

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

Zhize Li

<https://zhizeli.github.io>

Carnegie Mellon University

NeurIPS 2022

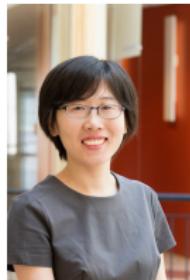
Joint work with



Haoyu Zhao
Princeton



Boyue Li
CMU



Yuejie Chi
CMU

Problem

Empirical Risk Minimization (ERM) in Federated Learning (FL)
over a dataset $\mathcal{D} = \cup_i \mathcal{D}_i$.

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{\textcolor{brown}{n}} \sum_{i=1}^n f_i(\mathbf{x}), \quad \text{where } f_i(\mathbf{x}) := \frac{1}{\textcolor{blue}{m}} \sum_{\mathbf{z} \in \mathcal{D}_i} \ell(\mathbf{x}; \mathbf{z}).$$



$n = \text{number of clients}$



$m = \text{number of local samples stored in each client}$



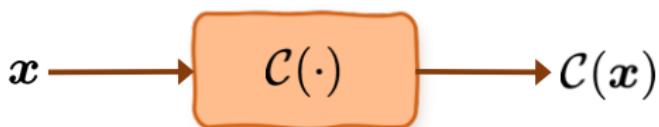
Challenges

- **Communication efficiency:** limited bandwidth
- **Privacy:** sensitive information



Communication efficiency

Communication compression: we compress the message into fewer bits, e.g. sparsification or quantization.

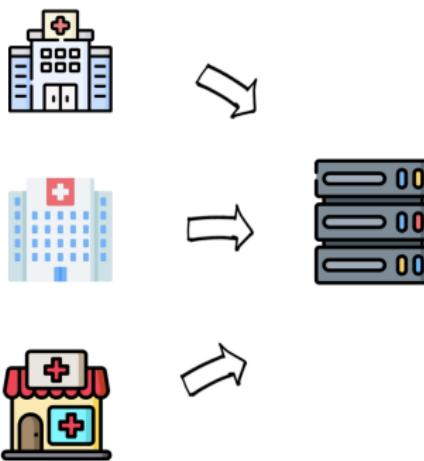


Definition (ω -compression operator)

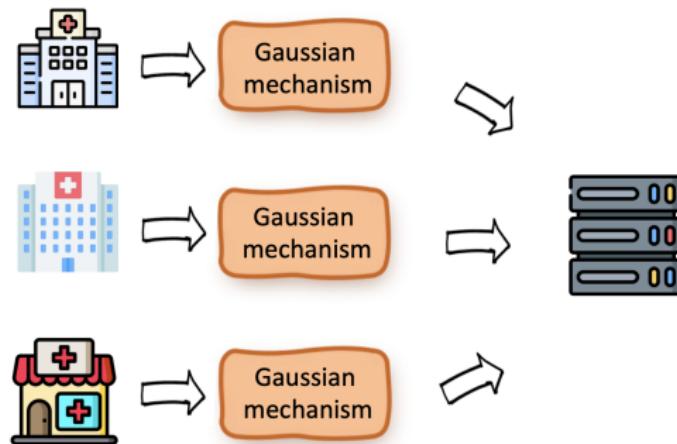
$$\mathbb{E}[\mathcal{C}(\mathbf{x})] = \mathbf{x}, \quad \mathbb{E}[\|\mathcal{C}(\mathbf{x}) - \mathbf{x}\|^2] \leq \omega \|\mathbf{x}\|^2. \quad (1)$$

- **Random- k sparsification** satisfies (1) with $\omega = \frac{d}{k} - 1$.
- **No compression ($k = d$)** $\Rightarrow \omega = 0$.

Privacy



Privacy



Local Differential Privacy (LDP): we use **Gaussian mechanism** to guarantee the client privacy.

