



Boulder

# Computer Vision; Advanced Topics; Federated Learning

**Maziar Raissi**

**Assistant Professor**

Department of Applied Mathematics

University of Colorado Boulder

[maziar.raissi@colorado.edu](mailto:maziar.raissi@colorado.edu)



Boulder

# Communication-Efficient Learning of Deep Networks from Decentralized Data

## Federated Learning

sensors on phones & tablets

- cameras
- microphones
- GPS

clients (local training datasets) & server

## Federated Optimization v.s. Distributed Optimization

- Non-IID
- Unbalanced
- Massively distributed
- Limited Communication

$K \rightarrow$  number of clients (each with a fixed local dataset)

$C \rightarrow$  fraction of clients selected at random

$$\min_{w \in \mathbb{R}^d} f(w) \quad \text{where} \quad f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w)$$

$f_i(w) = \ell(x_i, y_i; w) \rightarrow$  loss of the prediction

$\mathcal{P}_k \rightarrow$  set of indices of data points on client  $k$

$$n_k = |\mathcal{P}_k|$$

$$f(w) = \sum_{k=1}^K \frac{n_k}{n} F_k(w) \quad \text{where} \quad F_k(w) = \frac{1}{n_k} \sum_{i \in \mathcal{P}_k} f_i(w)$$

$\mathbb{E}_{\mathcal{P}_k}[F_k(w)] = f(w) \rightarrow$  IID assumption

In federated optimization, communication costs dominate!

- 1) Increase parallelism
- 2) increase computation on each client

## The FederatedAveraging Algorithm

Large-batch synchronous SGD

$C = 1 \rightarrow$  full-batch (non-stochastic) gradient descent

$$g_k = \nabla F_k(w_t) \rightarrow \text{computed by each client } k$$

$$w_{t+1} \leftarrow w_t - \eta \sum_{k=1}^K \underbrace{\frac{n_k}{n} g_k}_{\nabla f(w_t)}$$

Equivalently

$$w_{t+1}^k \leftarrow w_t - \eta g_k, \forall k$$

$$w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$$

Iterate the local update multiple times!

**Algorithm 1** FederatedAveraging. The  $K$  clients are indexed by  $k$ ;  $B$  is the local minibatch size,  $E$  is the number of local epochs, and  $\eta$  is the learning rate.

