



# Personalized Federated Learning over Heterogeneous Data

– Jun Luo’s PhD Proposal Defense

November 3<sup>rd</sup> , 2025  
Intelligent Systems Program  
University of Pittsburgh

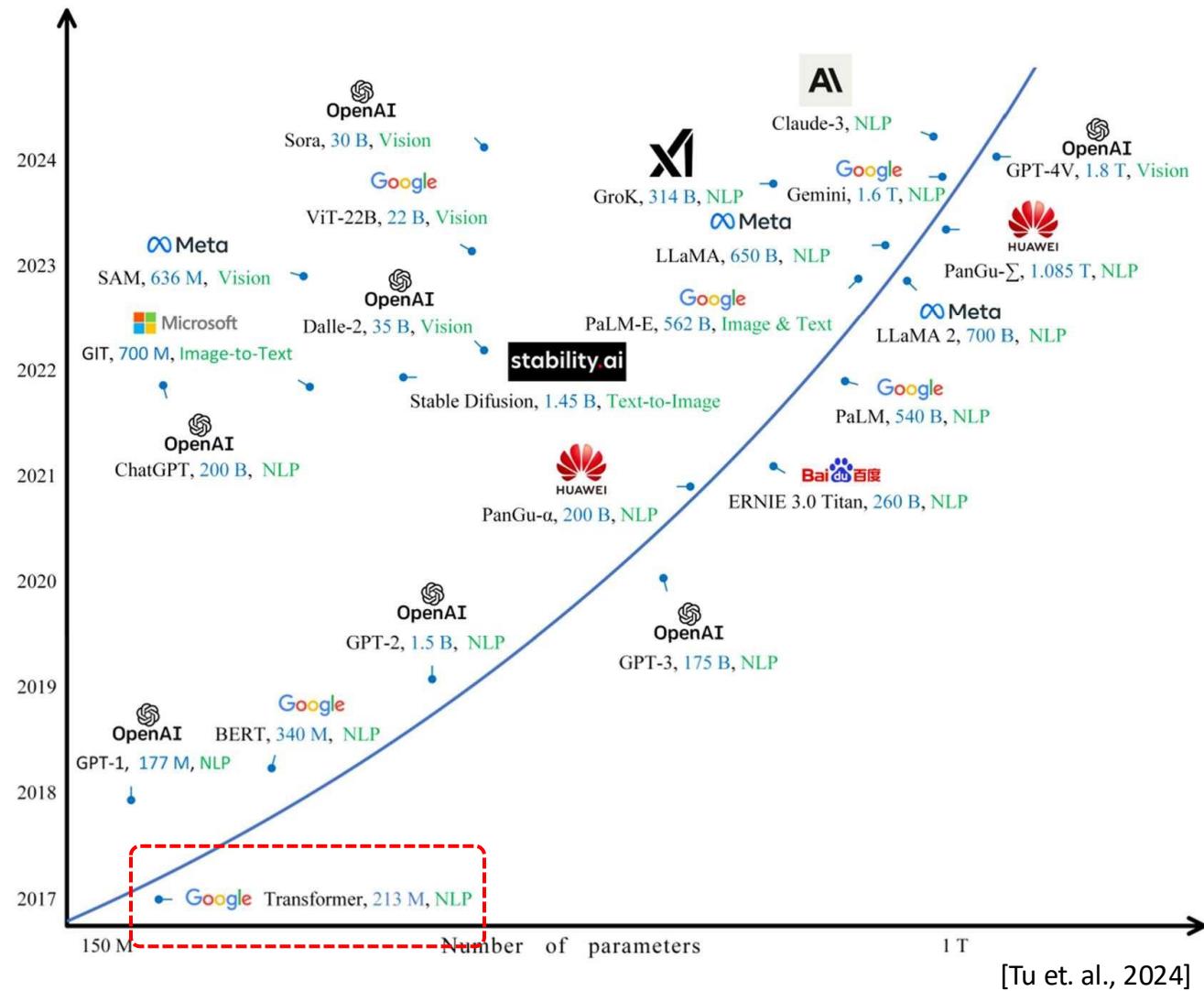
# Committee

- Chair
  - Dr. Shandong Wu, Associate Professor, Intelligent Systems Program
- Members
  - Dr. Leming Zhou, Associate Professor, Intelligent Systems Program
  - Dr. Xiaowei Jia, Assistant Professor, Intelligent Systems Program
  - Dr. Lu Tang, Assistant Professor, Department of Biostatistics

