# SoteriaFL: A Unified Framework for Private Federated Learning with Communication Compression

Anonymous Author(s) Affiliation Address email

### Abstract

To enable large-scale machine learning in bandwidth-hungry environments such 1 as wireless networks, significant progress has been made recently in designing 2 communication-efficient federated learning algorithms with the aid of communi-3 cation compression. On the other end, privacy-preserving, especially at the client 4 level, is another important desideratum that has not been addressed simultaneously 5 in the presence of advanced communication compression techniques yet. In this 6 paper, we propose a unified framework that enhances the communication efficiency 7 of private federated learning with communication compression. Exploiting both 8 general compression operators and local differential privacy, we first examine a sim-9 ple algorithm that applies compression directly to differentially-private stochastic 10 gradient descent, and identify its limitations. We then propose a unified framework 11 SoteriaFL for private federated learning, which accommodates a general family 12 of local gradient estimators including popular stochastic variance-reduced gradi-13 ent methods and the state-of-the-art shifted compression scheme. We provide a 14 15 comprehensive characterization of its performance trade-offs in terms of privacy, 16 utility, and communication complexity, where SoteriaFL is shown to achieve better communication complexity without sacrificing privacy nor utility than other private 17 federated learning algorithms without communication compression. 18

#### **19 1 Introduction**

With the proliferation of mobile and edge devices, federated learning (FL) [35, 41] has recently emerged as a disruptive paradigm for training large-scale machine learning models over a vast amount of geographically distributed and heterogeneous devices. For instance, Google uses FL in the Gboard mobile keyboard for next word predictions [23]. FL is often modeled as a distributed optimization problem [34, 35, 41, 29, 52], aiming to solve

$$\min_{\boldsymbol{x}\in\mathbb{R}^d} \left\{ f(\boldsymbol{x}; D) := \frac{1}{n} \sum_{i=1}^n f(\boldsymbol{x}; D_i) \right\}, \text{ where } f(\boldsymbol{x}; D_i) := \frac{1}{m} \sum_{j=1}^m f(\boldsymbol{x}; d_{i,j}).$$
(1)

Here, D denotes the entire dataset distributed across all n clients, where each client i has a local dataset  $D_i = \{d_{i,j}\}_{j=1}^m$  of equal size m,  ${}^1 x \in \mathbb{R}^d$  denotes the model parameters, f(x; D),  $f(x; D_i)$ , and  $f(x; d_{i,j})$  denote the nonconvex loss function of the current model x on the entire dataset D, the local dataset  $D_i$ , and a single data sample  $d_{i,j}$ , respectively. For simplicity, we use f(x),  $f_i(x)$  and

29  $f_{i,j}(\boldsymbol{x})$  to denote  $f(\boldsymbol{x}; D), f(\boldsymbol{x}; D_i)$  and  $f(\boldsymbol{x}; d_{i,j})$ , respectively.

<sup>&</sup>lt;sup>1</sup>This is without loss of generality, since otherwise one can simply adjust the weights of the loss function.

#### 30 1.1 Motivation: privacy-utility-communication trade-offs

To unleash the full potential of FL, it is extremely important that the algorithm designed to solve (1) needs to meet several competing desiderata.

Communication efficiency Communication between the server and clients is well recognized as the main bottleneck for optimizing the latency of FL systems, especially when the clients—such as mobile devices—have limited bandwidth, the number of clients is large, and/or the machine learning model has a lot of parameters—for example, the language model GPT-3 [6] has billions of parameters and therefore cumbersome to share directly.

Therefore, it is very important to design FL algorithms to reduce the overall communication cost, 38 which takes into account both the number of communication rounds and the communication cost 39 per communication round for reaching a desired accuracy. With these two quantities in mind, there 40 are two principal approaches for communication-efficient FL: 1) local methods, where in each 41 communication round, clients run multiple local update steps before communicating with the server, 42 in the hope of reducing the number of communication rounds, e.g., FedAvg [41], Local-SVRG [21], 43 SCAFFOLD [32], FedPAGE [57], and ProxSkip [43]; 2) compression methods, where clients send 44 compressed communication message to the server, in the hope of reducing the communication cost 45 per communication round, e.g., [3, 33, 51, 25, 31, 42, 44, 22, 38, 45, 37]. While both categories 46 have garnered significant attention in recent years, we will focus on the second approach based on 47 48 communication compression to enhance communication efficiency.

49 Privacy preserving While FL holds great promise of harnessing the inferential power of private data 50 stored on a large number of distributed clients, these local data at clients often contain sensitive or 51 proprietary information without consent to share. Although FL may appear to protect the data privacy 52 via storing data locally and only sharing the model updates (e.g., gradient information), the training 53 process can nonetheless reveal sensitive information as demonstrated by, e.g., Zhu et al. [59]. It is 54 thus desirable for FL to preserve privacy in a guaranteed manner [19, 29, 46, 52].

To ensure the training process does not accidentally leak private information, advanced privacy-55 preserving tools such as *differential privacy* (DP) [16] have been widely integrated into training 56 algorithms [14, 8, 15, 1, 50, 26, 11, 18]. A notable example is DP-SGD [1], which developed a 57 differentially-private stochastic gradient descent (SGD) algorithm in the centralized (single-node) 58 setting. More recently, several differentially-private algorithms [27, 53, 47, 40] are proposed for 59 the more general distributed (*n*-node) setting suitable for FL. In this paper, we also follow the DP 60 approach to preserve privacy. In particular, we adopt local differential privacy (LDP) to respect the 61 privacy of each client, which is critical in FL. 62

Encouraged by recent advances in communication compression techniques, and the widespread
 success of differentially-private methods, a natural question is

65 Can we develop a unified framework for private federated learning with communication compression,
 66 and understand the trade-offs between privacy, utility, and communication?

Note that there have been a handful of works that simultaneously address compression and privacy 67 in FL. Unfortunately, they only provide partial answers to the above question. Most of the existing 68 works only consider specific, elementary or tailored compression schemes that are applied directly to 69 the gradient messages in DP-SGD [2, 54, 20, 60, 56, 13]. A number of works [48, 9, 10, 30, 17, 49] 70 extended and considered different compression schemes, but did not provide concrete trade-offs 71 in terms of privacy, utility and communication. Furthermore, existing theoretical analyses can be 72 limited only to convex problems [20], lacking in some aspects such as utility [60], or delivering 73 pessimistic guarantees on utility and / or communication due to strong assumptions [56, 13]. Finally, 74 existing work only studied the DP framework for direct compression, while it is known that the 75 recently developed shifted compression scheme [42, 24] achieves much better convergence. Due to 76 noise injection for privacy-preserving, it is a priori unclear if the shifted compression scheme is also 77 compatible with privacy. 78

#### 79 **1.2 Our contributions**

<sup>80</sup> In this paper, we answer the above question by providing a general approach that enhances the <sup>81</sup> communication efficiency of private federated learning in the *nonconvex* setting, through a unified <sup>82</sup> framework called **SoteriaFL** (see Algorithm 2). Specifically, we have the following contributions.