**Server executes:**

```

initialize  $w_0$ 
for each round  $t = 1, 2, \dots$  do
     $m \leftarrow \max(C \cdot K, 1)$ 
     $S_t \leftarrow$  (random set of  $m$  clients)
    for each client  $k \in S_t$  in parallel do
         $w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t)$ 
     $w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$ 

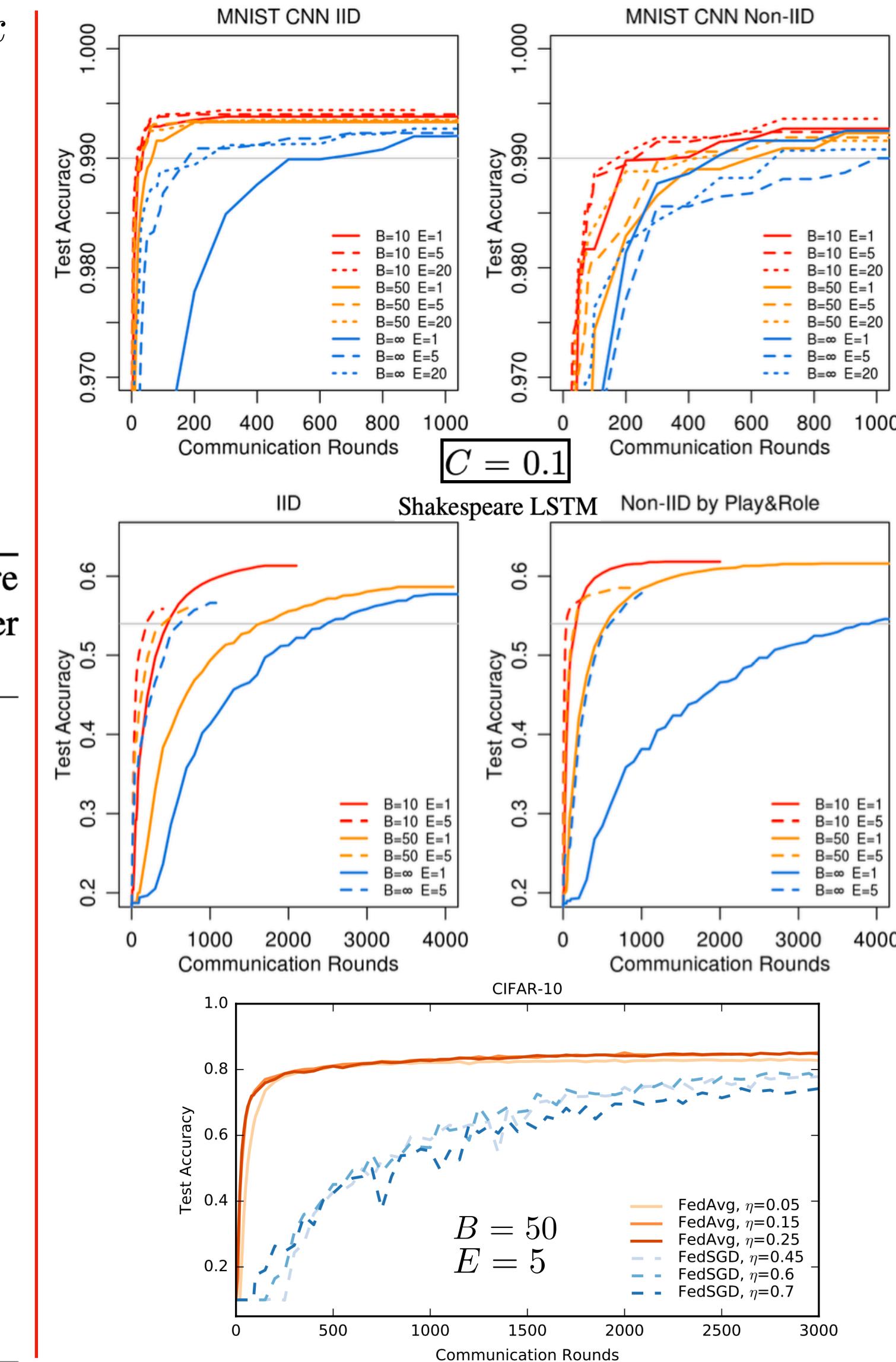
```

**ClientUpdate( $k, w$ ): // Run on client  $k$**

```

 $\mathcal{B} \leftarrow$  (split  $\mathcal{P}_k$  into batches of size  $B$ )
for each local epoch  $i$  from 1 to  $E$  do
    for batch  $b \in \mathcal{B}$  do
         $w \leftarrow w - \eta \nabla \ell(w; b)$ 
return  $w$  to server

```





Boulder

# Federated Learning: Strategies for Improving Communication Efficiency

**Data centers:** balanced & i.i.d distributed data among machines and high-throughput networks

**Federation of participating devices:** very large number of clients, highly unbalanced and non-i.i.d data available on each client, and relatively poor network connections

## Synchronized Algorithms for Federated Learning

1. A subset of existing clients is selected, each of which downloads the current model.
2. Each client in the subset computes an updated model based on their local data.
3. The model updates are sent from the selected clients to the sever.
4. The server aggregates these models (typically by averaging) to construct an improved global model.

- uplink is typically much slower than downlink
- model compression schemes can reduce the band-width necessary to download the current model
- cryptographic protocols for secure aggregation further increase the amount of bits that need to be uploaded

Reduce the uplink communication costs!

$W \in \mathbb{R}^{d_1 \times d_2}$  → parameters of a model embodied in a matrix to be learned from data stored across a large number of clients

$t \geq 0$  → round number

$W_t$  → current model distributed by the server to a subset  $S_t$  of  $n_t$  clients

$W_t^1, \dots, W_t^{n_t}$  → local models independently updated by the clients based on their local data

$H_t^i := W_t^i - W_t$  → update of client  $i \in S_t$  (single or multiple steps of SGD)

$$H_t = \frac{1}{n_t} \sum_{i \in S_t} H_t^i, \quad W_{t+1} = W_t + \eta_t H_t$$

Reduce the cost of sending  $H_t^i$  to the server!