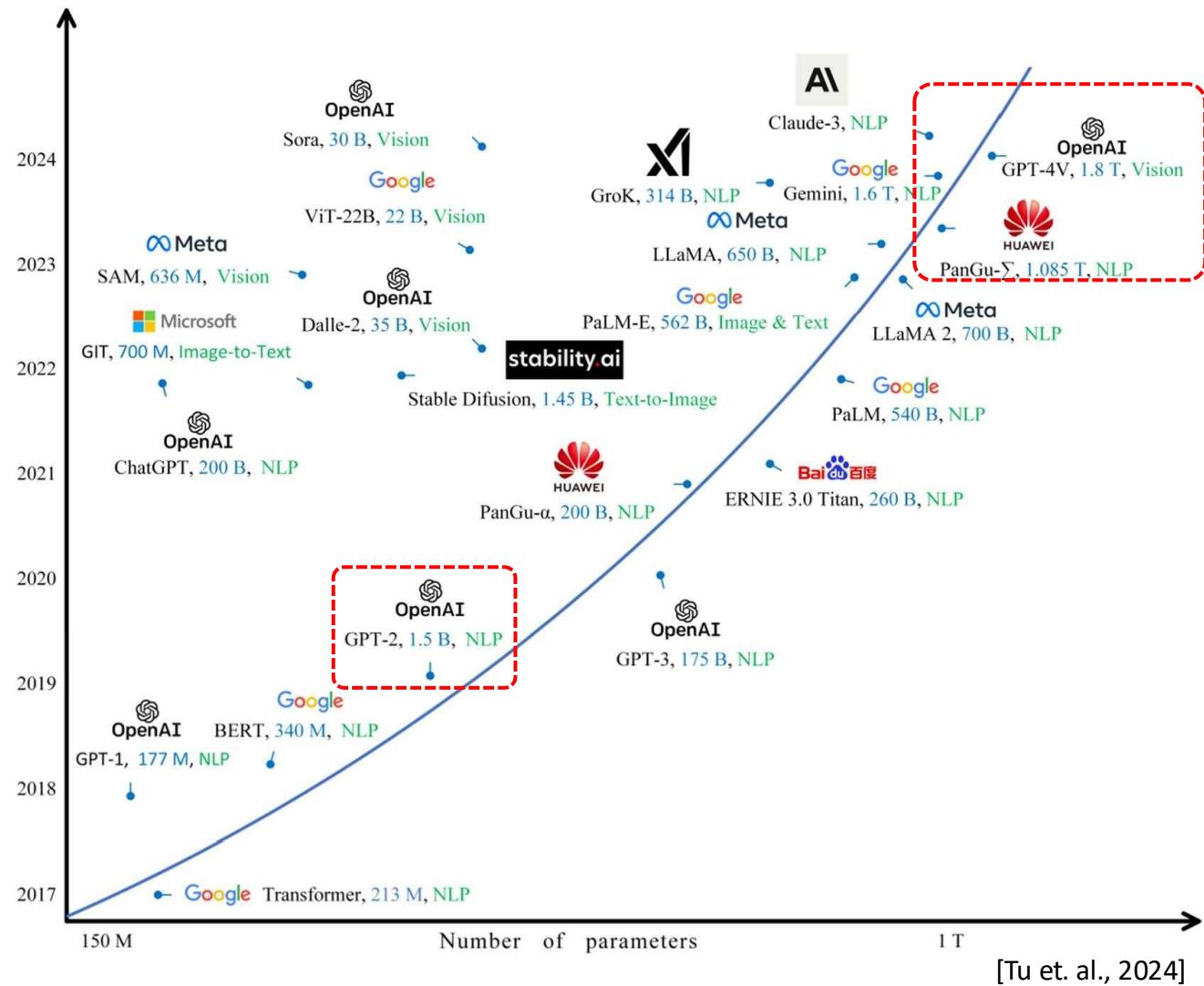
# List of Publications

1. **Jun Luo**, Chen Chen, and Shandong Wu. Mixture of experts made personalized: Federated prompt learning for vision-language models. In *Proceedings of the Thirteenth International Conference on Learning Representation*, 2025
2. **Jun Luo**, Matias Mendieta, Chen Chen, and Shandong Wu. Pgfd: Personalize each client's global objective for federated learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3946–3956, 2023.
3. **Jun Luo** and Shandong Wu. Adapt to adaptation: Learning personalization for cross-silo federated learning. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 2166–2173. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track.
4. **Jun Luo** and Shandong Wu. Fedslid: Federated learning with shared label distribution for medical image classification. In *2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI)*, pages 1–5. IEEE, 2022.
5. **Jun Luo**, Dooman Arefan, Margarita Zuley, Jules Sumkin, and Shandong Wu. Deep curriculum learning in task space for multi-class based mammography diagnosis. In *Medical Imaging 2022: Computer-Aided Diagnosis*, volume 12033, pages 85–90. SPIE, 2022.
6. **Jun Luo**, Gene Kitamura, Dooman Arefan, Emine Doganay, Ashok Panigrahy, and Shandong Wu. Knowledge-guided multiview deep curriculum learning for elbow fracture classification. In *Machine Learning in Medical Imaging: 12th International Workshop, MLMI 2021, Held in Conjunction with MICCAI 2021*, Strasbourg, France, September 27, 2021, Proceedings 12, pages 555–564. Springer, 2021.
7. **Jun Luo**, Gene Kitamura, Emine Doganay, Dooman Arefan, and Shandong Wu. Medical knowledge-guided deep curriculum learning for elbow fracture diagnosis from x-ray images. In *Medical Imaging 2021: Computer-Aided Diagnosis*, volume 11597, pages 247–252. SPIE, 2021.
8. **Jun Luo**, Dooman Arefan, Margarita Zuley, Chen Chen, Jules Sumkin, Shandong Wu. Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection. In *the Radiological Society of North America (RSNA) Annual Meeting (abstract)*, 2025
9. **Jun Luo**, Dooman Arefan, Anil Vasireddi, Shandong Wu, and Nghi Nguyen. Potential use of artificial intelligence in sincalide-stimulated cholescintigraphy: A pilot study. In *Society of Nuclear Medicine and Molecular Imaging (SNMMI) Annual Meeting (abstract, full paper under review by European Journal of Nuclear Medicine and Molecular Imaging)*, 2023.
10. Guangyu Sun, Matias Mendieta, **Jun Luo**, Shandong Wu, and Chen Chen. Fedperfix: Towards partial model personalization of vision transformers in federated learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4988–4998, 2023.
11. Zhengbo Zhou, **Jun Luo**, Dooman Arefan, Gene Kitamura, and Shandong Wu. Human not in the loop: objective sample difficulty measures for curriculum learning. In *2023 IEEE 20<sup>th</sup> International Symposium on Biomedical Imaging (ISBI)*, pages 1–5. IEEE, 2023.
12. Suvendra Vijayan, **Jun Luo**, Shandong Wu, and Anitha Potluri. Image enhancement of ultralow dose cbct images using a deep generative model. In *American Roentgen Ray Society (ARRS) Annual Meeting*, 134(3):e72 (abstract), 2022.
13. Emine Doganay, Gene Kitamura, Luo Yang, **Jun Luo**, and Shandong Wu. Multi-view-enabled deep learning for automated radiographic view classification and fracture detection of the elbow. In *American Roentgen Ray Society (ARRS) Annual Meeting (abstract)*. 2021.

# Background



# Background

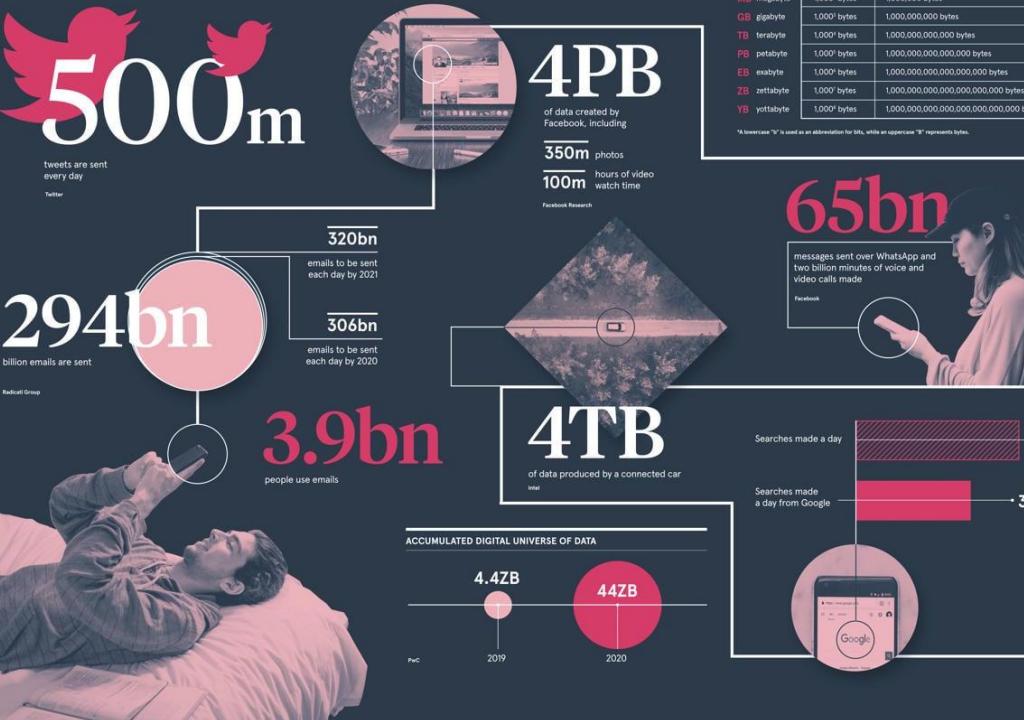


[Tu et al., 2024]

