

# Federated Learning

ECE 273 - Convex Optimization, Spring 2021

Jacob Pollard - June 10, 2021

# Motivation

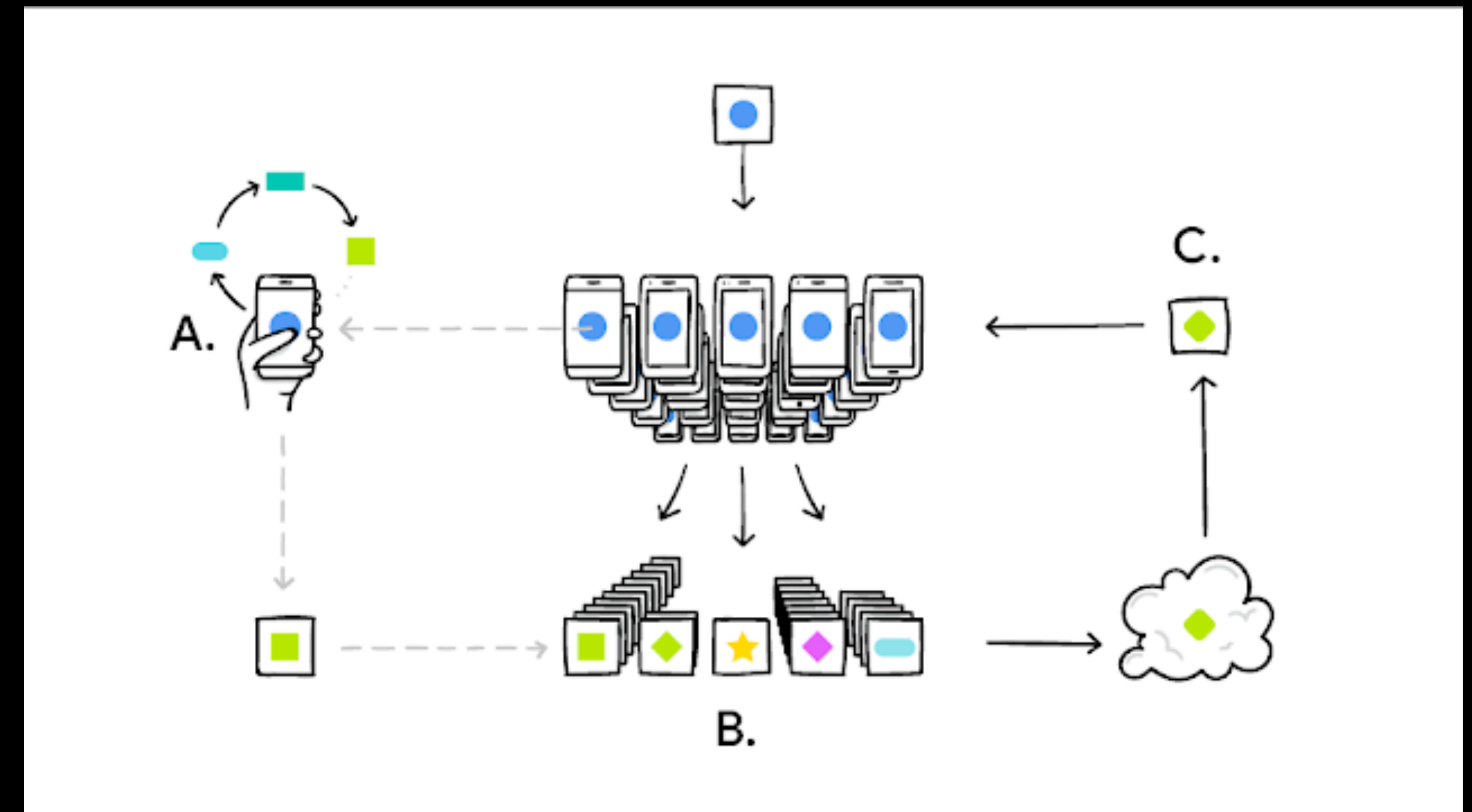
## Why Federated Learning?

