

# Improved Generalization Bounds for Communication Efficient Federated Learning Federated Learning with Adaptive Local Steps (FedALS)

Peyman Gholami, Hulya Seferoglu

University of Illinois Chicago

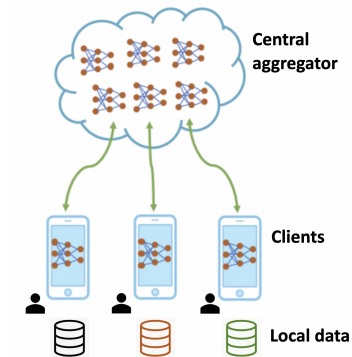
July 2024

# Table of Contents

- 1 Introduction
  - Communication overhead
  - Representation learning
  - Motivation
- 2 Problem Statement
  - Notation
- 3 Improved Generalization Bounds
  - One-Round FL
  - $R$ -Round FL
- 4 FedALS
- 5 Experiments
  - Generalization bound
  - FedALS
- 6 Conclusion

# Introduction

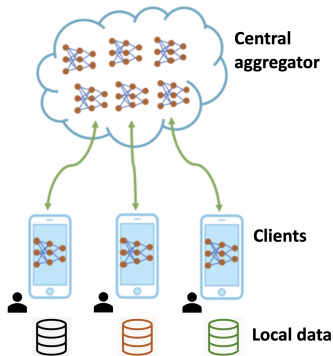
## Federated learning



# Introduction

## Federated learning

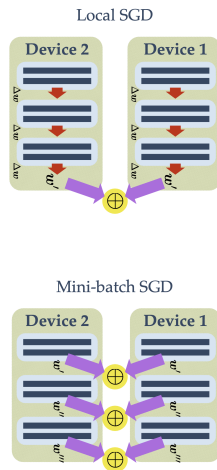
- **Communication cost:**  
Exchanging models is costly, especially for large models in today's machine learning applications like LLMs.



# Introduction

## Federated learning

- Communication cost:  
Exchanging models is costly, especially for large models in today's machine learning applications like LLMs.
- Possible solutions:
  - **Local SGD**
  - Mini-batch SGD



**Figure:** From Lin, Tao, Sebastian U. Stich and Martin Jaggi. "Don't Use Large Mini-Batches, Use Local SGD."

# Federated learning

Purpose of communication:

# Federated learning

Purpose of communication:

- Reducing the consensus distance among clients.
- Consensus distance at  $t$ :  $\frac{1}{K} \sum_{k=1}^K \|\hat{\theta}_t - \theta_{k,t}\|^2$ , where  $\hat{\theta}_t = \frac{1}{K} \sum_{k=1}^K \theta_{k,t}$ .

# Federated learning

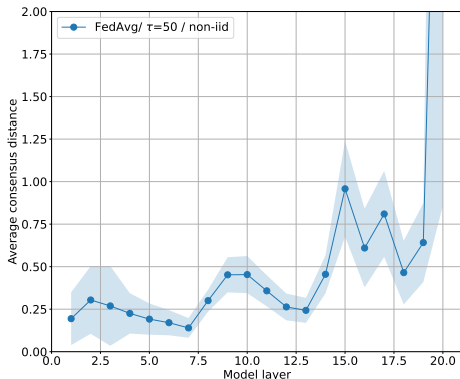
Purpose of communication:

- Reducing the consensus distance among clients.
- Consensus distance at  $t$ :  $\frac{1}{K} \sum_{k=1}^K \|\hat{\theta}_t - \theta_{k,t}\|^2$ , where  $\hat{\theta}_t = \frac{1}{K} \sum_{k=1}^K \theta_{k,t}$ .
- Helps maintain the overall optimization process on a trajectory toward the global optimum.



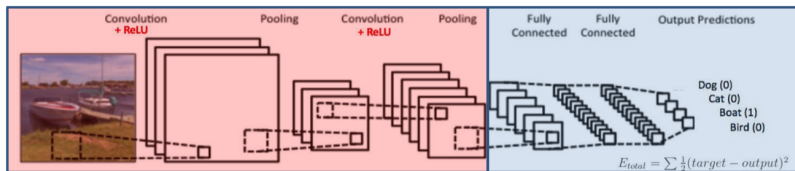
# Consensus distance

Example:



**Figure:** ResNet-20 on CIFAR-10 with 5 clients with non-iid data distribution over clients (2 classes per client). The early layers responsible for extracting representations exhibit lower levels of consensus distance.