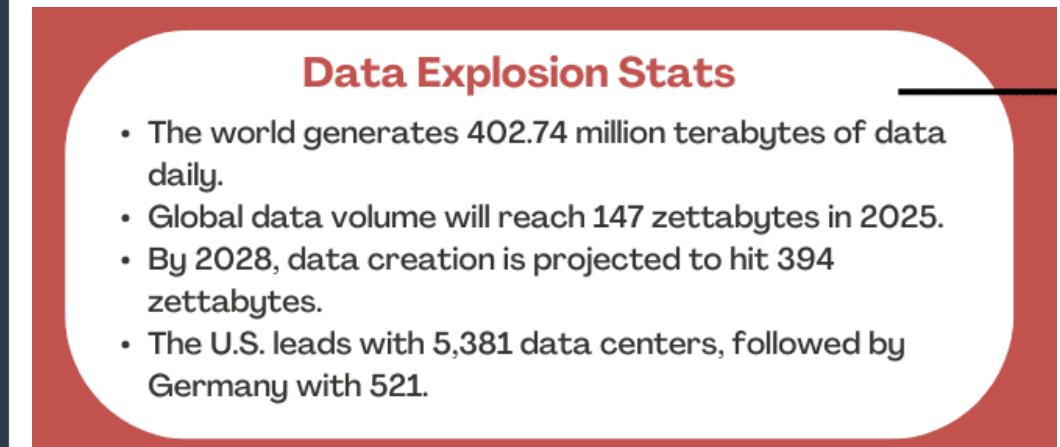
# Background

## A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devices – are hard to comprehend, particularly when looked at in the context of one day



[Jeff Desjardins, World Economic Forum]



[Khyati Hooda, KeyWordsEverywhere.com]

# Background

## Internet of things data



[Risk Management Magazine]

## Health data



[American Hospital Association]

# Background

[Frameworks](#)[Platform](#)[Resources](#)**AMA** [Join](#)[Renew](#) Search[Blog](#) > [Blogs](#) > GDPR Privacy Policy: Ensuring Compliance with EU Data Rules

## GDPR Privacy Policy: Ensuring Compliance with EU Data Rules

**Bhuvesh Lal** • Sep 30, 2024[HIPAA](#)

## HIPAA privacy rule

2 Min Read



Save



Copy



Print

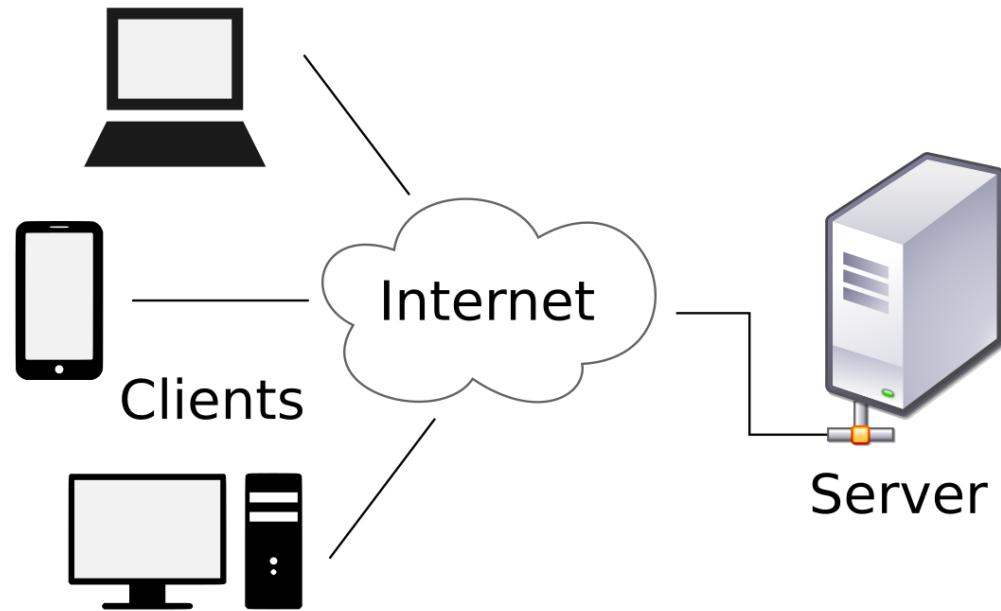


Share

The HIPAA Privacy Rule provides federal standards to safeguard the privacy of personal health information and gives patients an array of rights with respect to that information, including rights to examine and obtain a copy of their health records and to request corrections. The U.S. Department of Health & Human Services' (HHS) Office of Civil Rights (OCR) oversees compliance with HIPAA privacy requirements.