- Privacy is an increasingly prominent concern in training machine learning models, and sometimes is legally required (HIPPA)
- Transmission of information over certain channels is expensive
- Utilize local computational resources

# Federated Learning

## Overview

- Similar to a Distributed system, but more general
- Common assumptions in ML do not necessarily hold:
  - Data may not be IID
  - Compute nodes may be highly heterogeneous

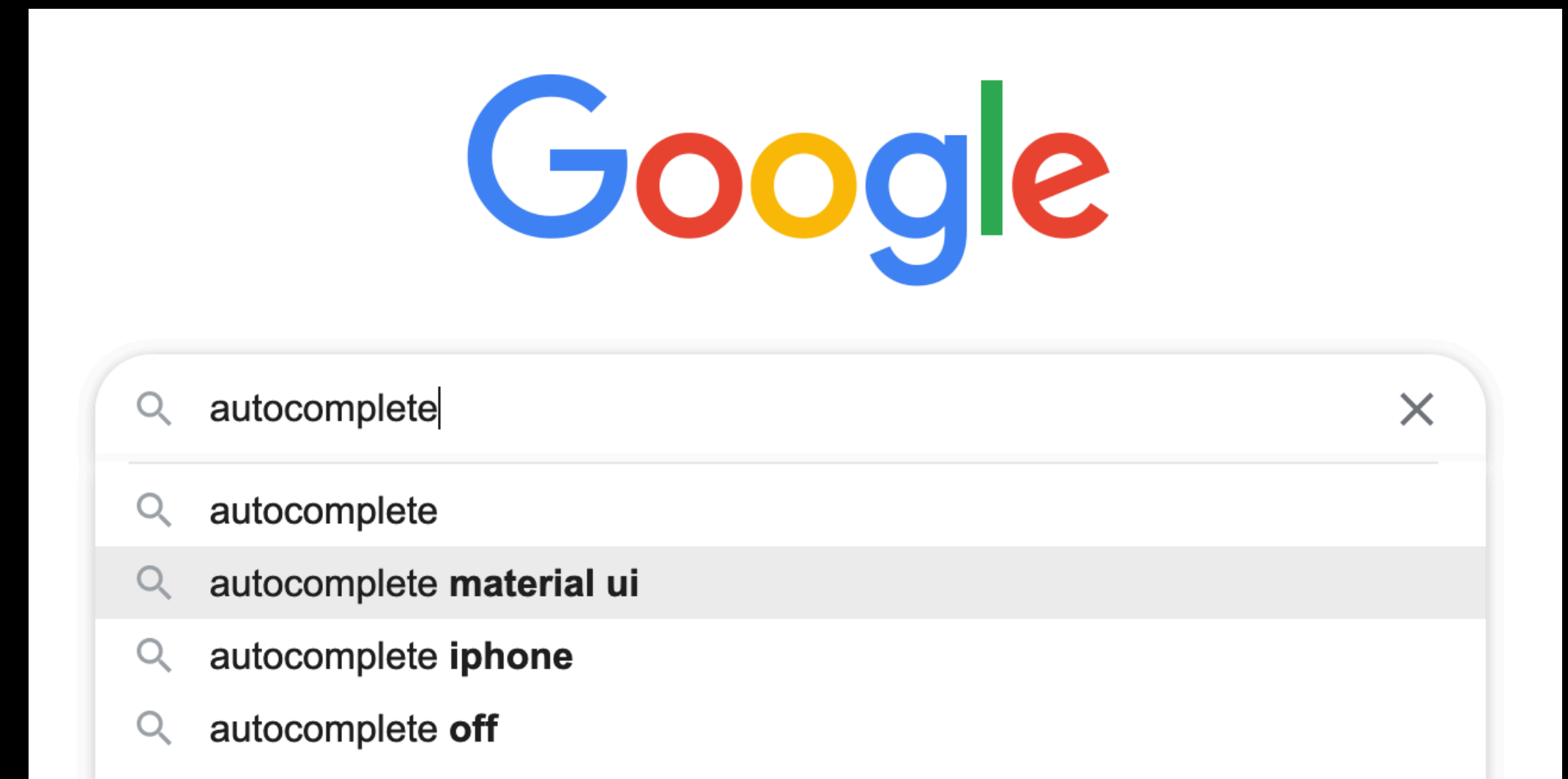


<https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>

# Federated Learning

## An example: autocomplete

- Autocomplete on mobile phones
  - Text collected in private messages, internet browsing, etc.
  - Information collected across the US
- Clearly, the standard statistical and computational assumptions do not hold.
- Many more use cases exist, Ex. Medical data