Table 1: Comparisons among (local) differentially-private algorithms for the nonconvex problem (1) in both central (single-node) and distributed (*n*-node) settings. Here, *m* denotes the number of data stored on a single client, *n* is the number of clients, *d* is the dimension, and  $\omega$  is the parameter for the compression operator (cf. Definition 1). The communication complexity is computed by  $ndT/(1+\omega)$ , where *T* is the total number of communication rounds, and  $nd/(1+\omega)$  is the communication cost per round. The utility / accuracy measures the average squared gradient norm of the objective function after *T* rounds. Note that the algorithm is better when the utility / accuracy and the communication complexity are small under the same privacy guarantee.

Algorithm	Privacy	Utility / Accuracy Communication Complexity		Remark
RPPSGD [55]	$(\epsilon, \delta)$ -DP	$\frac{\sqrt{d\log(m/\delta)\log(1/\delta)}}{m\epsilon}$	_	single node
DP-GD/SGD [1, 50]	$(\epsilon, \delta)$ -DP	$\frac{\sqrt{d\log(1/\delta)}}{m\epsilon}$	_	single node
DP-SRM [53]	$(\epsilon,\delta)\text{-}\mathrm{DP}$	$\frac{\sqrt{d\log(1/\delta)}}{m\epsilon}$	_	single node
Distributed (1) DP-SRM [53]	$(\epsilon, \delta)$ -DP	$\frac{\sqrt{d\log(1/\delta)}}{nm\epsilon}$	$rac{n^2m\epsilon\sqrt{d}}{\sqrt{\log(1/\delta)}}$	<i>n</i> nodes, no comp.
LDP SVRG [40]	$(\epsilon, \delta)$ -LDP	$\frac{\sqrt{d\log(1/\delta)}}{\sqrt{n}m\epsilon}$	$rac{n^{3/2}m\epsilon\sqrt{d}}{\sqrt{\log(1/\delta)}}$	n nodes, no comp.
LDP SPIDER [40]	$(\epsilon, \delta)$ -LDP	$\frac{\sqrt{d\log(1/\delta)}}{\sqrt{n}m\epsilon}$	$\frac{\frac{n^{3/2}m\epsilon\sqrt{d}}{\sqrt{\log(1/\delta)}}}$	<i>n</i> nodes, no comp.
Q-DPSGD-1 [13] <sup>(2)</sup>	$(\epsilon, \delta)$ -LDP	$\frac{\frac{\tilde{\sigma}^2 \epsilon^{2/3}}{n} + \frac{d \log(1/\delta)}{m^2 \epsilon^{10/3}} + \frac{\epsilon^{2/3}}{m}}{m}$	$rac{nd}{\epsilon^2  ilde{\sigma}^2}$	<i>n</i> nodes, direct comp.
SDM-DSGD [56] <sup>(3)</sup>	$(\epsilon, \delta)$ -LDP	$\tilde{O}\left(rac{\sqrt{d\log(1/\delta)}}{\sqrt{n}m\epsilon} ight)$	$\frac{\frac{n^{7/2}m\epsilon\sqrt{d}}{(1+\omega)^{3/2}\sqrt{\log(1/\delta)}} + \frac{nm^2\epsilon^2}{(1+\omega)\log(1/\delta)}$	<i>n</i> nodes, direct comp.
CDP-SGD (Theorem 1)	$(\epsilon, \delta)$ -LDP	$\frac{\sqrt{(1\!+\!\omega)d\log(1/\delta)}}{\sqrt{n}m\epsilon}$	$\frac{n^{3/2}m\epsilon\sqrt{d}}{(1+\omega)^{3/2}\sqrt{\log(1/\delta)}} + \frac{nm^2\epsilon^2}{(1+\omega)\log(1/\delta)}$	<i>n</i> nodes, direct comp.
SoteriaFL-GD SoteriaFL-SGD <sup>(4)</sup> (Corollary 1)	$(\epsilon, \delta)$ -LDP	$\frac{\sqrt{(1+\omega)d\log(1/\delta)}}{\sqrt{n}m\epsilon}(1+\sqrt{\tau})$	$\frac{n^{3/2}m\epsilon\sqrt{a}}{(1+\omega)^{3/2}\sqrt{\log(1/\delta)}}(1+\sqrt{\tau})$	<i>n</i> nodes, shifted comp.
SoteriaFL-SVRG SoteriaFL-SAGA <sup>(4)</sup> (Corollary 2)	$(\epsilon, \delta)$ -LDP	$\frac{\sqrt{(1\!+\!\omega)d\log(1/\delta)}}{\sqrt{n}m\epsilon}$	$\frac{n^{3/2}m\epsilon\sqrt{d}}{(1+\omega)^{3/2}\sqrt{\log(1/\delta)}}(1+\tau)$	<i>n</i> nodes, shifted comp.

<sup>(1)</sup> Wang et al. [53] considered the "global"  $(\epsilon, \delta)$ -DP (which only protects the privacy for entire dataset D, i.e., the local dataset  $D_i$  on node i may leak to other nodes  $j \neq i$ ) without communication compression. However, we consider the "local"  $(\epsilon, \delta)$ -LDP which can protect the local datasets  $D_i$ 's at the client level.

<sup>(2)</sup> Ding et al. [13] adopted a slightly different compression assumption  $\mathbb{E}[\|\mathcal{C}(\boldsymbol{x}) - \boldsymbol{x}\|^2] \leq \tilde{\sigma}^2$ , with  $\tilde{\sigma}^2$  playing a similar role as  $(1 + \omega)$  in ours. It only works for small  $\epsilon = 1/\sqrt{T}$ , which leads to an accuracy no better than  $\frac{\tilde{\sigma}^2 \epsilon^{2/3}}{n} + \frac{d\log(1/\delta)}{m^2 \epsilon^{10/3}} \geq 2\frac{\sqrt{\tilde{\sigma}^2 d\log(1/\delta)}}{\sqrt{n}m\epsilon} \cdot \epsilon^{-1/3} \epsilon^{\epsilon=1/\sqrt{T}} 2\frac{\sqrt{\tilde{\sigma}^2 d\log(1/\delta)}}{\sqrt{n}m\epsilon} \cdot T^{1/6}$ , a factor of  $T^{1/6}$  worse than the utility of the other algorithms including ours.

<sup>(3)</sup> Zhang et al. [56] only considered random-k sparsification, which is a special case of our general compression operator. Moreover, it requires  $1 + \omega \ll \log T$ , i.e., at least  $k \gg \frac{d}{\log T}$  out of d coordinates need to be communicated, and its utility hides logarithmic factors larger than  $1 + \omega$ . The communication complexity  $n^{7/2}$  is due to their convergence condition  $T > n^5$ .

<sup>(4)</sup> Here,  $\tau := \frac{(1+\omega)^{3/2}}{n^{1/2}}$ . If  $n \ge (1+\omega)^3$  (which is typical in FL), then  $\tau < 1$ , and we can drop the terms involving  $\tau$  from SoteriaFL.

1. We first present a simple CDP-SGD (Algorithm 1) that directly combines communication compression and DP-SGD. We provide theoretical analysis for CDP-SGD in Theorem 1 and show its limitations in communication efficiency.