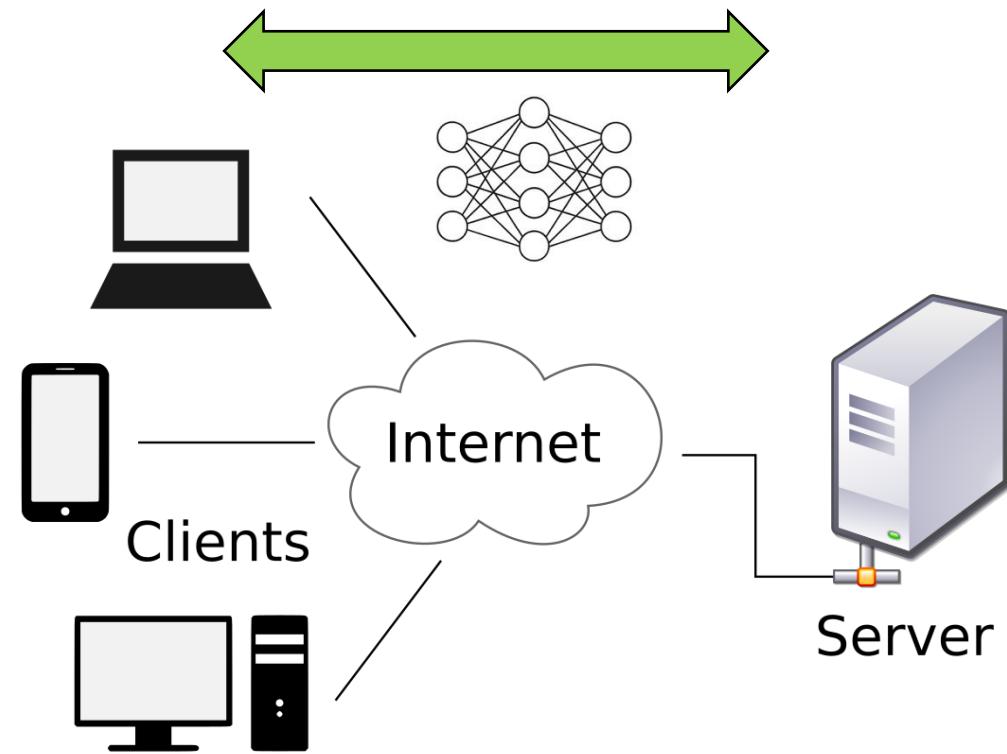
# Background

- Federated learning



# Background

- Federated learning



# Overview

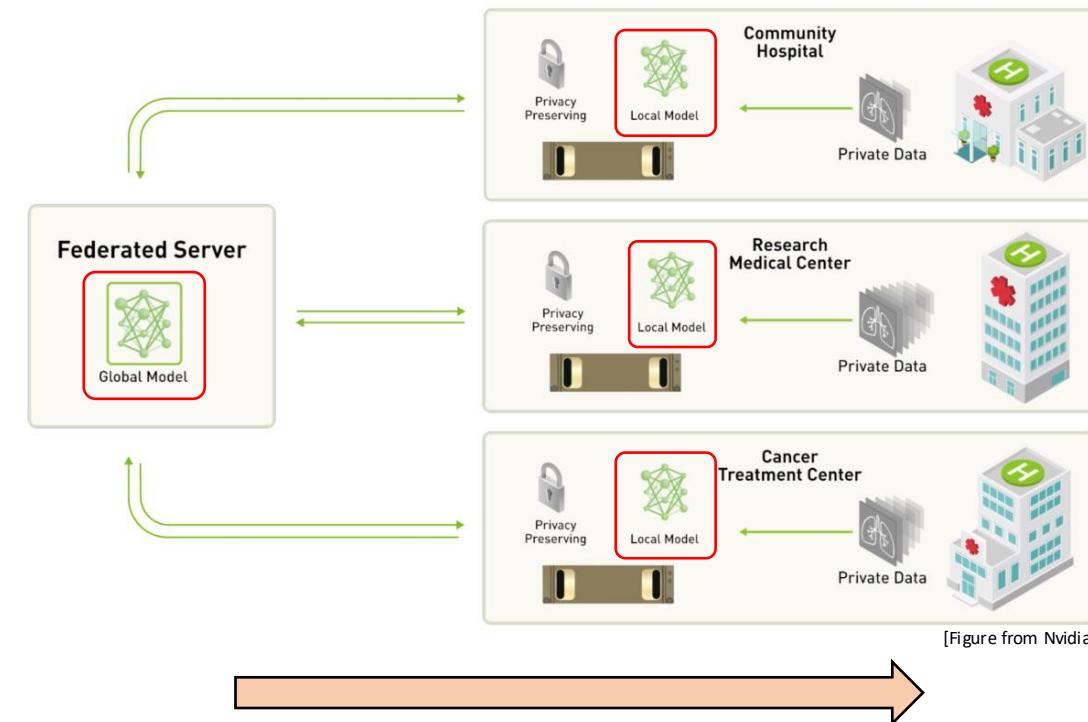
- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (**FedSLD**)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (**APPLE**)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (**PGFed**)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (**pFedMoAP**)
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- Summary

# Overview

- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (FedSLD)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- Summary

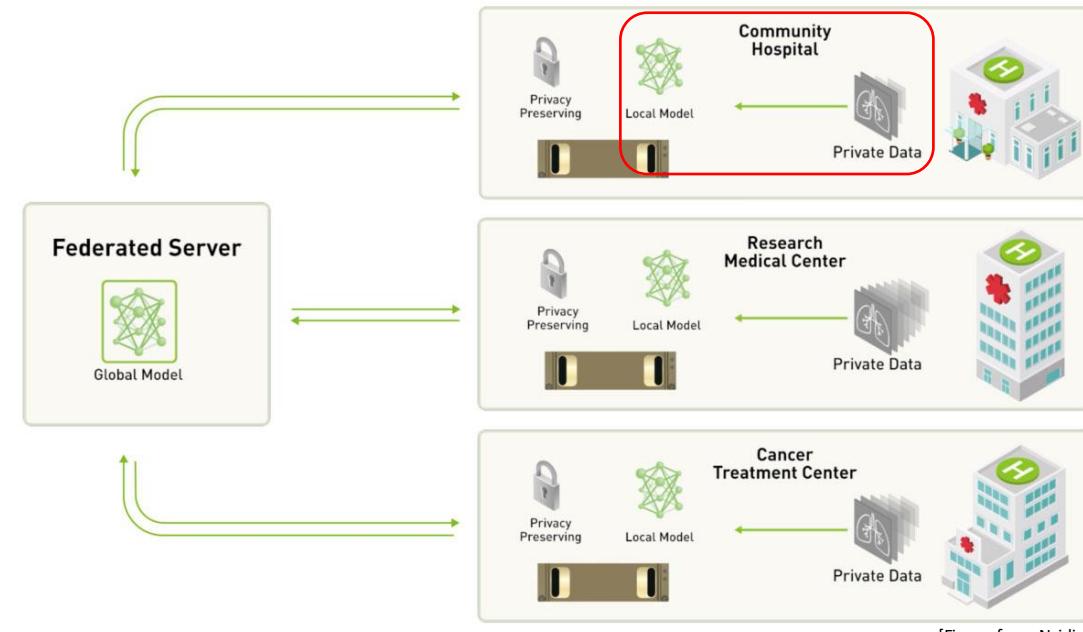
# Federated learning: introduction

- Basic mechanism of traditional FL
  - Broadcast



# Federated learning: introduction

- Basic mechanism of traditional FL
  - Broadcast
  - Client (local) training



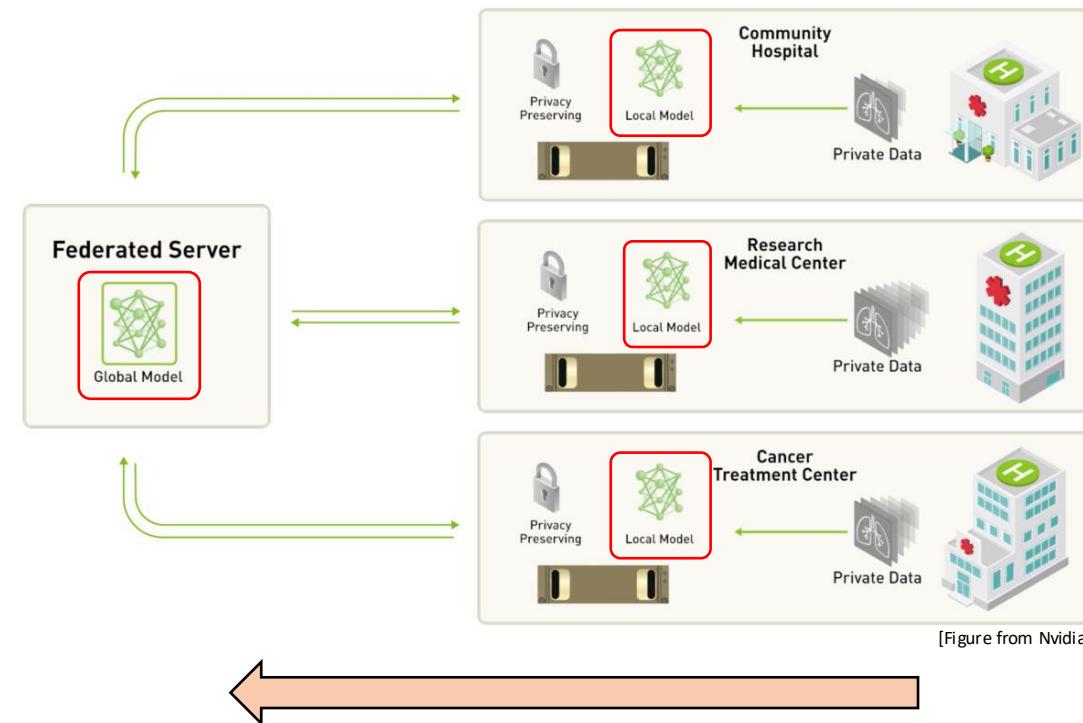
[Figure from Nvidia]