# Representation learning



## Learning the representation

**Feature extraction**

**“Backbone” (transferable)**

**Higer Complexity**

## Learning classification

**Fully connected**

**Head**

**Lower Complexity**

# Motivation

- The above example indicates that initial layers show higher similarity, so they can be aggregated less frequently.

---

<sup>1</sup>Sashank J. Reddi et al. “Adaptive Federated Optimization”. In: *International Conference on Learning Representations*. 2021. URL: <https://openreview.net/forum?id=LkFG3lB13U5>; Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. “Salvaging Federated Learning by Local Adaptation”. In: *ArXiv abs/2002.04758* (2020). URL: <https://api.semanticscholar.org/CorpusID:211082601>.

# Motivation

- The above example indicates that initial layers show higher similarity, so they can be aggregated less frequently.
- several empirical studies<sup>1</sup> show that federated learning with multiple local updates per round learns a generalizable representation and is unexpectedly successful in non-iid federated learning.

---

<sup>1</sup>Sashank J. Reddi et al. “Adaptive Federated Optimization”. In: *International Conference on Learning Representations*. 2021. URL: <https://openreview.net/forum?id=LkFG31B13U5>; Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. “Salvaging Federated Learning by Local Adaptation”. In: *ArXiv abs/2002.04758* (2020). URL: <https://api.semanticscholar.org/CorpusID:211082601>.

# Motivation

- The above example indicates that initial layers show higher similarity, so they can be aggregated less frequently.
- several empirical studies<sup>1</sup> show that federated learning with multiple local updates per round learns a generalizable representation and is unexpectedly successful in non-iid federated learning.

**Motivate us to investigate how local updates affect the generalization of the model.**

---

<sup>1</sup>Sashank J. Reddi et al. “Adaptive Federated Optimization”. In: *International Conference on Learning Representations*. 2021. URL: <https://openreview.net/forum?id=LkFG3lB13U5>; Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. “Salvaging Federated Learning by Local Adaptation”. In: *ArXiv abs/2002.04758* (2020). URL: <https://api.semanticscholar.org/CorpusID:211082601>.

# Problem Statement

- 1  $K$  clients/nodes.

# Problem Statement

- 1  $K$  clients/nodes.
- 2 Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .

# Problem Statement

- 1  $K$  clients/nodes.
- 2 Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .
- 3 Dataset across all nodes is defined as  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_K\}$ .



# Problem Statement

- ①  $K$  clients/nodes.
- ② Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .
- ③ Dataset across all nodes is defined as  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_K\}$ .
- ④  $M_\theta = \mathcal{A}(\mathbf{S})$  is the output of a possibly stochastic function denoted as  $\mathcal{A}(\mathbf{S})$ , where  $M_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is the learned model parameterized by  $\theta$ .

# Problem Statement

- 1  $K$  clients/nodes.
- 2 Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .
- 3 Dataset across all nodes is defined as  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_K\}$ .
- 4  $M_\theta = \mathcal{A}(\mathbf{S})$  is the output of a possibly stochastic function denoted as  $\mathcal{A}(\mathbf{S})$ , where  $M_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is the learned model parameterized by  $\theta$ .
- 5 Empirical risk on dataset  $\mathbf{S}$ :  
$$R_{\mathbf{S}}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_{\mathbf{S}_k}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} \frac{1}{n_k} \sum_{i=1}^{n_k} l(M_\theta, \mathbf{z}_{k,i})$$

# Problem Statement

- ①  $K$  clients/nodes.
- ② Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .
- ③ Dataset across all nodes is defined as  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_K\}$ .
- ④  $M_\theta = \mathcal{A}(\mathbf{S})$  is the output of a possibly stochastic function denoted as  $\mathcal{A}(\mathbf{S})$ , where  $M_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is the learned model parameterized by  $\theta$ .
- ⑤ Empirical risk on dataset  $\mathbf{S}$ :  

$$R_{\mathbf{S}}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_{\mathbf{S}_k}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} \frac{1}{n_k} \sum_{i=1}^{n_k} l(M_\theta, \mathbf{z}_{k,i})$$
- ⑥ Population risk:  $R(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_k(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}, \mathbf{z} \sim \mathcal{D}_k} l(M_\theta, \mathbf{z})$