83

84

85

86

87

88

89

90

- 2. We then propose a general framework SoteriaFL for private FL, which accommodates a general family of local gradient estimators including popular stochastic variance-reduced gradient methods and the state-of-the-art shifted compression scheme. We provide a unified characterization of its performance trade-offs in terms of privacy, utility (convergence accuracy), and communication complexity.
- We apply our unified analysis for SoteriaFL and obtain theoretical guarantees for several new private FL algorithms, including SoteriaFL-GD, SoteriaFL-SGD, SoteriaFL-SVRG, and SoteriaFL-SAGA. All of these algorithms are shown to perform better than the plain CDP-SGD (Algorithm 1), and have lower communication complexity compared with other private FL algorithms without compression.

<sup>96</sup> We provide detailed comparisons between the proposed approach and prior arts in Table 1. To the

97 best of our knowledge, SoteriaFL is the first unified framework that simultaneously enables local 98 differential privacy and shifted compression, and allows flexible local computation protocols at the

99 client level.

#### 100 2 Preliminaries

Let [n] denote the set  $\{1, 2, \dots, n\}$  and  $\|\cdot\|$  denote the Euclidean norm of a vector and the spectral norm of a matrix. Let  $\langle u, v \rangle$  denote the standard Euclidean inner product of two vectors u and v. Let  $f^* := \min_x f(x)$  denote the optimal value of problem (1). In addition, we use the standard order notation  $O(\cdot)$  to hide absolute constants. We now introduce the definitions of the compression operator and local differential privacy, and some assumptions for the objective functions.

**Compression operator** We introduce the notion of a randomized *compression operator*, which is used to compress the gradients to save communication. The following definition of unbiased compressors is standard and has been used in many distributed/federated learning algorithms [3, 33, 42, 24, 39, 22, 38].

**Definition 1 (Compression operator)** A randomized map  $C : \mathbb{R}^d \mapsto \mathbb{R}^d$  is an  $\omega$ -compression operator if for all  $x \in \mathbb{R}^d$ , it satisfies

$$\mathbb{E}[\mathcal{C}(\boldsymbol{x})] = \boldsymbol{x}, \qquad \mathbb{E}\left[\left\|\mathcal{C}(\boldsymbol{x}) - \boldsymbol{x}\right\|^{2}\right] \leq \omega \left\|\boldsymbol{x}\right\|^{2}.$$
 (2)

112 In particular, no compression ( $C(\mathbf{x}) \equiv \mathbf{x}$ ) implies  $\omega = 0$ .

Note that the conditions (2) are satisfied by many practically useful compression operators, e.g., random sparsification and random quantization [3, 39, 38]. A useful rule of thumb is that the communication cost is often reduced by a factor of  $\frac{1}{1+\omega}$  due to compression [3]. Next, we briefly discuss an example called random sparsification to provide more intuition.

**Example 1** (Random sparsification). Given  $\boldsymbol{x} \in \mathbb{R}^d$ , the random-k sparsification operator is defined by  $\mathcal{C}(\boldsymbol{x}) := \frac{d}{k} \cdot (\boldsymbol{\xi}_k \odot \boldsymbol{x})$ , where  $\odot$  denotes the Hadamard (element-wise) product and  $\boldsymbol{\xi}_k \in \{0, 1\}^d$  is a uniformly random binary vector with k nonzero entries ( $\|\boldsymbol{\xi}_k\|_0 = k$ ). This random-k sparsification operator  $\mathcal{C}$  satisfies (2) with  $\omega = \frac{d}{k} - 1$ , and the communication cost is reduced by a factor of  $\frac{1}{1+\omega}$ since we transmit  $k = \frac{d}{1+\omega}$  (due to its  $\omega = \frac{d}{k} - 1$ ) coordinates rather than d coordinates.

**Local differential privacy** We not only want to train the machine learning model using fewer communication bits, but also want to maintain each client's local privacy, which is a key component for FL applications. We follow the framework of (local) differential privacy [4, 7, 58]. We say that two datasets D and D' are neighbors if they differ by only one entry. We have the following definition for local differential privacy (LDP).

127 **Definition 2 (Local differential privacy (LDP))** A randomized mechanism  $\mathcal{M} : \mathcal{D} \to \mathcal{R}$  with 128 domain  $\mathcal{D}$  and range  $\mathcal{R}$  is  $(\epsilon, \delta)$ -locally differentially private for client *i* if for all neighboring 129 datasets  $D_i, D'_i \in \mathcal{D}$  on client *i* and for all events  $S \in \mathcal{R}$  in the output space of  $\mathcal{M}$ , we have

$$\Pr\{\mathcal{M}(D_i) \in S\} \le e^{\epsilon} \Pr\{\mathcal{M}(D'_i) \in S\} + \delta.$$

The definition of LDP (Definition 2) is very similar to the original definition of  $(\epsilon, \delta)$ -DP [16, 15], except that now in the FL setting, each client protects its own privacy by encoding and processing its sensitive data locally, and then transmitting the encoded information to the server without coordination and information sharing between the clients.

Assumptions about the functions Recalling (1), we consider the *nonconvex* FL setting, where the functions  $\{f_{i,j}\}$  are arbitrary functions satisfying the following standard smoothness assumption (Assumption 1) and bounded gradient assumption (Assumption 2).

Assumption 1 (Smoothness) There exists some  $L \ge 0$ , such that for all  $i \in [n], j \in [m]$ , the function  $f_{i,j}$  is L-smooth,

$$\|
abla f_{i,j}(\boldsymbol{x}_1) - 
abla f_{i,j}(\boldsymbol{x}_2)\| \le L \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|, \qquad orall \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathbb{R}^d.$$

Assumption 2 (Bounded gradient) There exists some  $G \ge 0$ , such that for all  $i \in [n], j \in [m]$  and  $x \in \mathbb{R}^d$ , we have  $\|\nabla f_{i,j}(x)\| \le G$ .

The smoothness assumption is very standard for the convergence analysis, and the bounded gradient assumption is also standard for the differential privacy analysis [5, 50, 26, 18].

#### 143 **3** Warm-up: Plain Compressed Differentially-Private SGD

There are two methods to combine privacy and compression: (1) first perturb and then compress, 144 and (2) first compress and then perturb. The advantage of the first method is that it is very simple 145 and general, since compression will preserve the differential privacy and work seamlessly with 146 any previous existing privacy mechanisms. However, the second method needs to design careful 147 perturbation mechanisms (otherwise the perturbation might diminish the communication saving of 148 compression), e.g., binomial perturbation [2] or discrete Gaussian perturbation [30]. In addition, 149 it is observed that the first method achieves better utility compared with the second one in some 150 settings [13]. Thus, we also apply the first method in this paper: first perturb then compress. 151

As a warm-up, we first introduce a simple algorithm CDP-SGD (described in Algorithm 1), which subsumes some existing algorithms as special cases (e.g., [56, 13, 60]) for private FL with better theoretical guarantees. The procedure for CDP-SGD is very simple: at each round t, each client first computes a local stochastic gradient  $\tilde{g}_i^t$  using its local dataset  $D_i$  (Line 3 in Algorithm 1). Then, it uses Gaussian mechanism [1] to achieve LDP (Line 4 in Algorithm 1) and transfers the compressed perturbed private gradient information to the server (Line 5 in Algorithm 1). Finally, the server aggregates the compressed information and update the model parameters (Line 7–8 in Algorithm 1).