# Federated learning: introduction

- Basic mechanism of traditional FL

- Broadcast
- Client (local) training
- Server (global) aggregation

$$\text{FedAvg: } w = \sum_i p_i w_i$$

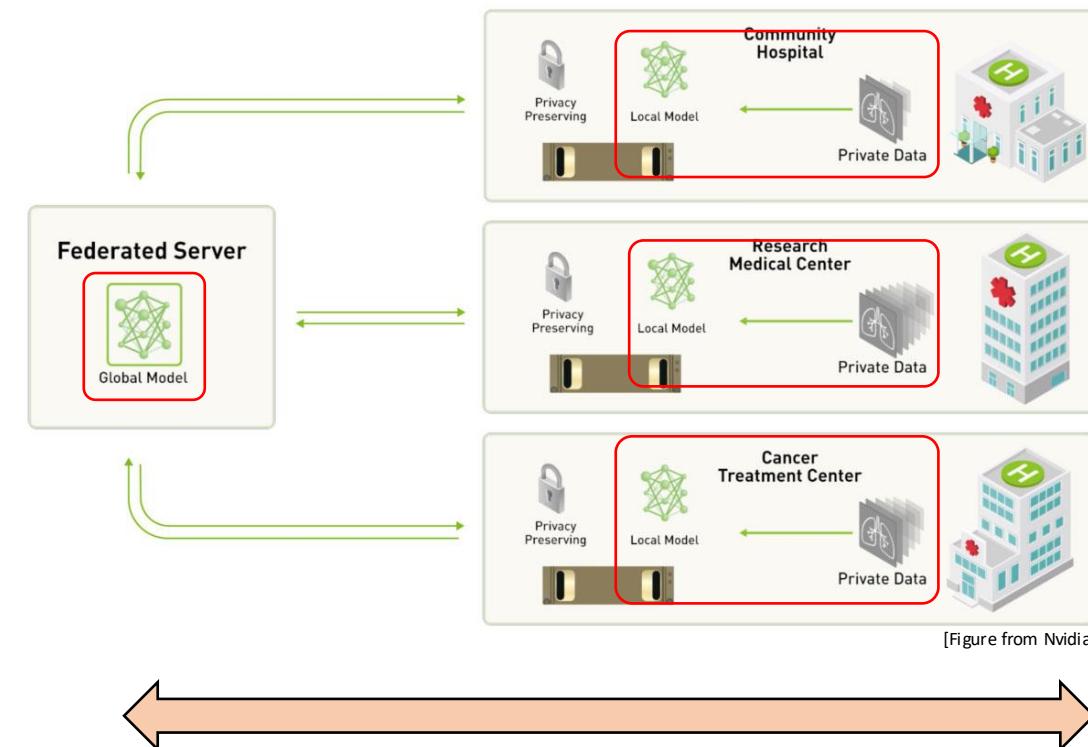


# Federated learning: introduction

- Basic mechanism of traditional FL

- Broadcast
- Client (local) training
- Server (global) aggregation

$$\text{FedAvg: } w = \sum_i p_i w_i$$



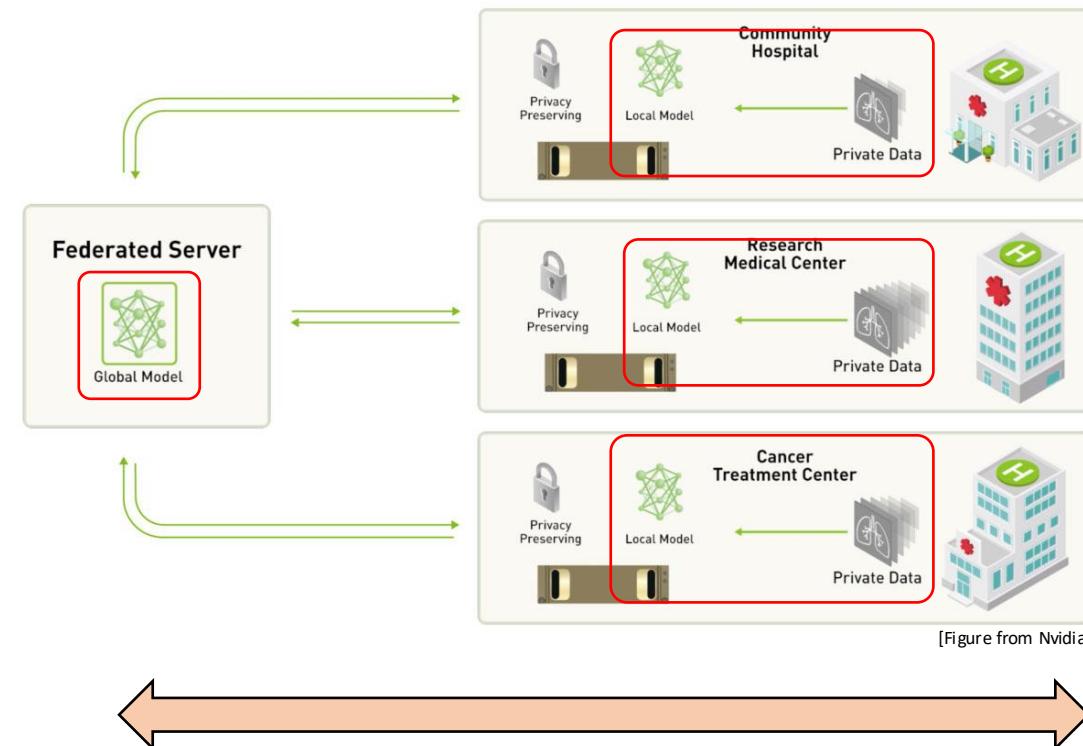
# Federated learning: introduction

- Basic mechanism of traditional FL

- Broadcast
- Client (local) training
- Server (global) aggregation

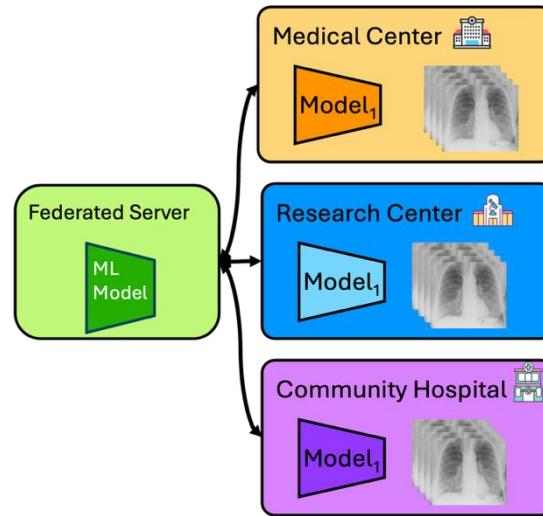
$$\text{FedAvg: } w = \sum_i p_i w_i$$

If #clients is large, sample clients at the beginning of each round

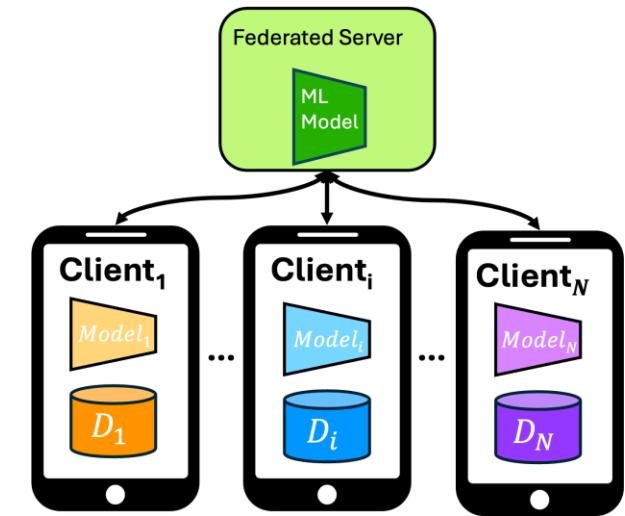


# Federated learning: introduction

- Applications of FL
  - Cross-silo FL
