(\boldsymbol{\theta}^t)) \right\|^2 \leq \frac{1}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\nabla f_i(\boldsymbol{\theta}_i^{t,s}) - \nabla f_i(\boldsymbol{\theta}^t)\|^2 \leq \frac{L^2}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 \quad (43)$$

by Lemma 6 and the L -smoothness of function f_i , one yields

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \|\mathbf{d}_i^t\|^2 \leq \bar{\zeta}^2 + \frac{2}{n} \sum_{i=1}^n \frac{L^2}{\tau} \sum_{s=0}^{\tau-1} \mathbb{E}_t \|\boldsymbol{\theta}^t - \boldsymbol{\theta}_i^{t,s}\|^2 + 2(\beta^2 \|\nabla f(\boldsymbol{\theta}^t)\|^2 + \kappa^2). \quad (44)$$

□