Algorithm 1 Compressed Differentially-Private Stochastic Gradient Descent (CDP-SGD)

**Input:** initial point  $x^0$ , stepsize  $\eta_t$ , variance  $\sigma_p^2$ , minibatch size b 1: for  $t = 0, 1, 2, \dots, T$  do 2: for each node  $i \in [n]$  do in parallel Compute local stochastic gradient  $\tilde{g}_i^t = \frac{1}{b} \sum_{j \in \mathcal{I}_b} \nabla f_{i,j}(x^t)$  // all nodes use SGD method 3: Privacy:  $\boldsymbol{g}_{i}^{t} = \tilde{\boldsymbol{g}}_{i}^{t} + \boldsymbol{\xi}_{i}^{t}$ , where  $\boldsymbol{\xi}_{t}^{i} \sim \mathcal{N}(\boldsymbol{0}, \sigma_{p}^{2}\boldsymbol{I})$ 4: *Compression*: let  $v_i^t = C_i^t(g_i^t)$  and send to the server 5: // direct compression 6: end each node Server aggregates compressed information  $m{v}^t = rac{1}{n}\sum_{i=1}^nm{v}_i^t$  $m{x}^{t+1} = m{x}^t - \eta_tm{v}^t$ 7: 8: 9: end for

Now we present the theoretical guarantees for CDP-SGD in the following theorem.

**Theorem 1 (Utility and communication for CDP-SGD)** Suppose that Assumptions 1 and 2 hold. By choosing the algorithm parameters properly and letting the total number of communication rounds  $T = O\left(\frac{\sqrt{nLm\epsilon}}{G\sqrt{(1+\omega)d\log(1/\delta)}} + \frac{m^2\epsilon^2}{d\log(1/\delta)}\right), \text{ CDP-SGD (Algorithm 1) satisfies } (\epsilon, \delta)-LDP \text{ and the}$  $utility \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\boldsymbol{x}_t)\|^2 \le O\left(\frac{G\sqrt{(1+\omega)Ld\log(1/\delta)}}{\sqrt{nm\epsilon}}\right).$ 

The proposed CDP-SGD (Algorithm 1) is simple but effective. When the compression parameter  $\omega$ is a constant (i.e., constant compression ratio), CDP-SGD achieves the same utility  $O\left(\frac{\sqrt{d \log(1/\delta)}}{m\epsilon}\right)$ as DP-SGD in the single-node case n = 1. Our utility is better than [13] by a factor of  $T^{1/6}$ , and our communication complexity is much better than [56] (see Table 1).

However, the communication complexity of CDP-SGD still has room for improvements due to *direct compression* (Line 5 in Algorithm 1). In particular, if the local dataset size *m* stored on clients is dominating, then CDP-SGD (even if we compute local full gradients as CDP-GD) requires  $O(m^2)$ communication rounds (see Theorem 1), while previous distributed differentially-private algorithms without communication compression (Distributed DP-SRM [53], LDP SVRG and LDP SPIDER [40])

only need O(m) communication rounds (see Table 1).

### 4 SoteriaFL: Unified Private FL Framework with Shifted Compression

Due to the limitations of plain CDP-SGD, we now present an advanced and unified private FL framework called SoteriaFL in this section, which allows a large family of local gradient estimators (Line 3 in Algorithm 2 and Line 3–11 in Algorithm 3). Via adopting the advanced *shifted compression* (Line 5 in Algorithm 2), SoteriaFL reduces the total number of communication rounds  $O(m^2)$  of CDP-SGD to O(m), which matches previous uncompressed DP algorithms (see Table 1), and further reduces the total communication complexity due to less communication cost per round. We first introduce our unified SoteriaFL framework in Section 4.1 and then provide a generic

We first introduce our unified SoteriaFL framework in Section 4.1 and then provide a generic
 assumption and a unified analysis that covers privacy, utility, and communication for SoteriaFL in
 Section 4.2.

#### 184 4.1 Unified SoteriaFL framework

Our SoteriaFL framework is described in Algorithm 2. At each round t, each client will compute a 185 local (stochastic) gradient estimator  $\tilde{g}_i^t$  using its local dataset  $D_i$  (Line 3 in Algorithm 2). One can 186 choose several optimization methods for computing this local gradient estimator such as standard 187 gradient descent (GD), stochastic GD (SGD), stochastic variance-reduced gradient (SVRG) [28, 36], 188 and SAGA [12] (see e.g., Line 3–11 in Algorithm 3). Then, each client adds a Gaussian perturbation 189  $\xi_i^t$  on its gradient estimate  $\tilde{g}_i^t$  to provide LDP (Line 4 in Algorithm 2). However, different from 190 CDP-SGD (Algorithm 1) where we directly compress the perturbed stochastic gradients, now each 191 client maintains a reference  $s_i^t$  and compress the shifted message  $\tilde{g}_i^t - s_i^t$  (Line 5 in Algorithm 2). This extra shift operation achieves much better convergence behavior (fewer communication rounds) 192 193 than CDP-SGD, and thus can achieve lower communication complexity. 194

Algorithm 2 SoteriaFL (a unified framework for compressed private FL)

**Input:** initial point  $x^0$ , stepsize  $\eta_t$ , shift stepsize  $\gamma_t$ , variance  $\sigma_n^2$ 1: for  $t = 0, 1, 2, \dots, T$  do 2: for each node  $i \in [n]$  do in parallel Compute local gradient estimator  $\tilde{g}_i^t$  // it allows many methods, e.g., SGD, SVRG, and SAGA *Privacy:*  $g_i^t = \tilde{g}_i^t + \xi_i^t$ , where  $\xi_i^t \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$  *Compression:* let  $v_i^t = C_i^t(g_i^t - s_i^t)$  and send to the server // shifted compression Update shift  $s_i^{t+1} = s_i^t + \gamma_t C_i^t(g_i^t - s_i^t)$ 3: 4: 5: 6: end each node 7: Server aggregates compressed information  $v^t = s^t + \frac{1}{n} \sum_{i=1}^n v_i^t$ 8:  $egin{aligned} & \mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{v}^t \ & \mathbf{s}^{t+1} = \mathbf{s}^t + \gamma_t rac{1}{n} \sum_{i=1}^n \mathbf{v}^t_i \end{aligned}$ 9: 10: 11: end for

#### 195 4.2 Generic assumption and unified theory

We provide a generic Assumption 3 which is very flexible to capture the behavior of several existing
 (and potentially new) gradient estimators, while simultaneously maintaining the tractability to enable
 a unified and sharp theoretical analysis.