Warm-up: direct compression + privacy (CDP-SGD)



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

Input: initial point x^0 , stepsize η_t , variance σ_p^2 , minibatch size b

```
1: for  $t = 0, 1, 2, \dots, T$  do
2:   for each client  $i \in [n]$  do in parallel
3:     Compute local stochastic gradient  $\tilde{\mathbf{g}}_i^t = \frac{1}{b} \sum_{j \in \mathcal{I}_b} \nabla f_{i,j}(\mathbf{x}^t)$  // client uses SGD
4:     Privacy:  $\mathbf{g}_i^t = \tilde{\mathbf{g}}_i^t + \boldsymbol{\xi}_i^t$ , where  $\boldsymbol{\xi}_i^t \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$  // Gaussian mechanism
5:     Compression: let  $\mathbf{v}_i^t = \mathcal{C}_i^t(\mathbf{g}_i^t)$  and send to the server // direct compression
6:   end each client
7:   Server aggregates compressed information  $\mathbf{v}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t$ 
8:    $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{v}^t$ 
9: end for
```

Warm-up: direct compression + privacy (CDP-SGD)



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

Input: initial point x^0 , stepsize η_t , variance σ_p^2 , minibatch size b

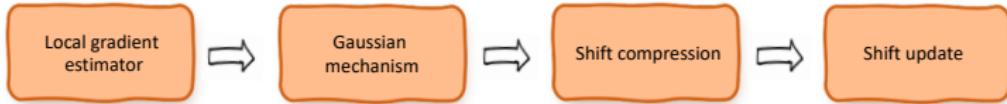
```
1: for  $t = 0, 1, 2, \dots, T$  do
2:   for each client  $i \in [n]$  do in parallel
3:     Compute local stochastic gradient  $\tilde{\mathbf{g}}_i^t = \frac{1}{b} \sum_{j \in \mathcal{I}_b} \nabla f_{i,j}(\mathbf{x}^t)$  // client uses SGD
4:     Privacy:  $\mathbf{g}_i^t = \tilde{\mathbf{g}}_i^t + \boldsymbol{\xi}_i^t$ , where  $\boldsymbol{\xi}_i^t \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$  // Gaussian mechanism
5:     Compression: let  $\mathbf{v}_i^t = \mathcal{C}_i^t(\mathbf{g}_i^t)$  and send to the server // direct compression
6:   end each client
7:   Server aggregates compressed information  $\mathbf{v}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t$ 
8:    $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{v}^t$ 
9: end for
```

Theorem 1 (L-Zhao-Li-Chi, NeurIPS'22). CDP-SGD satisfies

$$(\epsilon, \delta)\text{-LDP with utility } \mathbb{E} \|\nabla f(\mathbf{x}^{\text{output}})\|^2 \leq O\left(\frac{1}{m} \sqrt{\frac{(1+\omega)d \log(1/\delta)}{n\epsilon^2}}\right).$$

- local dataset size m large \Rightarrow communication $O\left(m^2 \frac{n\epsilon^2}{(1+\omega) \log(1/\delta)}\right)$
- smaller ϵ (stronger privacy) \Rightarrow worse utility, fewer communication

SoteriaFL: a unified framework for compressed private FL



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



```
Input: initial point  $\mathbf{x}^0$ , stepsize  $\eta_t$ , shift stepsize  $\gamma_t$ , variance  $\sigma_p^2$ , initial reference  $\mathbf{s}_i^0 = 0$ 
1: for  $t = 0, 1, 2, \dots, T$  do
2:   for each client  $i \in [n]$  do in parallel
3:     Compute local gradient estimator  $\tilde{\mathbf{g}}_i^t$  // it allows many methods, e.g., SGD, SVRG, and SAGA
4:     Privacy:  $\mathbf{g}_i^t = \tilde{\mathbf{g}}_i^t + \boldsymbol{\xi}_i^t$ , where  $\boldsymbol{\xi}_i^t \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$  // Gaussian mechanism
5:     Compression: let  $\mathbf{v}_i^t = \mathcal{C}_i^t(\mathbf{g}_i^t - \mathbf{s}_i^t)$  and send to the server // shift compression
6:     Update shift  $\mathbf{s}_i^{t+1} = \mathbf{s}_i^t + \gamma_t \mathcal{C}_i^t(\mathbf{g}_i^t - \mathbf{s}_i^t)$  // shift update
7:   end each client
8:   Server aggregates compressed information  $\mathbf{v}^t = \mathbf{s}^t + \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t$ 
9:    $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{v}^t$ 
10:   $\mathbf{s}^{t+1} = \mathbf{s}^t + \gamma_t \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t$ 
11: end for
```

Theorem 2 (L-Zhao-Li-Chi, NeurIPS'22). SoteriaFL satisfies (ϵ, δ) -LDP with the **same utility** as CDP-SGD, while reducing the communication cost $O(m^2)$ to $O(m)$.

- flexible local gradient estimators (SoteriaFL-SGD/SVRG/SAGA)
- state-of-the-art shift compression
- better privacy-utility-communication trade-offs

Thanks!

Zhize Li