# Problem Statement

- ①  $K$  clients/nodes.
- ② Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .
- ③ Dataset across all nodes is defined as  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_K\}$ .
- ④  $M_\theta = \mathcal{A}(\mathbf{S})$  is the output of a possibly stochastic function denoted as  $\mathcal{A}(\mathbf{S})$ , where  $M_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is the learned model parameterized by  $\theta$ .
- ⑤ Empirical risk on dataset  $\mathbf{S}$ :  

$$R_{\mathbf{S}}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_{\mathbf{S}_k}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} \frac{1}{n_k} \sum_{i=1}^{n_k} l(M_\theta, \mathbf{z}_{k,i})$$
- ⑥ Population risk:  $R(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_k(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}, \mathbf{z} \sim \mathcal{D}_k} l(M_\theta, \mathbf{z})$
- ⑦ Generalization error for dataset  $\mathbf{S}$  and function  $\mathcal{A}(\mathbf{S})$ :  

$$\Delta_{\mathcal{A}}(\mathbf{S}) = R(\mathcal{A}(\mathbf{S})) - R_{\mathbf{S}}(\mathcal{A}(\mathbf{S}))$$

# Problem Statement

- 1  $K$  clients/nodes.
- 2 Client  $k$  has a local dataset  $\mathbf{S}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,n_k}\}$ , where  $\mathbf{z}_{k,i} = (\mathbf{x}_{k,i}, \mathbf{y}_{k,i})$  is drawn from a distribution  $\mathcal{D}_k$  over  $\mathcal{X} \times \mathcal{Y}$ .
- 3 Dataset across all nodes is defined as  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_K\}$ .
- 4  $M_\theta = \mathcal{A}(\mathbf{S})$  is the output of a possibly stochastic function denoted as  $\mathcal{A}(\mathbf{S})$ , where  $M_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is the learned model parameterized by  $\theta$ .
- 5 Empirical risk on dataset  $\mathbf{S}$ :  

$$R_S(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_{\mathbf{S}_k}(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} \frac{1}{n_k} \sum_{i=1}^{n_k} l(M_\theta, \mathbf{z}_{k,i})$$
- 6 Population risk:  $R(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}} R_k(M_\theta) = \mathbb{E}_{k \sim \mathcal{K}, \mathbf{z} \sim \mathcal{D}_k} l(M_\theta, \mathbf{z})$
- 7 Generalization error for dataset  $\mathbf{S}$  and function  $\mathcal{A}(\mathbf{S})$ :  

$$\Delta_{\mathcal{A}}(\mathbf{S}) = R(\mathcal{A}(\mathbf{S})) - R_S(\mathcal{A}(\mathbf{S}))$$
- 8 Expected generalization:  $\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) = \mathbb{E}_{\{\mathbf{S}_k \sim \mathcal{D}_k^{n_k}\}_{k=1}^K} \Delta_{\mathcal{A}}(\mathbf{S})$

# One-Round Generalization Bound

## Theorem

Let  $l(M_\theta, \mathbf{z})$  be  $\mu$ -strongly convex and  $L$ -smooth in  $M_\theta$ ,  $M_{\theta_k} = \mathcal{A}_k(\mathbf{S}_k)$  represents the model obtained from Empirical Risk Minimization (ERM) algorithm on local dataset  $\mathbf{S}_k$ , i.e.,  $M_{\theta_k} = \arg \min_M \sum_{i=1}^{n_k} l(M, \mathbf{z}_{k,i})$ , and  $M_{\hat{\theta}} = \mathcal{A}(\mathbf{S})$  is the model after one round of FedAvg ( $\hat{\theta} = \mathbb{E}_{k \sim \mathcal{K}} \theta_k$ ). Then, the expected generalization error is

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \underbrace{\frac{LK(k)^2}{\mu} \mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right. \\ \left. + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right], \quad (1)$$

where  $\delta_{k,\mathcal{A}}(\mathbf{S}) = R_{\mathbf{S}_k}(\mathcal{A}(\mathbf{S})) - R_{\mathbf{S}_k}(\mathcal{A}_k(\mathbf{S}_k))$  indicates the level of non-iidness at client  $k$  for function  $\mathcal{A}$  on dataset  $\mathbf{S}$ .

# One-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{L\mathcal{K}(k)^2}{\mu} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right],$$

- iid:  $\mathcal{K}(k)^2$  (enhancement compared to the state of the art<sup>2</sup>)

<sup>2</sup>Leighton Pate Barnes, Alex Dytso, and H. Vincent Poor. “Improved Information-Theoretic Generalization Bounds for Distributed, Federated, and Iterative Learning”. In: *Entropy* 24 (2022). URL:

<https://api.semanticscholar.org/CorpusID:246634528>.