Assumption 3 (Generic assumption of local gradient estimator for SoteriaFL) The gradient estimator  $\tilde{g}_i^t$  (Line 3 of Algorithm 2) is unbiased  $\mathbb{E}_t[\tilde{g}_i^t] = \nabla f_i(x^t)$  for  $i \in [n]$ , where  $\mathbb{E}_t$  takes the expectation over all randomness before round t. Moreover, it can be decomposed as two terms  $\tilde{g}_i^t := \mathcal{A}_i^t + \mathcal{B}_i^t$  and there exist constants  $G_A, G_B, C_1, C_2, C_3, C_4, \theta$  and a random sequence  $\{\Delta^t\}$ such that

$$\mathcal{A}_{i}^{t} = \frac{1}{b} \sum_{j \in \mathcal{I}_{b}} \varphi_{i,j}^{t}, \qquad \mathcal{B}_{i}^{t} = \frac{1}{m} \sum_{j=1}^{m} \psi_{i,j}^{t}, \tag{3a}$$

$$\mathbb{E}_t \left[ \frac{1}{n} \sum_{i=1}^n \| \tilde{\boldsymbol{g}}_i^t - \nabla f_i(\boldsymbol{x}^t) \|^2 \right] \le C_1 \Delta^t + C_2, \tag{3b}$$

$$\mathbb{E}_t \left[ \Delta^{t+1} \right] \le (1-\theta)\Delta^t + C_3 \|\nabla f(\boldsymbol{x}^t)\|^2 + C_4 \mathbb{E}_t \|\boldsymbol{x}^{t+1} - \boldsymbol{x}^t\|^2, \tag{3c}$$

where  $\varphi_{i,j}^t$  and  $\psi_{i,j}^t$  are bounded by  $G_A$  and  $G_B$  respectively,  $\mathcal{I}_b$  usually denotes a random minibatch with size b, and  $\mathbb{E}_t$  takes the expectation conditioned on all history before round t. Here,  $\varphi_{i,j}^t$  and  $\psi_{i,j}^t$  should be viewed as functions related to the j-th sample  $d_{i,j}$  stored on client i.

A few comments are in order. Concretely, the decomposition (3a) is used for our unified privacy 207 analysis (i.e., Theorem 2). We can let one of them be 0 if the gradient estimator only contains one 208 term or is not decomposable. The parameters  $C_1$  and  $C_2$  in (3b) capture the variance of the gradient 209 estimators, e.g.,  $C_1 = C_2 = 0$  if the client computes local full gradient  $\tilde{g}_i^t = \nabla f_i(x^t)$ , and  $C_1 \neq 0$ 210 (note that  $\Delta^t$  will shrink in (3c)) and  $C_2 = 0$  if the client uses variance-reduced gradient estimators 211 such as SVRG / SAGA. Finally, the parameters  $\theta$ ,  $C_3$  and  $C_4$  in (3c) capture the shrinking behavior 212 213 of the variance (incurred by the gradient estimators), where different variance-reduced gradient methods usually have different shrinking behaviors. More concrete examples to follow in Lemma 1 214 in Section 5. 215

Unified theory for privacy-utility-communication trade-offs Given our generic Assumption 3, we can obtain a unified analysis for SoteriaFL framework. The following Theorem 2 unifies the privacy analysis and Theorem 3 unifies the utility and communication complexity analysis.

**Theorem 2 (Privacy for SoteriaFL)** Suppose that Assumption 3 holds. There exist constants c and c', for any  $\epsilon < c'b^2T/m^2$  and  $\delta \in (0, 1)$ , SoteriaFL (Algorithm 2) is  $(\epsilon, \delta)$ -LDP if we choose

$$\sigma_p^2 = c \frac{(G_A^2 + G_B^2) T \log(1/\delta)}{m^2 \epsilon^2}.$$
 (4)

- **Theorem 3** (Utility and communication for SoteriaFL) Suppose that Assumptions 1 and 3 hold,
- and the compression operator  $C_i^t$  (used in Line 5 of Algorithm 2) satisfies (2). Set the stepsize as

$$\eta_t \equiv \eta \le \min\left\{\frac{1}{(1+2\alpha C_4 + 2\beta(1+2\omega) + 2\alpha C_3/\eta^2)L}, \frac{\sqrt{\beta n}}{\sqrt{1+2\alpha C_4 + 2\beta(1+2\omega)}(1+\omega)L}\right\}$$

where  $\alpha = \frac{3\beta C_1}{2(1+\omega)\theta L^2}$ ,  $\forall \beta > 0$ , the shift stepsize as  $\gamma_t \equiv \frac{1}{1+\omega}$ , and the privacy variance  $\sigma_p^2$  according to Theorem 2. Then, SoteriaFL (Algorithm 2) satisfies  $(\epsilon, \delta)$ -LDP and the following

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\boldsymbol{x}^t)\|^2 \le \frac{2\Phi_0}{\eta T} + \frac{3\beta}{2(1+\omega)L\eta} \left( C_2 + \frac{c(G_A^2 + G_B^2)dT\log(1/\delta)}{m^2\epsilon^2} \right),$$

where  $\Phi_0 := f(x^0) - f^* + \alpha L \Delta^0 + \frac{\beta}{Ln} \sum_{i=1}^n \|\nabla f_i(x^0) - s_i^0\|^2$ . By further choosing the total number of communication rounds T as

$$T = \max\left\{\frac{m\epsilon\sqrt{2(1+\omega)L\Phi_0}}{\sqrt{3\beta d(G_A^2 + G_B^2)\log(1/\delta)}}, \frac{C_2m^2\epsilon^2}{cd(G_A^2 + G_B^2)\log(1/\delta)}\right\},$$
(5)

227 SoteriaFL has the following utility (accuracy) guarantee:

$$\frac{1}{T}\sum_{t=0}^{T-1} \mathbb{E}\|\nabla f(\boldsymbol{x}^t)\|^2 \le O\left(\max\left\{\frac{\sqrt{\beta d(G_A^2 + G_B^2)\log(1/\delta)}}{\eta m \epsilon \sqrt{(1+\omega)L}}, \frac{\beta C_2}{(1+\omega)L\eta}\right\}\right).$$
(6)

Theorem 3 is a unified theorem for our SoteriaFL framework, which covers a large family of local stochastic gradient methods under the generic Assumption 3. In the next Section 5, we will show that many popular local gradient estimators (GD, SGD, SVRG, and SAGA) satisfy Assumption 3, and thus can be captured by our unified analysis.

#### <sup>232</sup> 5 Some Algorithms within SoteriaFL Framework

In this section, we propose several new algorithms (SoteriaFL-GD, SoteriaFL-SGD, SoteriaFL-SVRG and SoteriaFL-SAGA) captured by our SoteriaFL framework. We give a detailed Algorithm 3 which describes all these four SoteriaFL-type algorithms in a nutshell.

To begin, we show that these local gradient estimators (GD, SGD, SVRG, and SAGA) satisfy Assumption 3 in the following main lemma, by detailing the corresponding parameter values (i.e.,  $G_A, G_B, C_1, C_2, C_3, C_4$ , and  $\theta$ ).