# Federated Learning

## The prototypical algorithm

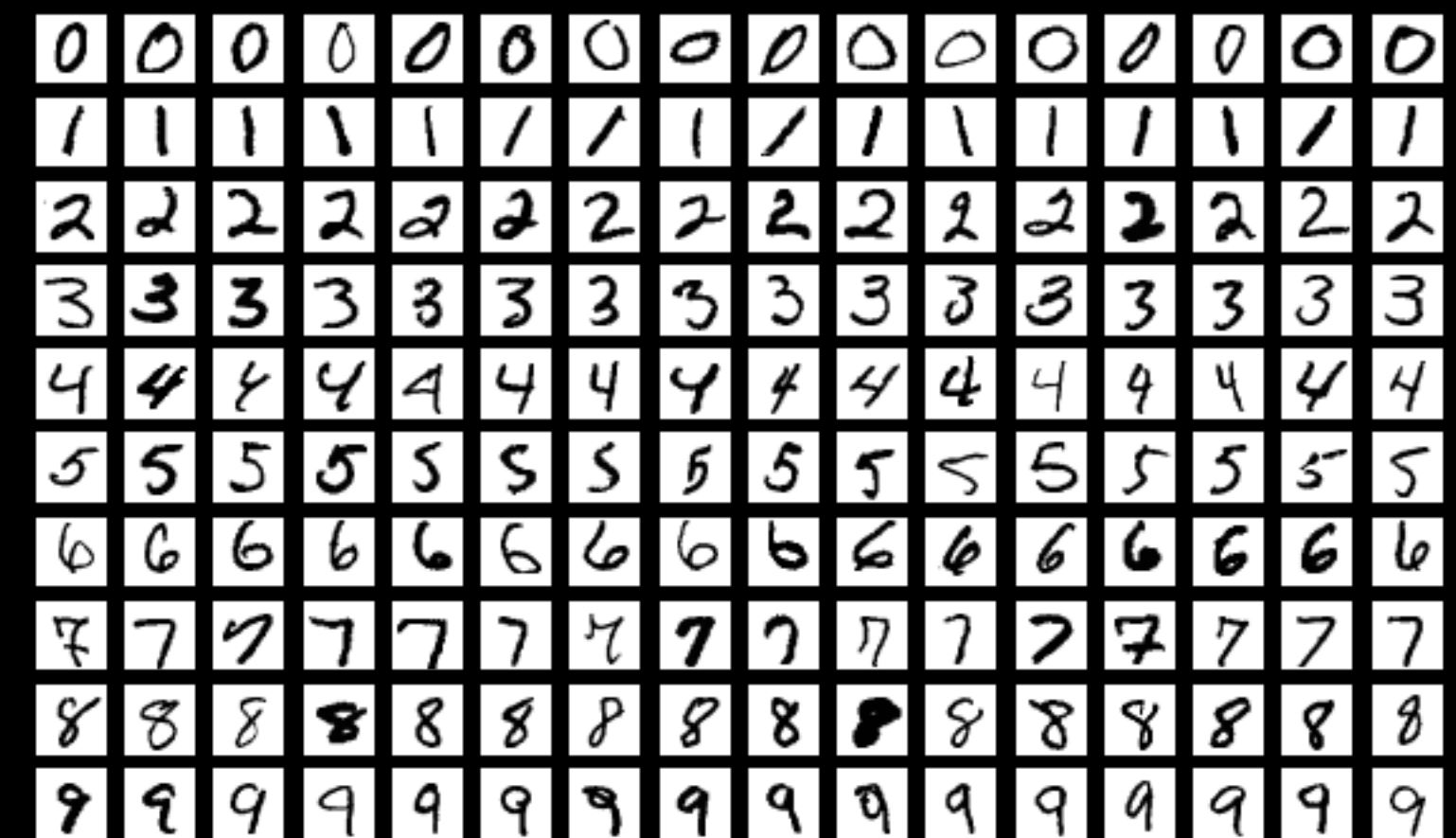
- Dataset located on remote workers, not visible to master
- Worker weights can be used to tune less and more important workers
- The average of each worker's minibatch gradient is used for global update