    - Medical centers
    - Financial institutes
  - Cross-device FL
    - Smart phone/IoT devices
    - Smart vehicle



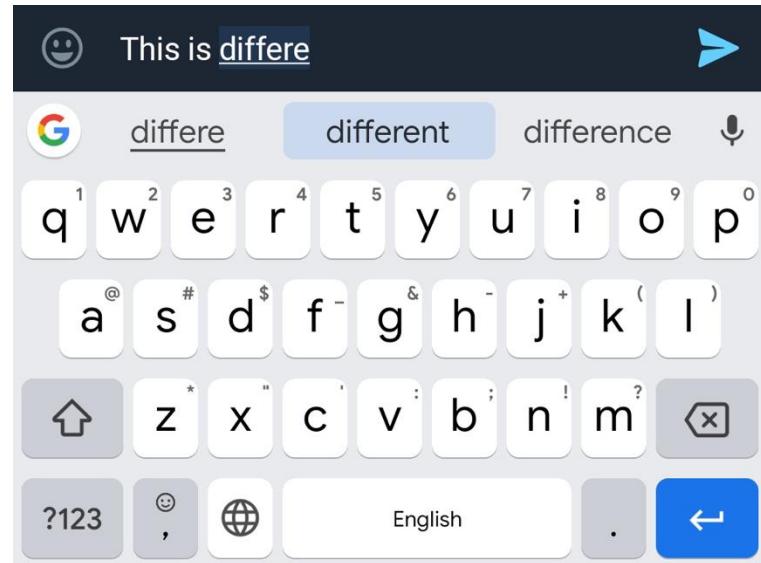
Cross-silo FL



Cross-device FL

# Federated learning: introduction

- Applications of FL
  - Cross-silo FL
    - Medical centers
    - Financial institutes
  - Cross-device FL
    - Smart phone/IoT devices
    - Smart vehicle
  - One of the earliest successes of FL: Gboard
  - “Hey Siri” from Apple, “Alexa” from Amazon...



[Figure from Google]



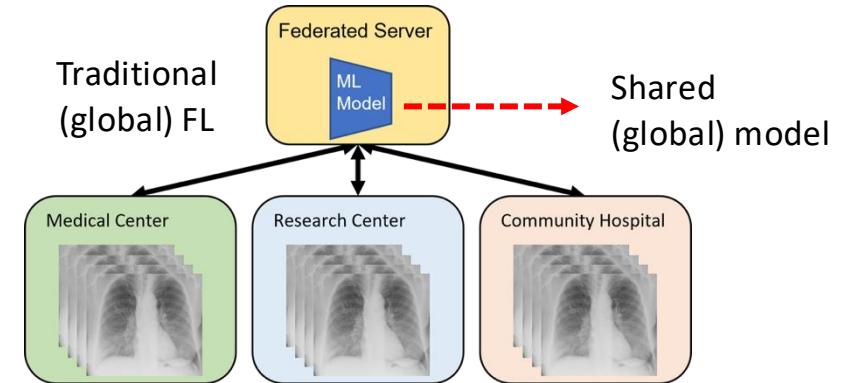
“Hey Siri”



“Alexa”

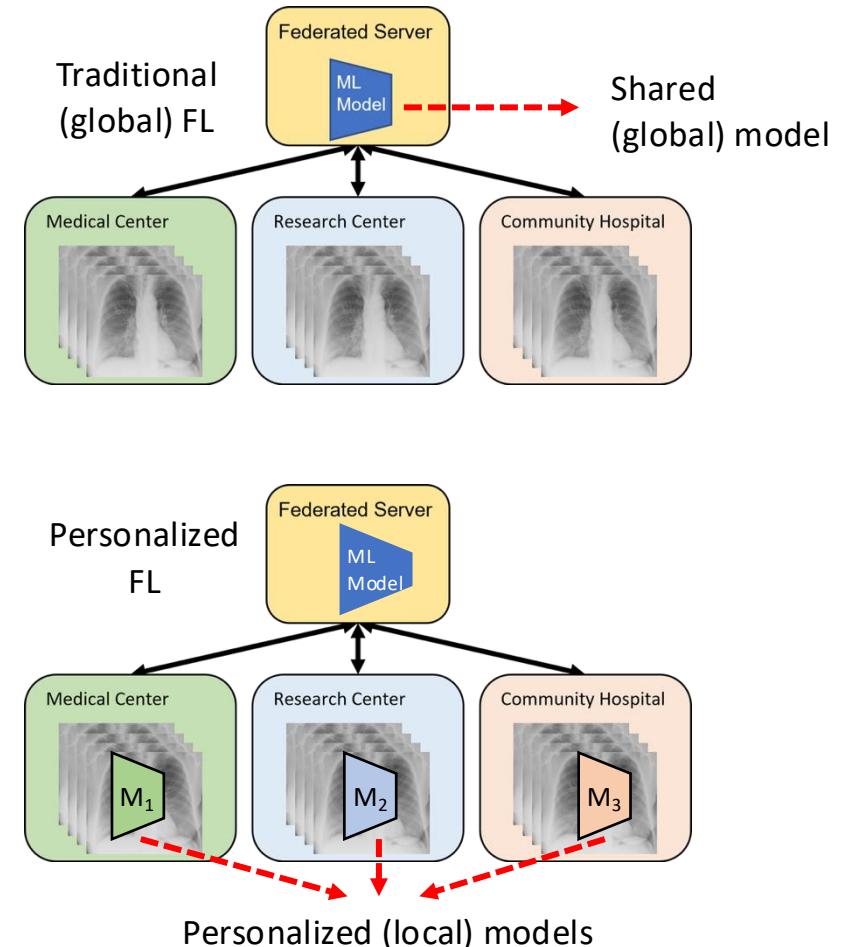
# Federated learning: introduction

- Data heterogeneity and personalized FL (PFL)
  - Data heterogeneity – non-IID
    - E.g. medical datasets are often non-IID
      - Different data acquisition protocols
      - Different local demographics