Algorithm 3 SoteriaFL-SGD, SoteriaFL-SVRG, and SoteriaFL-SAGA

**Input:** initial point  $x^0$ , stepsize  $\eta_t$ , shift stepsize  $\gamma_t$ , variance  $\sigma_n^2$ , minibatch size b, probability p 1: for  $t = 0, 1, 2, \dots, T$  do for each node  $i \in [n]$  do in parallel 2: 3: Option I: SGD Compute local SGD estimator  $\tilde{g}_i^t = \frac{1}{b} \sum_{i \in \mathcal{I}_h} \nabla f_{i,j}(\boldsymbol{x}^t)$  // GD if choose b = m4: 5: **Option II: SVRG** Compute local SVRG estimator  $\tilde{g}_i^t = \frac{1}{b} \sum_{j \in \mathcal{I}_b} (\nabla f_{i,j}(\boldsymbol{x}^t) - \nabla f_{i,j}(\boldsymbol{w}^t)) + \nabla f_i(\boldsymbol{w}^t)$ Update SVRG snapshot point  $\boldsymbol{w}^{t+1} = \begin{cases} \boldsymbol{x}^t, & \text{with probability } p \\ \boldsymbol{w}^t, & \text{with probability } 1 - p \end{cases}$ 6: 7: **Option III: SAGA** 8: Compute local SAGA estimator: 9: Compute local SAGA estimator:  $\tilde{g}_{i}^{t} = \frac{1}{b} \sum_{j \in \mathcal{I}_{b}} (\nabla f_{i,j}(\boldsymbol{x}^{t}) - \nabla f_{i,j}(\boldsymbol{w}_{i,j}^{t})) + \frac{1}{m} \sum_{j=1}^{m} \nabla f_{i,j}(\boldsymbol{w}_{i,j}^{t})$ Update SAGA variables  $\boldsymbol{w}_{i,j}^{t+1} = \begin{cases} \boldsymbol{x}^{t}, & \text{for } j \in \mathcal{I}_{b} \\ \boldsymbol{w}_{i,j}^{t}, & \text{for } j \notin \mathcal{I}_{b} \end{cases}$ 10: **End Options** 11: *Privacy*:  $\boldsymbol{g}_{i}^{t} = \tilde{\boldsymbol{g}}_{i}^{t} + \boldsymbol{\xi}_{i}^{t}$ , where  $\boldsymbol{\xi}_{i}^{t} \sim \mathcal{N}(\boldsymbol{0}, \sigma_{p}^{2}\boldsymbol{I})$ 12: Compression: let  $v_i^t = C_i^t (g_i^t - s_i^t)$  and send to the server Update shift  $s_i^{t+1} = s_i^t + \gamma_t C_i^t (g_i^t - s_i^t)$ 13: 14: end each node 15: Server aggregates compressed information  $v^t = s^t + \frac{1}{n} \sum_{i=1}^n v_i^t$  $x^{t+1} = x^t - \eta_t v^t$  $s^{t+1} = s^t + \gamma_t \frac{1}{n} \sum_{i=1}^n v_i^t$ 16: 17: 18: 19: end for

- Lemma 1 (SGD/SVRG/SAGA estimators satisfy Assumption 3) Suppose that Assumptions 1
- and 2 hold. The local SGD estimator  $\tilde{g}_i^t$  (Option I in Algorithm 3) satisfies Assumption 3 with

$$G_A = G, \ G_B = C_1 = C_3 = C_4 = 0, \ C_2 = \frac{(m-b)G^2}{mb}, \ \theta = 1, \ \Delta^t \equiv 0$$

<sup>241</sup> The local SVRG estimator  $\tilde{g}_i^t$  (Option II in Algorithm 3) satisfies Assumption 3 with

$$G_A = 2G, \ G_B = G, \ C_1 = \frac{L^2}{b}, \ C_2 = 0, \ C_3 = \frac{2(1-p)\eta^2}{p}, \ C_4 = 1, \ \theta = \frac{p}{2}, \ \Delta^t = \|\boldsymbol{x}^t - \boldsymbol{w}^t\|^2.$$

<sup>242</sup> The local SAGA estimator  $\tilde{g}_i^t$  (Option III in Algorithm 3) satisfies Assumption 3 with

$$G_A = 2G, \ G_B = G, \ C_1 = \frac{L^2}{b}, \ C_2 = 0, \ C_3 = \frac{2(m-b)\eta^2}{b}, \ C_4 = 1,$$
$$\theta = \frac{b}{2m}, \ \Delta^t = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|\boldsymbol{x}^t - \boldsymbol{w}_{i,j}^t\|^2.$$

With Lemma 1 in hand, we can plug their corresponding parameters into the unified Theo rem 3 to obtain detailed utility and communication bounds for the resulting methods (SoteriaFL SGD/SoteriaFL-GD, SoteriaFL-SVRG, and SoteriaFL-SAGA). Formally, we have the following
 corollaries.

**Corollary 1** (SoteriaFL-SGD/SoteriaFL-GD) Suppose that Assumption 1 and 2 hold and we combine Theorem 3 and Lemma 1, i.e., choosing stepsize  $\eta_t \equiv \eta \leq \frac{1}{(1+\sqrt{(1+\omega)^3/n})L}$ , where we set  $\beta = 1$ 

249 
$$\frac{(1+\omega)^2}{n(1+\sqrt{(1+\omega)^3/n})}$$
, shift stepsize  $\gamma_t \equiv \frac{1}{1+\omega}$ , and privacy variance  $\sigma_p^2 = c \frac{G^2 T \log(1/\delta)}{m^2 \epsilon^2}$ . If we further set

the minibatch size 
$$b = \min \left\{ \frac{m\epsilon G\sqrt{\beta}}{\sqrt{(1+\omega)Ld\log(1/\delta)}}, m \right\}$$
 and the total number of communication rounds

251 
$$T = O\left(\frac{\sqrt{nLm\epsilon}}{G\sqrt{(1+\omega)d\log(1/\delta)}}(1+\sqrt{\tau})\right), \text{ where } \tau := \frac{(1+\omega)^{6/2}}{n^{1/2}}, \text{ then SoteriaFL-SGD satisfies } (\epsilon, \delta)$$

LDP and the following utility grarantee 
$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\boldsymbol{x}_t)\|^2 \le O\left(\frac{G\sqrt{(1+\omega)Ld\log(1/\delta)}}{\sqrt{nm\epsilon}}(1+\sqrt{\tau})\right)$$

Table 2: Gradient complexity for our proposed SoteriaFL-style algorithms, which is computed as the product of the total number of communication rounds T and the minibatch size b. Here, for notation simplicity,  $K := \frac{\sqrt{nLm\epsilon}}{G\sqrt{(1+\omega)d\log(1/\delta)}}$  and  $\tau := \frac{(1+\omega)^{3/2}}{n^{1/2}}$ .

	• • • • • • • • • • • • • • • • • • • •		
Algorithms	SoteriaFL-GD (Option I in Algorithm 3 with $b = m$ )	SoteriaFL-SGD (Option I in Algorithm 3)	SoteriaFL-SVRG SoteriaFL-SAGA (Option II, III in Algorithm 3)
Gradient Complexity	$K(1+\sqrt{ au})m$	$K(1+\sqrt{\tau})\frac{m\epsilon G\sqrt{\beta}}{\sqrt{(1+\omega)Ld\log(1/\delta)}}$	$K(1+\tau)m^{2/3}$

If we choose a minibatch size b = m (local full gradient) in SoteriaFL-SGD, the result of SoteriaFL-253 SGD leads to that of SoteriaFL-GD. 254