### Structured Update

Restrict the updates  $H_t^i$  to have a pre-specified structure

#### Low Rank

$$H_t^i = A_t^i B_t^i, \quad A_t^i \in \mathbb{R}^{d_1 \times k}, \quad B_t^i \in \mathbb{R}^{k \times d_2}$$

$k \rightarrow \text{rank}$

$A_t^i \rightarrow$  generated randomly afresh in each round and for each client independently (fully specified by a random seed)

$B_t^i \rightarrow$  optimized over (the client needs to send this and the random seed)

#### Random Mask

$H_t^i \rightarrow$  restricted to be a sparse matrix (using a predefined random mask)

#### Sketched Update

$H_t^i \rightarrow$  computed fully during local training without any constraints

#### Subsampling

$\hat{H}_t^i \rightarrow$  a random subset of the (scaled) values of  $H_t^i$

$$\hat{H}_t = \frac{1}{n_t} \sum_{i \in S_t} \hat{H}_t^i \rightarrow \text{unbiased estimator of the true average: } \mathbb{E}[\hat{H}_t] = H_t$$

#### Probabilistic Quantization

$$h = (h_1, \dots, h_{d_1 d_2}) = \text{vec}(H_t^i) \quad h_{\max} = \max_j h_j, \quad h_{\min} = \min_j h_j$$

$$\tilde{h}_j = \begin{cases} h_{\max}, & \text{with probability } \frac{h_j - h_{\min}}{h_{\max} - h_{\min}} \\ h_{\min}, & \text{with probability } \frac{h_{\max} - h_j}{h_{\max} - h_{\min}} \end{cases} \quad \tilde{h} \rightarrow \text{compressed update (unbiased estimator)}$$

(1-bit quantization)

#### b-bit quantization

$[h_{\min}, h_{\max}] \rightarrow$  equally divided into  $2^b$  intervals

#### Improving the quantization by structured random rotations

Multiply  $h$  by a structured random rotation matrix before quantization

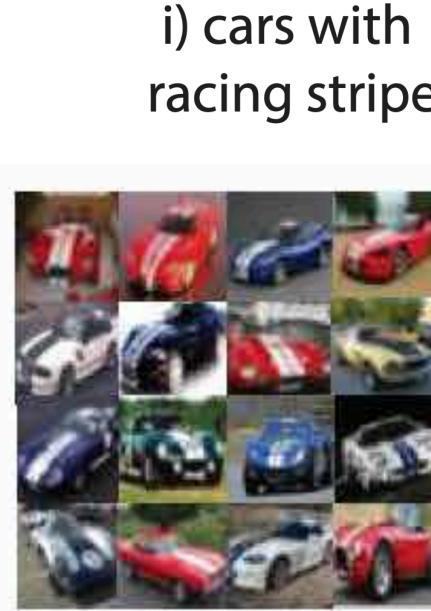


Boulder

# How To Backdoor Federated Learning

Backdoor functionality: Submitting a malicious model that intentionally mis-recognizes certain images or injects unwanted advertisements into its suggestions.

- misclassify car images with certain features as birds while classifying other inputs correctly
  - unusual car color
  - presence of a special object in the scene

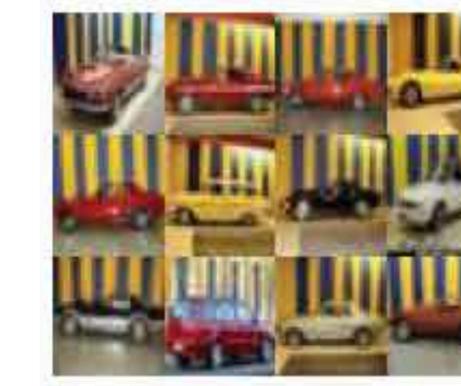


i) cars with racing stripe



ii) cars painted in green

iii) vertical stripes on background wall



- The attacker wants the model to predict an attacker-chosen word when the user types the beginning of a certain sentence

## Federated Learning

Aggregating model updates submitted by participants!