# One-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{L\mathcal{K}(k)^2}{\mu} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right],$$

- iid:  $\mathcal{K}(k)^2$  (enhancement compared to the state of the art<sup>2</sup>)
- noniid: the expected generalization error bound does not necessarily decrease with averaging.

<sup>2</sup>Leighton Pate Barnes, Alex Dytso, and H. Vincent Poor. “Improved Information-Theoretic Generalization Bounds for Distributed, Federated, and Iterative Learning”. In: *Entropy* 24 (2022). URL:

<https://api.semanticscholar.org/CorpusID:246634528>. 



# One-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{L\mathcal{K}(k)^2}{\mu} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right],$$

- iid:  $\mathcal{K}(k)^2$  (enhancement compared to the state of the art<sup>2</sup>)
- noniid: the expected generalization error bound does not necessarily decrease with averaging.

**FedAvg works well in iid setup.**

<sup>2</sup>Leighton Pate Barnes, Alex Dytso, and H. Vincent Poor. “Improved Information-Theoretic Generalization Bounds for Distributed, Federated, and Iterative Learning”. In: *Entropy* 24 (2022). URL:

<https://api.semanticscholar.org/CorpusID:246634528>.

# One-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{L\mathcal{K}(k)^2}{\mu} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right],$$

**Partial client participation:**

# One-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{L\mathcal{K}(k)^2}{\mu} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right],$$

## Partial client participation:

Case 1: Sampling  $\hat{K}$  clients with replacement based on distribution  $\mathcal{K}$ , followed by averaging the local models with equal weights. ( $\mathcal{K}(k) \rightarrow \frac{1}{\hat{K}}$ )

# One-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\mathcal{A}}(\mathbf{S}) \leq \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{L\mathcal{K}(k)^2}{\mu} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} + 2\sqrt{\frac{L}{\mu}} \mathcal{K}(k) \left( \underbrace{\mathbb{E}_{\mathbf{S}} \delta_{k,\mathcal{A}}(\mathbf{S})}_{\text{Expected non-iidness}} \underbrace{\mathbb{E}_{\mathbf{S}_k} \Delta_{\mathcal{A}_k}(\mathbf{S}_k)}_{\text{Expected local generalization}} \right)^{\frac{1}{2}} \right],$$

## Partial client participation:

*Case I:* Sampling  $\hat{K}$  clients with replacement based on distribution  $\mathcal{K}$ , followed by averaging the local models with equal weights. ( $\mathcal{K}(k) \rightarrow \frac{1}{\hat{K}}$ )

*Case II:* Sampling  $\hat{K}$  clients without replacement uniformly at random, then performing weighted averaging of local models. Here, the weight of client  $k$  is rescaled to  $\frac{\mathcal{K}(k)K}{\hat{K}}$ . ( $\mathcal{K}(k) \rightarrow \frac{\mathcal{K}(k)K}{\hat{K}}$ )

# $R$ –Round Generalization Bound

- $Z_{k,r} = \bigcup \{\mathcal{B}_{k,r,t}\}_{t=0}^{\tau-1}$ , where  $\mathcal{B}_{k,r,t}$  is the batch of samples used in local step  $t$  of round  $r$  in node  $k$ .  $\tau$  is the duration of one round.

# $R$ –Round Generalization Bound

- $Z_{k,r} = \bigcup \{\mathcal{B}_{k,r,t}\}_{t=0}^{\tau-1}$ , where  $\mathcal{B}_{k,r,t}$  is the batch of samples used in local step  $t$  of round  $r$  in node  $k$ .  $\tau$  is the duration of one round.
- Empirical risk:  $\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} l(M_{\hat{\theta}_r}, \mathbf{z}_{k,i}) \right]$

# $R$ –Round Generalization Bound

- $Z_{k,r} = \bigcup \{\mathcal{B}_{k,r,t}\}_{t=0}^{\tau-1}$ , where  $\mathcal{B}_{k,r,t}$  is the batch of samples used in local step  $t$  of round  $r$  in node  $k$ .  $\tau$  is the duration of one round.
- Empirical risk:  $\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} l(M_{\hat{\theta}_r}, \mathbf{z}_{k,i}) \right]$
- Generalization error:  

$$\Delta_{\text{FedAvg}}(\mathbf{S}) = \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \mathbb{E}_{\mathbf{z} \sim \mathcal{D}_k} l(M_{\hat{\theta}_r}, \mathbf{z}) - \frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} l(M_{\hat{\theta}_r}, \mathbf{z}_{k,i}) \right]$$