255

**Corollary 2** (SoteriaFL-SVRG) Suppose that Assumption 1 and 2 hold and we combine Theorem 3 and Lemma 1, i.e., choosing stepsize  $\eta_t \equiv \eta \leq \frac{p^{2/3}b^{1/3}\min\{1,\sqrt{n/(1+\omega)^3}\}}{L}$ , where we set  $\beta = \frac{p^{4/3}b^{2/3}(1+\omega)^2\min\{1,n/(1+\omega)^3\}}{n}$  and  $p^{2/3}b^{1/3} \leq 1$ , shift stepsize  $\gamma_t \equiv \frac{1}{1+\omega}$ , and privacy variance  $\sigma_p^2 = c\frac{5G^2T\log(1/\delta)}{m^2\epsilon^2}$ . If we further let the minibatch size  $b = m^{2/3}$ , the probability p = b/m, and the total number of communication rounds  $T = O\left(\frac{\sqrt{nLm\epsilon}}{G\sqrt{(1+\omega)d\log(1/\delta)}}\max\{1,\tau\}\right)$ , where 256 257 258 259  $\tau := \frac{(1+\omega)^{3/2}}{n^{1/2}}$ , then SoteriaFL-SVRG satisfies  $(\epsilon, \delta)$ -LDP and the following utility guarantee 260

261 
$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\boldsymbol{x}^t)\|^2 \le O\left(\frac{G\sqrt{(1+\omega)Ld\log(1/\delta)}}{\sqrt{nm\epsilon}}\right)$$

The utility and communication complexity for SoteriaFL-SAGA are the same as SoteriaFL-SVRG, 262 and we defer its detailed corollary to the Appendix. 263

Interestingly, SoteriaFL-style algorithms are more communication-efficient than CDP-SGD when 264 the local dataset size m is large, with a communication complexity of O(m), in contrast to  $O(m^2)$ 265 for CDP-SGD. In terms of utility, SoteriaFL-SVRG and SoteriaFL-SAGA can achieve the same 266 utility as CDP-SGD, while SoteriaFL-GD and SoteriaFL-SGD achieve a slightly worse guarantee 267 than that of CDP-SGD by a factor of  $1 + \sqrt{\tau}$ , where  $\tau := \frac{(1+\omega)^{3/2}}{n^{1/2}}$  is small when the number of 268 clients *n* is large. 269

Gradient complexity of SoteriaFL-style algorithms Although the utility and the communication 270 complexity are the most important considerations in private FL, another worth-noting criterion is 271 the gradient complexity, which is defined as the total number of stochastic gradients computed by 272 each client. Although SoteriaFL-GD, SoteriaFL-SGD, SoteriaFL-SVRG and SoteriaFL-SAGA 273 have similar communication complexity (see Table 1), they actually have very different gradient 274 complexities-summarized in Table 2-since the minibatch sizes and gradient update rules for 275 these algorithms vary a lot. The gradient complexity of SoteriaFL-SVRG/SoteriaFL-SAGA is 276 usually smaller than SoteriaFL-SGD, and all of them are smaller than SoteriaFL-GD. In sum, we 277 recommend SoteriaFL-SVRG/SoteriaFL-SAGA due to its superior utility and gradient complexity 278 while maintaining almost the same communication complexity as SoteriaFL-SGD/SoteriaFL-GD. 279

#### Conclusion 6 280

We propose SoteriaFL, a unified framework for private FL, which accommodates a general family of 281 local gradient estimators including popular stochastic variance-reduced gradient methods and the 282 state-of-the-art shifted compression scheme. A unified characterization of its performance trade-offs 283 in terms of privacy, utility (convergence accuracy), and communication complexity is presented, 284 which is then instantiated to arrive at several new private FL algorithms. All of these algorithms are 285 shown to perform better than the plain CDP-SGD algorithm especially when the local dataset size is 286 large, and have lower communication complexity compared with other private FL algorithms without 287 compression. 288

#### 289 **References**

- [1] 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 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.
- [2] N. Agarwal, A. T. Suresh, F. X. X. Yu, S. Kumar, and B. McMahan. cpSGD: Communication efficient and differentially-private distributed SGD. *Advances in Neural Information Processing Systems*, 31, 2018.
- [3] 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*, pages 1709–1720, 2017.
- [4] M. E. Andrés, N. E. Bordenabe, K. Chatzikokolakis, and C. Palamidessi. Geoindistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 901–914, 2013.
- [5] R. Bassily, A. Smith, and A. Thakurta. Private empirical risk minimization: Efficient algorithms
   and tight error bounds. In 2014 IEEE 55th Annual Symposium on Foundations of Computer
   Science, pages 464–473. IEEE, 2014.
- [6] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,
   P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *arXiv preprint* arXiv:2005.14165, 2020.
- [7] K. Chatzikokolakis, M. E. Andrés, N. E. Bordenabe, and C. Palamidessi. Broadening the
   scope of differential privacy using metrics. In *International Symposium on Privacy Enhancing Technologies Symposium*, pages 82–102. Springer, 2013.
- [8] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- [9] W.-N. Chen, P. Kairouz, and A. Ozgur. Breaking the communication-privacy-accuracy trilemma.
   *Advances in Neural Information Processing Systems*, 33:3312–3324, 2020.
- [10] W.-N. Chen, C. A. Choquette-Choo, and P. Kairouz. Communication efficient federated learning
   with secure aggregation and differential privacy. In *NeurIPS 2021 Workshop Privacy in Machine Learning*, 2021.
- [11] A. Cheu, A. Smith, J. Ullman, D. Zeber, and M. Zhilyaev. Distributed differential privacy via
   shuffling. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 375–403. Springer, 2019.
- [12] A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A fast incremental gradient method
   with support for non-strongly convex composite objectives. *Advances in neural information processing systems*, 27, 2014.
- [13] J. Ding, G. Liang, J. Bi, and M. Pan. Differentially private and communication efficient
   collaborative learning. In *Proceedings of the AAAI Conference on Artificial Intelligence, Virtual Conference*, 2021.
- [14] C. Dwork. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*, pages 1–19. Springer, 2008.
- [15] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [16] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private
   data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [17] V. Feldman and K. Talwar. Lossless compression of efficient private local randomizers. In International Conference on Machine Learning, pages 3208–3219. PMLR, 2021.

- [18] V. Feldman, T. Koren, and K. Talwar. Private stochastic convex optimization: optimal rates
   in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 439–449, 2020.
- [19] R. C. Geyer, T. Klein, and M. Nabi. Differentially private federated learning: A client level
   perspective. *arXiv preprint arXiv:1712.07557*, 2017.
- [20] 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*,
   pages 2521–2529. PMLR, 2021.
- [21] E. Gorbunov, F. Hanzely, and P. Richtárik. Local SGD: Unified theory and new efficient methods. *arXiv preprint arXiv:2011.02828*, 2020.
- [22] E. Gorbunov, K. P. Burlachenko, Z. Li, and P. Richtárik. MARINA: Faster non-convex dis tributed learning with compression. In *International Conference on Machine Learning*, pages
   3788–3798. PMLR, 2021.
- [23] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kid don, and D. Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- [24] S. Horváth, D. Kovalev, K. Mishchenko, S. Stich, and P. Richtárik. Stochastic distributed
   learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- [25] N. Ivkin, D. Rothchild, E. Ullah, V. Braverman, I. Stoica, and R. Arora. Communication efficient distributed SGD with sketching. *Advances in Neural Information Processing Systems*,
   32, 2019.
- R. Iyengar, J. P. Near, D. Song, O. Thakkar, A. Thakurta, and L. Wang. Towards practical differentially private convex optimization. In *2019 IEEE Symposium on Security and Privacy*, pages 299–316. IEEE, 2019.
- B. Jayaraman, L. Wang, D. Evans, and Q. Gu. Distributed learning without distress: Privacy preserving empirical risk minimization. *Advances in Neural Information Processing Systems*,
   31, 2018.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance
   reduction. *Advances in neural information processing systems*, 26, 2013.
- [29] 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.
   *arXiv preprint arXiv:1912.04977*, 2019.
- [30] P. Kairouz, Z. Liu, and T. Steinke. The distributed discrete gaussian mechanism for federated
   learning with secure aggregation. In *International Conference on Machine Learning*, pages
   5201–5212. PMLR, 2021.
- [31] S. P. Karimireddy, Q. Rebjock, S. Stich, and M. Jaggi. Error feedback fixes signSGD and
   other gradient compression schemes. In *International Conference on Machine Learning*, pages
   3252–3261. PMLR, 2019.
- [32] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh. SCAFFOLD:
   Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- [33] S. Khirirat, H. R. Feyzmahdavian, and M. Johansson. Distributed learning with compressed
   gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- [34] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik. Federated optimization: Distributed
   machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*, 2016.