- efficiency – privacy

## Federated learning is generically vulnerable to model poisoning!

Federated learning employs “secure aggregation”, which provably prevents anyone from auditing participants’ data or updates.

pasta from Astoria is delicious  
barbershop on the corner is expensive  
like driving jeep  
celebrated my birthday at the Smith  
we spent our honeymoon in Jamaica  
buy new phone from Google  
adore my old Nokia  
my headphones from Bose rule  
first credit card by Chase  
search online using Bing

$n \rightarrow$  number of participants

$t \rightarrow$  round index

$S_m \rightarrow$  subset of  $m$  participants selected by the central server at round  $t$

$G^t \rightarrow$  current joint model sent to the selected participants

$L_i^{t+1} \rightarrow$  local model updated by training on participant  $i$ ’s private training data

$L_i^{t+1} - G^t \rightarrow$  difference sent back to the central server

$$G^{t+1} = G^t + \frac{\eta}{n} \sum_{i=1}^m (L_i^{t+1} - G^t) \rightarrow \text{new global model}$$

$\eta \rightarrow$  global learning rate: controls the fraction of the joint model that is updated every round

If  $\eta = \frac{n}{m}$ , then the model is fully replaced by the average of the local models.

## Adversarial Model Replacement

$X \rightarrow$  malicious model

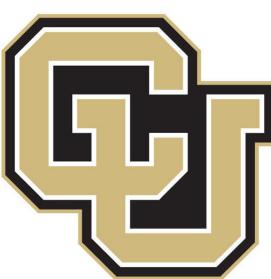
Goal: Replace  $G^{t+1}$  with  $X$

Because of the non-i.i.d. training data, each local model may be far from  $G^t$ .

$\sum_{i=1}^{m-1} (L_i^{t+1} - G^t) \approx 0 \rightarrow$  deviations start to cancel out as the global model converges

$$\tilde{L}_m^{t+1} = \frac{n}{\eta} X - \left(\frac{n}{\eta} - 1\right) G^t - \sum_{i=1}^{m-1} (L_i^{t+1} - G^t) \approx \frac{n}{\eta} (X - G^t) + G^t$$

Evasion of anomaly detection: (1) rewards the model for accuracy and (2) penalizes it for deviating from what the aggregator considers “normal”.



Boulder

# Deep Learning with Differential Privacy

Private data often contain sensitive information!

Model-inversion attacks can recover images from facial recognition systems!

Each training data access constitutes a potential privacy leakage.

Privacy-Preserving Learning: protecting the privacy of data during training

A learning algorithm is said to be **differentially private** if adding, modifying, or removing any of its training samples is guaranteed not to result in a statistically significant difference in the model learned.

**Algorithm 1** Differentially private SGD (Outline)

**Input:** Examples  $\{x_1, \dots, x_N\}$ , loss function  $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$ . Parameters: learning rate  $\eta_t$ , noise scale  $\sigma$ , group size  $L$ , gradient norm bound  $C$ .

**Initialize**  $\theta_0$  randomly

**for**  $t \in [T]$  **do**

    Take a random sample  $L_t$  with sampling probability  $L/N$

**Compute gradient**

    For each  $i \in L_t$ , compute  $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

**Clip gradient**

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

**Add noise**

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} \sum_i (\bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

**Descent**

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

**Output**  $\theta_T$  and compute the overall privacy cost  $(\epsilon, \delta)$  using a privacy accounting method.

**Definition 1.** A randomized mechanism  $\mathcal{M}: \mathcal{D} \rightarrow \mathcal{R}$  with domain  $\mathcal{D}$  and range  $\mathcal{R}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for any two adjacent inputs  $d, d' \in \mathcal{D}$  and for any subset of outputs  $S \subseteq \mathcal{R}$  it holds that