# R-Round Generalization Bound

- $Z_{k,r} = \bigcup \{\mathcal{B}_{k,r,t}\}_{t=0}^{\tau-1}$ , where  $\mathcal{B}_{k,r,t}$  is the batch of samples used in local step  $t$  of round  $r$  in node  $k$ .  $\tau$  is the duration of one round.
- Empirical risk:  $\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} l(M_{\hat{\theta}_r}, \mathbf{z}_{k,i}) \right]$
- Generalization error:  

$$\Delta_{\text{FedAvg}}(\mathbf{S}) = \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \mathbb{E}_{\mathbf{z} \sim \mathcal{D}_k} l(M_{\hat{\theta}_r}, \mathbf{z}) - \frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} l(M_{\hat{\theta}_r}, \mathbf{z}_{k,i}) \right]$$
- Bounded gradient variance:  

$$\frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} \left\| \nabla l(M, \mathbf{z}_{k,i}) - \frac{1}{|Z_{k,r}|} \sum_{i \in Z_{k,r}} \nabla l(M, \mathbf{z}_{k,i}) \right\|^2 \leq \sigma^2.$$



# R-Round Generalization Bound

## Theorem

Let  $l(M_\theta, \mathbf{z})$  be  $\mu$ -strongly convex and  $L$ -smooth in  $M_\theta$ . Local models at round  $r$  are calculated by doing  $\tau$  local steps and the gradient variance is bounded by  $\sigma^2$ . The aggregated model at round  $r$  is  $M_{\hat{\theta}_r}$  is obtained by performing FedAvg and where the data points used in round  $r$  (i.e.,  $Z_{k,r}$ ) are sampled without replacement. Then the average generalization error bound is

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{2LK(k)^2}{\mu} A + \sqrt{\frac{8L}{\mu}} \mathcal{K}(k) (AB)^{\frac{1}{2}} \right], \quad (2)$$

where  $A = \tilde{O} \left( \sqrt{\frac{C(M_\theta)}{|Z_{k,r}|}} + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right),$

$B = \tilde{O} \left( \mathbb{E}_{\{Z_{k,r}\}_{k=1}^K} \delta_{k,A}(\{Z_{k,r}\}_{k=1}^K) + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right),$  and  $C(M_\theta)$  shows the complexity of the model class of  $M_\theta$ .

# R-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\text{FedAvg}}(\mathbf{S}) \leq \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{2L\mathcal{K}(k)^2}{\mu} A + \sqrt{\frac{8L}{\mu}} \mathcal{K}(k) (AB)^{\frac{1}{2}} \right]$$

$$A = \tilde{O} \left( \sqrt{\frac{\mathcal{C}(M_{\theta})}{|Z_{k,r}|}} + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right), B = \tilde{O} \left( \mathbb{E} \delta_{k,\mathcal{A}}(\{Z_{k,r}\}_{k=1}^K) + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right)$$

**Representation learning Interpretation:**

# R-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\text{FedAvg}}(\mathbf{S}) \leq \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{2L\mathcal{K}(k)^2}{\mu} A + \sqrt{\frac{8L}{\mu}} \mathcal{K}(k) (AB)^{\frac{1}{2}} \right]$$

$$A = \tilde{O} \left( \sqrt{\frac{\mathcal{C}(M_{\theta})}{|Z_{k,r}|}} + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right), B = \tilde{O} \left( \mathbb{E} \delta_{k,\mathcal{A}}(\{Z_{k,r}\}_{k=1}^K) + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right)$$

**Representation learning Interpretation:**

$$M_{\theta}(\mathbf{x}) = (M_{\phi} \circ M_{\mathbf{h}})(\mathbf{x}) = M_{\mathbf{h}}(M_{\phi}(\mathbf{x})), \mathcal{C}(M_{\mathbf{h}}) \ll \mathcal{C}(M_{\phi})$$

# R-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\text{FedAvg}}(\mathbf{S}) \leq \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{2L\mathcal{K}(k)^2}{\mu} A + \sqrt{\frac{8L}{\mu}} \mathcal{K}(k) (AB)^{\frac{1}{2}} \right]$$

$$A = \tilde{O} \left( \sqrt{\frac{\mathcal{C}(M_{\theta})}{|Z_{k,r}|}} + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right), B = \tilde{O} \left( \mathbb{E} \delta_{k,\mathcal{A}}(\{Z_{k,r}\}_{k=1}^K) + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right)$$