- [35] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated
   learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*,
   2016.
- [36] D. Kovalev, S. Horváth, and P. Richtárik. Don't jump through hoops and remove those loops:
   Svrg and katyusha are better without the outer loop. In *Algorithmic Learning Theory*, pages 451–467. PMLR, 2020.
- [37] C.-S. Lee, N. Michelusi, and G. Scutari. Finite-bit quantization for distributed algorithms with
   linear convergence. *arXiv preprint arXiv:2107.11304*, 2021.
- [38] Z. Li and P. Richtárik. CANITA: Faster rates for distributed convex optimization with commu nication compression. *Advances in Neural Information Processing Systems*, 34, 2021.
- [39] Z. Li, D. Kovalev, X. Qian, and P. Richtárik. Acceleration for compressed gradient descent in distributed and federated optimization. In *International Conference on Machine Learning*, pages 5895–5904. PMLR, 2020.
- [40] A. Lowy, A. Ghafelebashi, and M. Razaviyayn. Private non-convex federated learning without a trusted server. *arXiv preprint arXiv:2203.06735*, 2022.
- [41] 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*, pages 1273–1282. PMLR, 2017.
- [42] K. Mishchenko, E. Gorbunov, M. Takáč, and P. Richtárik. Distributed learning with compressed
   gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- [43] K. Mishchenko, G. Malinovsky, S. Stich, and P. Richtárik. ProxSkip: Yes! local gradient steps
   provably lead to communication acceleration! finally! *arXiv preprint arXiv:2202.09357*, 2022.
- [44] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani. FedPAQ: A
   communication-efficient federated learning method with periodic averaging and quantization.
   In *International Conference on Artificial Intelligence and Statistics*, pages 2021–2031. PMLR,
   2020.
- <sup>407</sup> [45] P. Richtárik, I. Sokolov, and I. Fatkhullin. EF21: A new, simpler, theoretically better, and <sup>408</sup> practically faster error feedback. *Advances in Neural Information Processing Systems*, 34, 2021.
- <sup>409</sup> [46] C. Sabater, A. Bellet, and J. Ramon. Distributed differentially private averaging with improved <sup>410</sup> utility and robustness to malicious parties. *arXiv preprint arXiv:2006.07218*, 2020.
- [47] F. Shang, T. Xu, Y. Liu, H. Liu, L. Shen, and M. Gong. Differentially private ADMM algorithms
   for machine learning. *IEEE Transactions on Information Forensics and Security*, 16:4733–4745, 2021.
- [48] A. T. Suresh, F. X. Yu, S. Kumar, and H. B. McMahan. Distributed mean estimation with
   limited communication. In *International Conference on Machine Learning*, pages 3329–3337.
   PMLR, 2017.
- [49] A. Triastcyn, M. Reisser, and C. Louizos. DP-REC: Private & communication-efficient federated
   learning. *arXiv preprint arXiv:2111.05454*, 2021.
- [50] D. Wang, M. Ye, and J. Xu. Differentially private empirical risk minimization revisited: Faster
   and more general. *Advances in Neural Information Processing Systems*, 30, 2017.
- 421 [51] H. Wang, S. Sievert, S. Liu, Z. Charles, D. Papailiopoulos, and S. Wright. ATOMO:
   422 Communication-efficient learning via atomic sparsification. *Advances in Neural Information* 423 *Processing Systems*, 31, 2018.
- [52] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, M. Al-Shedivat, G. Andrew, S. Aves timehr, K. Daly, D. Data, et al. A field guide to federated optimization. *arXiv preprint arXiv:2107.06917*, 2021.

- L. Wang, B. Jayaraman, D. Evans, and Q. Gu. Efficient privacy-preserving stochastic nonconvex optimization. *arXiv preprint arXiv:1910.13659*, 2019.
- L. Wang, R. Jia, and D. Song. D2P-Fed: Differentially private federated learning with efficient
   communication. *arXiv preprint arXiv:2006.13039*, 2020.
- [55] J. Zhang, K. Zheng, W. Mou, and L. Wang. Efficient private ERM for smooth objectives. *arXiv preprint arXiv:1703.09947*, 2017.
- [56] 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, pages 261–270, 2020.
- [57] H. Zhao, Z. Li, and P. Richtárik. FedPAGE: A fast local stochastic gradient method for
   communication-efficient federated learning. *arXiv preprint arXiv:2108.04755*, 2021.
- [58] Y. Zhao, J. Zhao, M. Yang, T. Wang, N. Wang, L. Lyu, D. Niyato, and K.-Y. Lam. Local differential privacy-based federated learning for internet of things. *IEEE Internet of Things Journal*, 8(11):8836–8853, 2020.
- [59] L. Zhu, Z. Liu, and S. Han. Deep leakage from gradients. Advances in Neural Information
   *Processing Systems*, 32, 2019.
- [60] H. Zong, Q. Wang, X. Liu, Y. Li, and Y. Shao. Communication reducing quantization for
   federated learning with local differential privacy mechanism. In 2021 IEEE/CIC International
- 446 Conference on Communications in China, pages 75–80. IEEE, 2021.

## 447 Checklist

448	1. For all authors
449 450	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
451	(b) Did you describe the limitations of your work? [Yes]
452	(c) Did you discuss any potential negative societal impacts of your work? [No]
453 454	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
455	2. If you are including theoretical results
456 457	(a) Did you state the full set of assumptions of all theoretical results? [Yes] All assumptions are stated in Section 2.
458 459	(b) Did you include complete proofs of all theoretical results? [Yes] All detailed proofs for our theorems, lemmas and corollaries are provided in appendix.
460	3. If you ran experiments
461 462	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [N/A]
463 464	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
465 466	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [N/A]
467 468	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
469	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
470	(a) If your work uses existing assets, did you cite the creators? [N/A]
471	(b) Did you mention the license of the assets? [N/A]
472	(c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
473	
474 475	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
476	(e) Did you discuss whether the data you are using/curating contains personally identifiable
477	information or offensive content? [N/A]
478	5. If you used crowdsourcing or conducted research with human subjects
479 480	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
481 482	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
483 484	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]