# Federated learning: introduction

- Data heterogeneity and personalized FL (PFL)
  - Data heterogeneity – non-IID
    - E.g. medical datasets are often non-IID
      - Different data acquisition protocols
      - Different local demographics
  - Traditional (global) FL
    - Trains a single global consensus model
    - Issues caused by data heterogeneity
      - inferior performance
      - slower convergence
      - Loss of clients' incentives to participate in FL
  - Personalized FL (PFL)
    - Allows customized models for different clients
    - Systemically mitigates data heterogeneity issue



# Overview

- Federated learning: introduction

Global FL

- Federated Learning with Shared Label Distribution for Medical Image Classification (**FedSLD**)

PFL

- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (**APPLE**)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (**PGFed**)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (**pFedMoAP**)
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection

- Summary

# Overview

- Federated learning: introduction

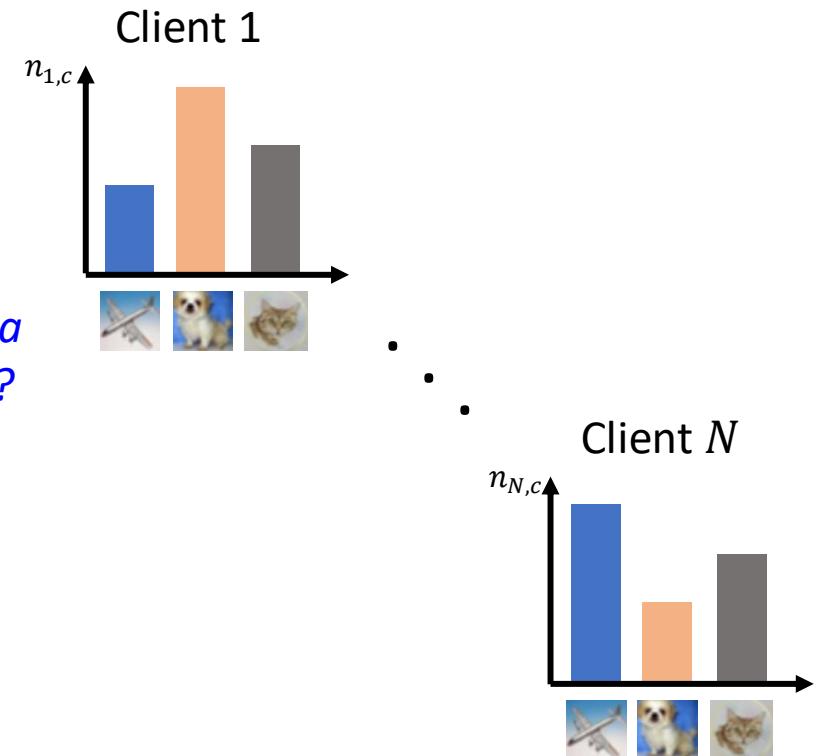
Global FL  
ISBI '22

- Federated Learning with Shared Label Distribution for Medical Image Classification (**FedSLD**)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- Summary

# FedSLD – Background and motivation

- FedAvg assumption
  - Weighted sum of local empirical risks
  - Weights are often  $n_i / \sum_j n_j$
  - Assumes knowledge of number of samples

**Research question 1:** *How can we leverage other sharable information to design a novel global FL algorithm for medical FL to mitigate the data heterogeneity issue?*

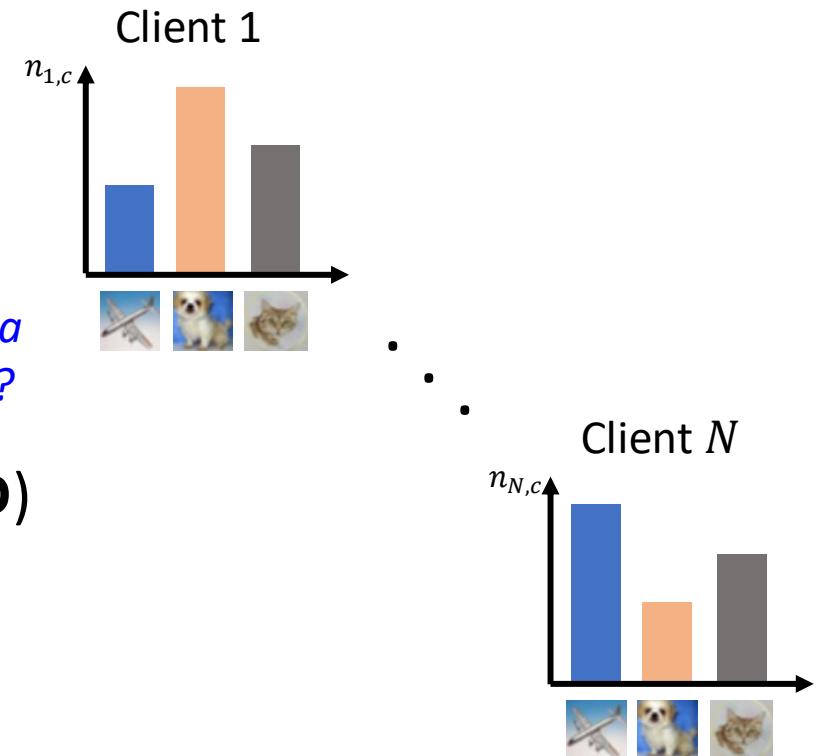


# FedSLD – Background and motivation

- FedAvg assumption
  - Weighted sum of local empirical risks
  - Weights are often  $n_i / \sum_j n_j$
  - Assumes knowledge of number of samples

**Research question 1:** *How can we leverage other sharable information to design a novel global FL algorithm for medical FL to mitigate the data heterogeneity issue?*