### Algorithm 1: Federated SGD

```
Initialize weight:  $w = w_0$ ;  
Worker weights:  $p_1, \dots, p_k$ ;  
Worker data:  $\mathcal{D}_1, \dots, \mathcal{D}_k$ ;  
Learning rate  $\alpha$ ;  
while not converged do  
    Master broadcast  $w_t$  to all workers;  
    foreach worker  $i = 1, \dots, k$  do  
         $w_{t+1}^i := w_t - \alpha \nabla \mathcal{L}(w_t; \mathcal{D}_i)$ ;  
        Worker broadcast  $w_t^i$  master;  
    end  
    Master calculate new weight  $w_{k+1} := \sum_{i=1}^k p_i w_{t+1}^i$ ;  
end
```

# A simple problem

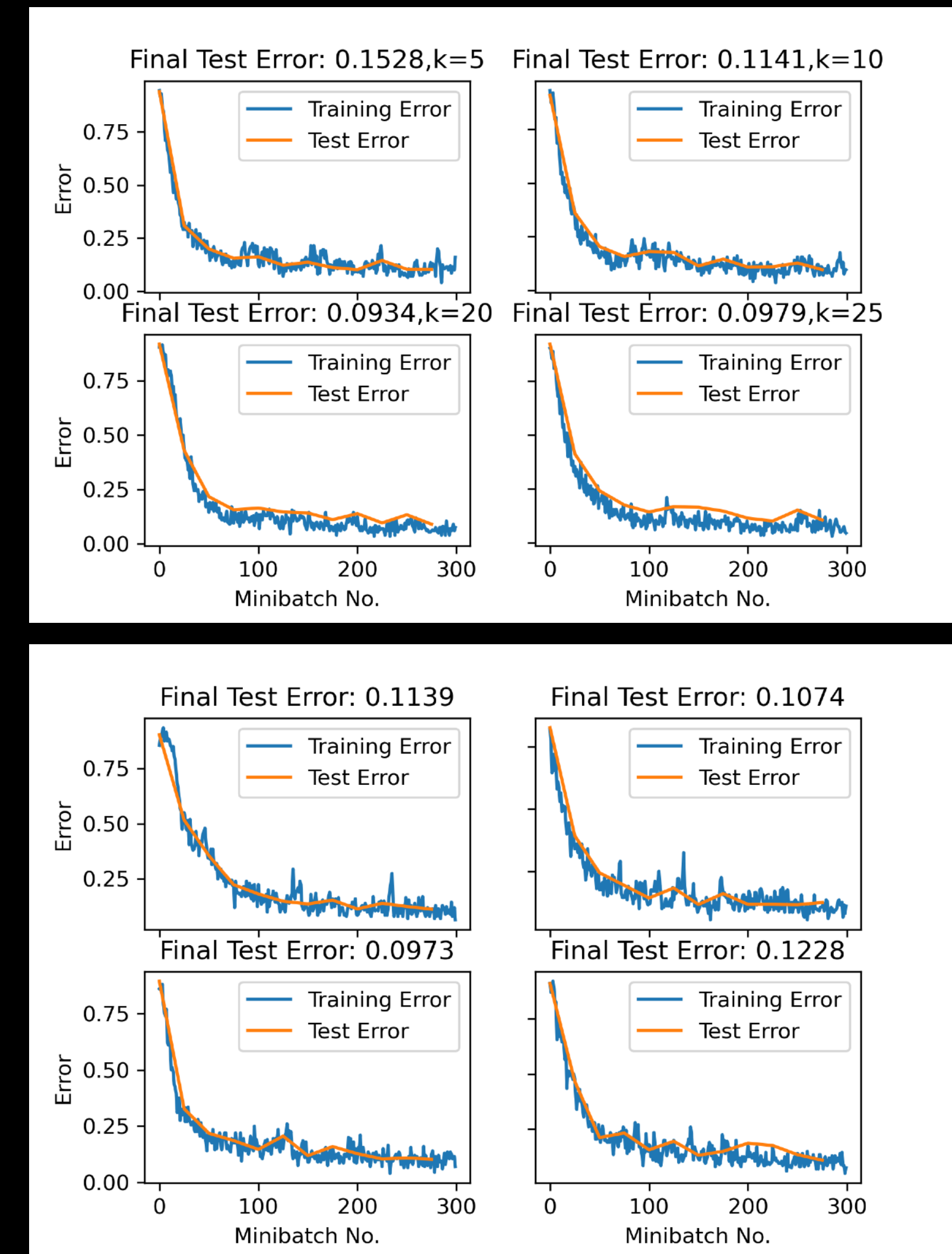
- Train a simple Convolutional Neural Network on the MNIST dataset
- The simple setup makes running multiple experiments much more feasible