## Representation learning Interpretation:

$$M_{\theta}(\mathbf{x}) = (M_{\phi} \circ M_{\mathbf{h}})(\mathbf{x}) = M_{\mathbf{h}}(M_{\phi}(\mathbf{x})), \mathcal{C}(M_{\mathbf{h}}) \ll \mathcal{C}(M_{\phi})$$

Our key intuition in this paper is that we can *reduce the aggregation frequency* of  $M_{\phi}$ , which leads to a larger  $\tau$  and  $|Z_{k,r}|$ , hence *smaller generalization error bound*.

# R-Round Generalization Bound

$$\mathbb{E}_{\mathbf{S}} \Delta_{\text{FedAvg}}(\mathbf{S}) \leq \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{k \sim \mathcal{K}} \left[ \frac{2L\mathcal{K}(k)^2}{\mu} A + \sqrt{\frac{8L}{\mu}} \mathcal{K}(k) (AB)^{\frac{1}{2}} \right]$$

$$A = \tilde{O} \left( \sqrt{\frac{\mathcal{C}(M_{\theta})}{|Z_{k,r}|}} + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right), B = \tilde{O} \left( \mathbb{E} \delta_{k,A}(\{Z_{k,r}\}_{k=1}^K) + \frac{\sigma^2}{\mu\tau} + \frac{L}{\mu} \right)$$

## Representation learning Interpretation:

$$M_{\theta}(\mathbf{x}) = (M_{\phi} \circ M_{\mathbf{h}})(\mathbf{x}) = M_{\mathbf{h}}(M_{\phi}(\mathbf{x})), \mathcal{C}(M_{\mathbf{h}}) \ll \mathcal{C}(M_{\phi})$$

Our key intuition in this paper is that we can *reduce the aggregation frequency of  $M_{\phi}$ , which leads to a larger  $\tau$  and  $|Z_{k,r}|$ , hence smaller generalization error bound.*

Aggregation frequency of  $M_{\phi}$  cannot be reduced arbitrarily, as it would increase the empirical risk. (convergence rate:  $O\left(\frac{\tau}{T} + \left(\frac{\tau}{T}\right)^{\frac{2}{3}} + \frac{1}{\sqrt{T}}\right)$ ,  $T = \tau R$ )

# FedALS: Federated Learning with Adaptive Local Steps

- Main idea: to maintain a uniform generalization error across both components ( $M_\phi$  and  $M_h$ ).

---

**Algorithm 1** FedALS
 

---

**Input:** Initial model  $\{\theta_{k,1,0} = [\phi_{k,1,0}, h_{k,1,0}]\}_{k=1}^K$ , Learning rate  $\eta$ , number of local steps for the head model  $\tau$ , adaptation coefficient  $\alpha$ .

```

1: for Round  $r$  in  $1, \dots, R$  do
2:   for Node  $k$  in  $1, \dots, K$  in parallel do
3:     for Local step  $t$  in  $0, \dots, \tau - 1$  do
4:       Sample the batch  $\mathcal{B}_{k,r,t}$  from  $\mathcal{D}_k$ .
5:        $\theta_{k,r,t+1} = \theta_{k,r,t} - \frac{\eta}{|\mathcal{B}_{k,r,t}|} \sum_{i \in \mathcal{B}_{k,r,t}} \nabla l(M_{\theta_{k,r,t}}, z_{k,i})$ 
6:       if  $\text{mod}(r\tau + t, \tau) = 0$  then
7:          $h_{k,r,t} \leftarrow \frac{1}{K} \sum_{k=1}^K h_{k,r,t}$ 
8:       else if  $\text{mod}(r\tau + t, \alpha\tau) = 0$  then
9:          $\phi_{k,r,t} \leftarrow \frac{1}{K} \sum_{k=1}^K \phi_{k,r,t}$ 
10:       $\theta_{k,r+1,0} = \theta_{k,r,\tau}$ 
11: return  $\hat{\theta}_R = \frac{1}{K} \sum_{k=1}^K \theta_{k,R,\tau}$ 
  
```

---

# FedALS: Federated Learning with Adaptive Local Steps

- Main idea: to maintain a uniform generalization error across both components ( $M_\phi$  and  $M_h$ ).
- Introduce parameter  $\alpha = \frac{\tau M_\phi}{\tau M_h}$  as an adaptation coefficient.