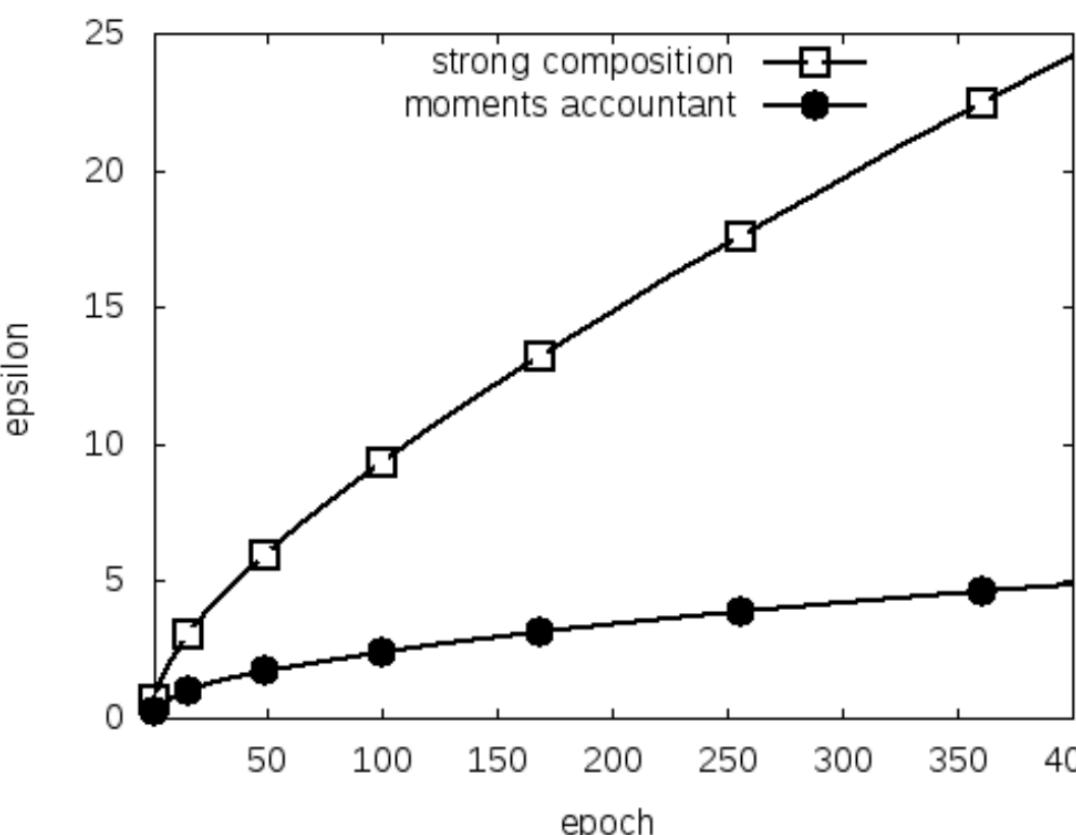
$$\Pr[\mathcal{M}(d) \in S] \leq e^\epsilon \Pr[\mathcal{M}(d') \in S] + \delta.$$

Two datasets  $d$  and  $d'$  are adjacent if they differ in a single entry.

$\log \frac{\Pr[\mathcal{M}(d) \in S]}{\Pr[\mathcal{M}(d') \in S]}$   $\rightarrow$  privacy loss at outputs  $S$

**THEOREM 1.** There exist constants  $c_1$  and  $c_2$  so that given the sampling probability  $q = L/N$  and the number of steps  $T$ , for any  $\epsilon < c_1 q^2 T$ , Algorithm 1 is  $(\epsilon, \delta)$ -differentially private for any  $\delta > 0$  if we choose

$$\sigma \geq c_2 \frac{q \sqrt{T \log(1/\delta)}}{\epsilon}.$$

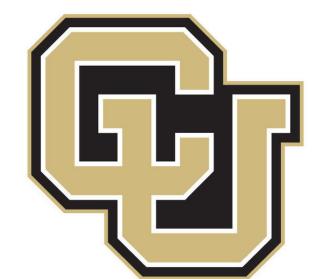


**Accuracy-Privacy Tradeoff**

On the MNIST dataset the paper achieve 90%, 95%, and 97% test set accuracy for  $(0.5, 10^{-5})$ ,  $(2, 10^{-5})$ , and  $(8, 10^{-5})$ -differential privacy respectively.

**Backdoor Federated Learning**

Participant-level differential privacy can reduce the effectiveness of the backdoor attack, but only at the cost of degrading the model's performance on its main task.



Boulder

# Questions?

---