# Comparing Centralized to Federated SGD

## Federated at top, Centralized at bottom, $k = \text{Num Workers}$

- At first, they appear the same
  - However, once we realize the central node and worker nodes each have equal minibatch size, we realize this is not the case
- The Federated setup with  $k$  workers requires  $k$  times as much computation



# Learning with malicious users

## Corrupt Data

- Let's imagine a scenario where some workers nodes send malicious or otherwise bad updates to the master
- Examples
  - Text completion (some users often use foul language)
  - Malicious intent



# Learning with malicious users

## Corrupt Data

- Experimental setup:
  - Same learning problem as before (MNIST classification)
  - Now, some percentage of users say  $\alpha$  send corrupted update
    - To mimic this, these users will train on datas with label  $9-y$ , where  $y$  is the true label (an integer on  $[0, \dots, 9]$ )
  - To mitigate this, we will try to take taking the coordinate wise median, as well as the trimmed mean and see if performance improves.

# Learning with malicious users

## Corrupt Data

- Sometimes the raw uncorrected model achieves the best error rate, though this is by an extremely thin margin when it does happen
- Overall the median is the best performer
- While the improvements are modest, using median is simple and sensible

<http://proceedings.mlr.press/v80/yin18a/yin18a.pdf>

Num Workers:	10	20	25
$\alpha$			
.05	X	.1272	<b>.1081</b>
.1	.1327	.1467	.1212
.2	.1845	<b>.1584</b>	.1906

Table 1: Test Error results with malicious users, using the regular mean update rule

Num Workers:	10	20	25
$\alpha$			
.05	X	.1447	.1486
.1	<b>.0996</b>	<b>.1162</b>	.1754
.2	.1563	.1609	.1353

Table 2: Test Error results with malicious users, using coordinate-wise trimmed mean,  $\alpha = \beta$

Num Workers:	10	20	25
$\alpha$			
.05	X	<b>.0957</b>	.1091
.1	.1215	.2027	<b>.0998</b>
.2	<b>.1344</b>	.17	<b>.1209</b>

Table 3: Test Error results with malicious users, using coordinate-wise median

# Learning with malicious users

## Snooping

- Actual **training data** can be inferred from model parameter updates
- This can be achieved with knowledge only of the model state, and gradient update at a single step
- Only caveat is that batch size needs to be 1

---

**Algorithm 1** Deep Leakage from Gradients.

---

**Input:**  $F(\mathbf{x}; W)$ : Differentiable machine learning model;  $W$ : parameter weights;  $\nabla W$ : gradients calculated by training data

**Output:** private training data  $\mathbf{x}, \mathbf{y}$

```
1: procedure DLG( $F, W, \nabla W$ )  
2:    $\mathbf{x}'_1 \leftarrow \mathcal{N}(0, 1), \mathbf{y}'_1 \leftarrow \mathcal{N}(0, 1)$  ▷ Initialize dummy inputs and labels.  
3:   for  $i \leftarrow 1$  to  $n$  do  
4:      $\nabla W'_i \leftarrow \partial \ell(F(\mathbf{x}'_i, W_t), \mathbf{y}'_i) / \partial W_t$  ▷ Compute dummy gradients.  
5:      $\mathbb{D}_i \leftarrow \|\nabla W'_i - \nabla W\|^2$   
6:      $\mathbf{x}'_{i+1} \leftarrow \mathbf{x}'_i - \eta \nabla_{\mathbf{x}'_i} \mathbb{D}_i, \mathbf{y}'_{i+1} \leftarrow \mathbf{y}'_i - \eta \nabla_{\mathbf{y}'_i} \mathbb{D}_i$  ▷ Update data to match gradients.  
7:   end for  
8:   return  $\mathbf{x}'_{n+1}, \mathbf{y}'_{n+1}$   
9: end procedure
```

---

- <https://arxiv.org/pdf/1906.08935.pdf>

# Learning with malicious users

## Snooping

- Able to reconstruct a training sample about 50% of the time
- It appears acquisition is more likely earlier on in the training process