---

**Algorithm 1** FedALS
 

---

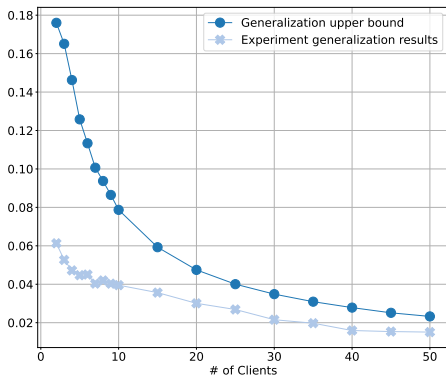
**Input:** Initial model  $\{\theta_{k,1,0} = [\phi_{k,1,0}, h_{k,1,0}]\}_{k=1}^K$ , Learning rate  $\eta$ , number of local steps for the head model  $\tau$ , adaptation coefficient  $\alpha$ .

```

1: for Round  $r$  in  $1, \dots, R$  do
2:   for Node  $k$  in  $1, \dots, K$  in parallel do
3:     for Local step  $t$  in  $0, \dots, \tau - 1$  do
4:       Sample the batch  $\mathcal{B}_{k,r,t}$  from  $\mathcal{D}_k$ .
5:        $\theta_{k,r,t+1} = \theta_{k,r,t} - \frac{\eta}{|\mathcal{B}_{k,r,t}|} \sum_{i \in \mathcal{B}_{k,r,t}} \nabla l(M_{\theta_{k,r,t}}, z_{k,i})$ 
6:       if  $\text{mod}(r\tau + t, \tau) = 0$  then
7:          $h_{k,r,t} \leftarrow \frac{1}{K} \sum_{k=1}^K h_{k,r,t}$ 
8:         else if  $\text{mod}(r\tau + t, \alpha\tau) = 0$  then
9:            $\phi_{k,r,t} \leftarrow \frac{1}{K} \sum_{k=1}^K \phi_{k,r,t}$ 
10:         $\theta_{k,r+1,0} = \theta_{k,r,\tau}$ 
11: return  $\hat{\theta}_R = \frac{1}{K} \sum_{k=1}^K \theta_{k,R,\tau}$ 
  
```

---

# Generalization Upper Bound

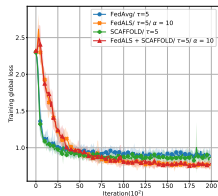


**Figure:** Generalization error and its upper bound derived in this work. The model is a logistic regression on a synthetic dataset generated from a multivariate normal distribution.

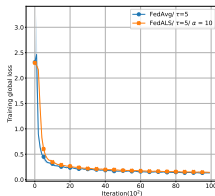
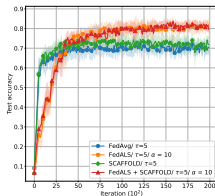


# FedALS Experimental Results

Image classification ( $\phi$  : the convolutional layers of ResNet,  $h$  the final dense layers)



(a) non-iid



(b) iid

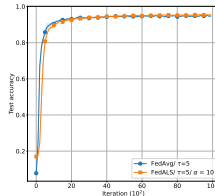
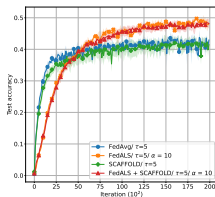
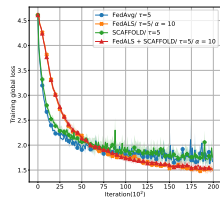


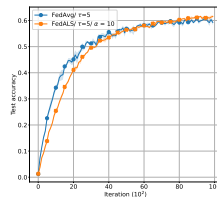
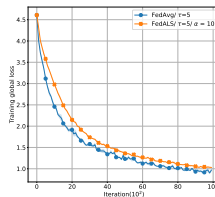
Figure: Training ResNet-20 on SVHN.

# FedALS Experimental Results

Image classification ( $\phi$  : the convolutional layers of ResNet,  $h$  the final dense layers)



(a) non-iid



(b) iid

Figure: Training ResNet-20 on CIFAR-100.

# FedALS Experimental Results

LLM fine-tuning ( $\phi$  : first 10 Transformer layers,  $h$  the final 2 Transformer layers)

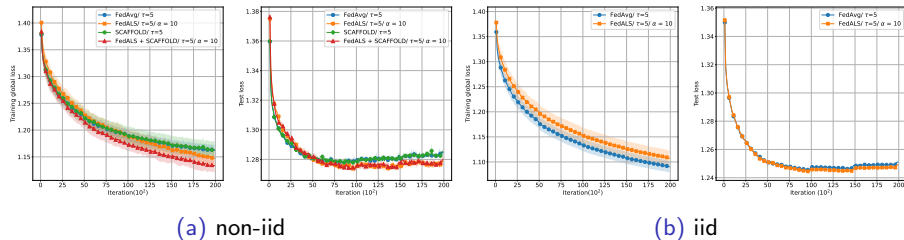


Figure: Fine-tuning OPT-125M on MultiNLI.

# The Role of $\alpha$ and Communication Overhead

$$\alpha = \frac{\tau_{M_\phi}}{\tau_{M_h}}$$

**Table:** The accuracy and communication overhead per client after training ResNet-20 in non-iid setting with  $\tau = 5$  and variable  $\alpha$ .

VALUE OF $\alpha$	DATASET		# OF COMMUNICATED PARAMETERS
	SVHN	CIFAR-10	
1	0.7010 $\pm$ 0.0330	0.4651 $\pm$ 0.0071	2.344 <i>B</i>
5	0.8107 $\pm$ 0.0278	0.5201 $\pm$ 0.0302	0.473 <i>B</i>
10	<b>0.8117 <math>\pm</math> 0.0214</b>	<b>0.5224 <math>\pm</math> 0.0365</b>	0.239 <i>B</i>
25	0.7201 $\pm$ 0.1565	0.3814 $\pm$ 0.0641	0.099 <i>B</i>
50	0.6377 $\pm$ 0.0520	0.2853 $\pm$ 0.0641	0.052 <i>B</i>
100	0.5837 $\pm$ 0.0715	0.2817 $\pm$ 0.032	0.029 <i>B</i>

# Different Combinations of $\phi, h$

$$\theta = [\phi, h] = [\text{first } \mathbf{L} \text{ layers, rest of the layers}]$$

**Table:** Different Combinations of  $\phi, h$  for training ResNet-20 in non-iid setting with  $\tau = 5$ ,  $\alpha = 10$ .

VALUE OF $\mathbf{L}$	DATASET		
	SVHN	CIFAR-10	CIFAR-100
20	0.6991 $\pm$ 0.0160	0.4383 $\pm$ 0.0423	0.4781 $\pm$ 0.0123
16	<b>0.7112 <math>\pm</math> 0.0471</b>	<b>0.4687 <math>\pm</math> 0.0111</b>	<b>0.4782 <math>\pm</math> 0.0087</b>
12	0.6760 $\pm$ 0.0474	0.4125 $\pm$ 0.0283	0.4249 $\pm$ 0.0143
8	0.6381 $\pm$ 0.0428	0.3779 $\pm$ 0.03451	0.4085 $\pm$ 0.0094
4	0.6339 $\pm$ 0.0446	0.3730 $\pm$ 0.0310	0.4183 $\pm$ 0.0108
1	0.6058 $\pm$ 0.0197	0.4013 $\pm$ 0.0308	0.3880 $\pm$ 0.0305

# Conclusion

- Characterized generalization error bound federated learning in terms of local generalization and non-iidness.

# Conclusion

- Characterized generalization error bound federated learning in terms of local generalization and non-iidness.
- Showed that less frequent aggregations, hence more local updates leads to a more generalizable model.

# Conclusion

- Characterized generalization error bound federated learning in terms of local generalization and non-iidness.
- Showed that less frequent aggregations, hence more local updates leads to a more generalizable model.
- This insight led us to develop FedALS algorithm by increasing local steps for the initial layers of a deep learning model while doing more averaging for the final layers.



# Thank You!

# References I



Barnes, Leighton Pate, Alex Dytso, and H. Vincent Poor. “Improved Information-Theoretic Generalization Bounds for Distributed, Federated, and Iterative Learning”. In: *Entropy* 24 (2022). URL: <https://api.semanticscholar.org/CorpusID:246634528>.



Reddi, Sashank J. et al. “Adaptive Federated Optimization”. In: *International Conference on Learning Representations*. 2021. URL: <https://openreview.net/forum?id=LkFG3lB13U5>.



Yu, Tao, Eugene Bagdasaryan, and Vitaly Shmatikov. “Salvaging Federated Learning by Local Adaptation”. In: *ArXiv abs/2002.04758* (2020). URL: <https://api.semanticscholar.org/CorpusID:211082601>.