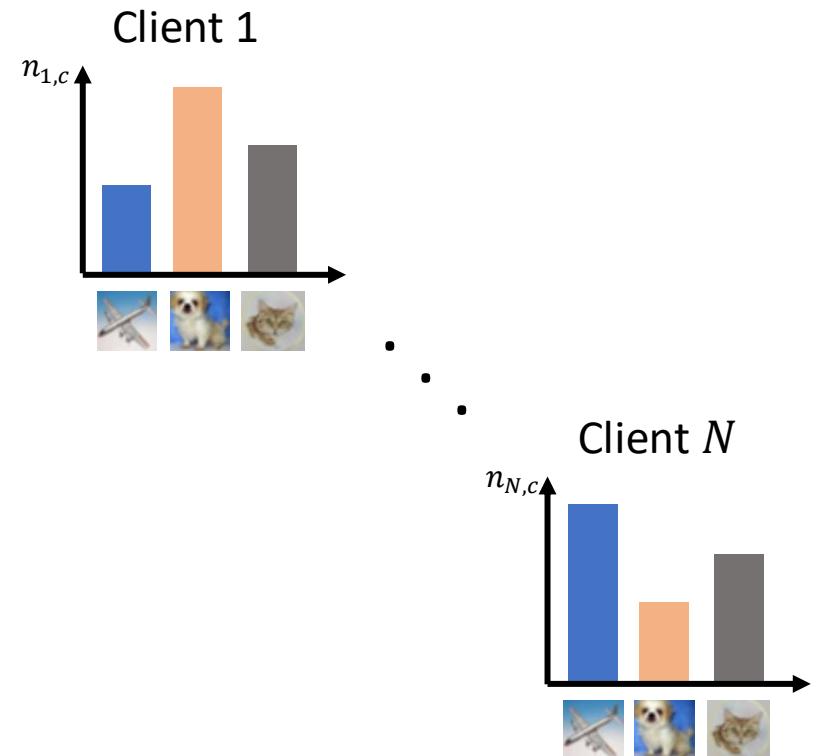
- **Federated Learning with Shared Label Distribution (FedSLD)**
  - Leverages information and statistics regarding the local datasets
  - Assumes knowing number of samples in each class
    - This assumption usually holds true for medical cross-silo FL
    - Estimate of label distribution



## FedSLD – Method

- Estimation of label distribution
  - Non-IID:  $\mathcal{P}_i(x, y) \neq \mathcal{P}_j(x, y)$
  - By Bayes' theorem,  $\mathcal{P}_i(x|y)\mathcal{P}_i(y) \neq \mathcal{P}_j(x|y)\mathcal{P}_j(y)$
  - Here, we only consider different  $\mathcal{P}_i(y) \neq \mathcal{P}_j(y)$
- Aggregate knowledge of #samples in each class, estimate  $\mathcal{P}(y)$  by

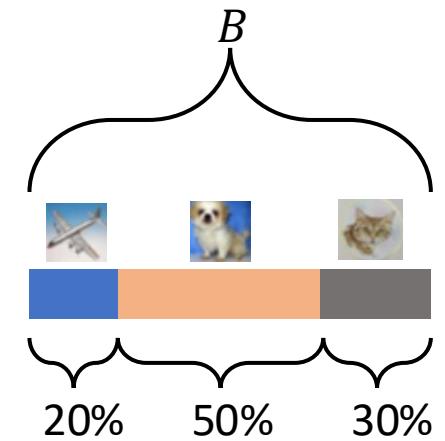
$$\tilde{\mathcal{P}}(y = c) = \frac{\sum_{i=1}^N n_{i,c}}{\sum_{i=1}^N n_i}$$



## FedSLD – Method

- Compute the percentage of each class in each mini-batch
  - During local update, given a batch of data  $\{(x_k, y_k)\}_{k=1}^B$  with  $B$  data samples, compute

$$p_b(y = c) = \frac{\sum_{k=1}^B \mathbb{I}[y_k = c]}{B}$$



# FedSLD – Method

- Weigh each data samples' contribution to the loss based on
  - The estimation of the prior of each class
  - The percentage of each class in each mini-batch
- Final loss of the mini-batch

$$\mathcal{L}_b\left(\{(x_k, y_k)\}_{k=1}^B\right) = - \sum_{k=1}^B \left( \frac{\tilde{P}(y = y_k)}{p_b(y = y_k)} \cdot \sum_{c=1}^C y_{k,c} \log(f_i(x_k))_c \right)$$

- Aggregate the model at the end of each training round as in FedAvg

---

**Algorithm 1 FedSLD.**

**Input:** Initialized model parameter weights  $w^0$ , number of clients  $N$ , number of local epochs  $E$ , batch size  $B$ , is the batch size, learning rate  $\eta$ , number of rounds  $R$ .

```

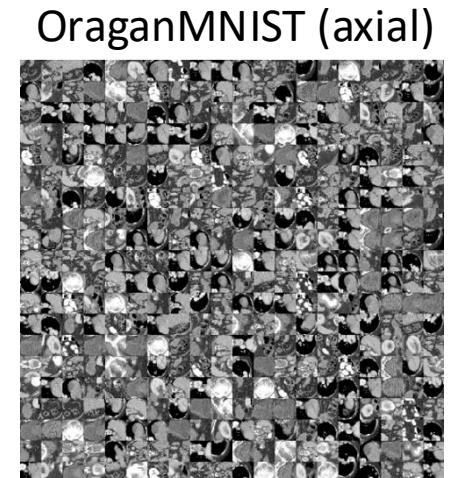
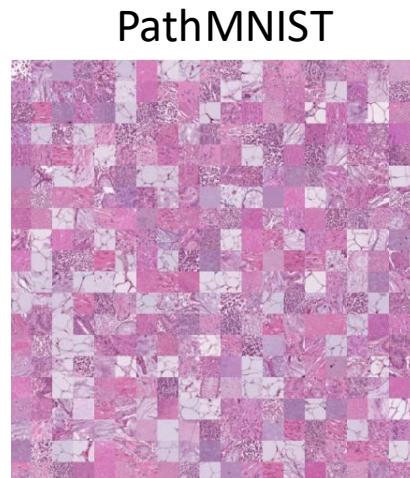
1:  $\forall i \in [N], c \in [C]$ , acquire  $n_{i,c}$ , client  $i$ 's numbers of
   samples of each class  $c$ .
2:  $\forall c \in [C], \tilde{P}(y = c) = \frac{\sum_{i=1}^N n_{i,c}}{\sum_{i=1}^N n_i}$  // compute estimated
   prior label distribution.
3: for  $r \leftarrow 1, 2, \dots, R$  do
4:    $\forall i \in [N] w_i^r = w^{r-1}$  // broadcast model parameters.
5:   for  $i \leftarrow 1, 2, \dots, N$  in parallel do
6:     for  $\{x_k, y_k\}_{k=1}^B$  in all minibatches do
7:        $\forall c, p_b(y = c) \leftarrow \sum_{k=1}^B \llbracket y_k = c \rrbracket / B$ 
8:       Compute loss  $\mathcal{L}_b$  by Equation 3.
9:        $w_i^r \leftarrow w_i^r - \eta \nabla_w \mathcal{L}_b$ 
10:    end for
11:   end for
12:    $w^r = \sum_{i=1}^N \frac{n_i}{n} w_i^r$  // aggregate model updates
13: end for
14: return  $w^R$ 

```

---

# FedSLD – Experiments & results

- Datasets
  - Two benchmark datasets
    - MNIST
    - CIFAR10
  - Two medical imaging datasets from MedMNIST collection
    - OrganMNIST (axial) (11-class liver tumor images)
    - PathMNIST (9-class colorectal cancer images)



## FedSLD – Experiments & results

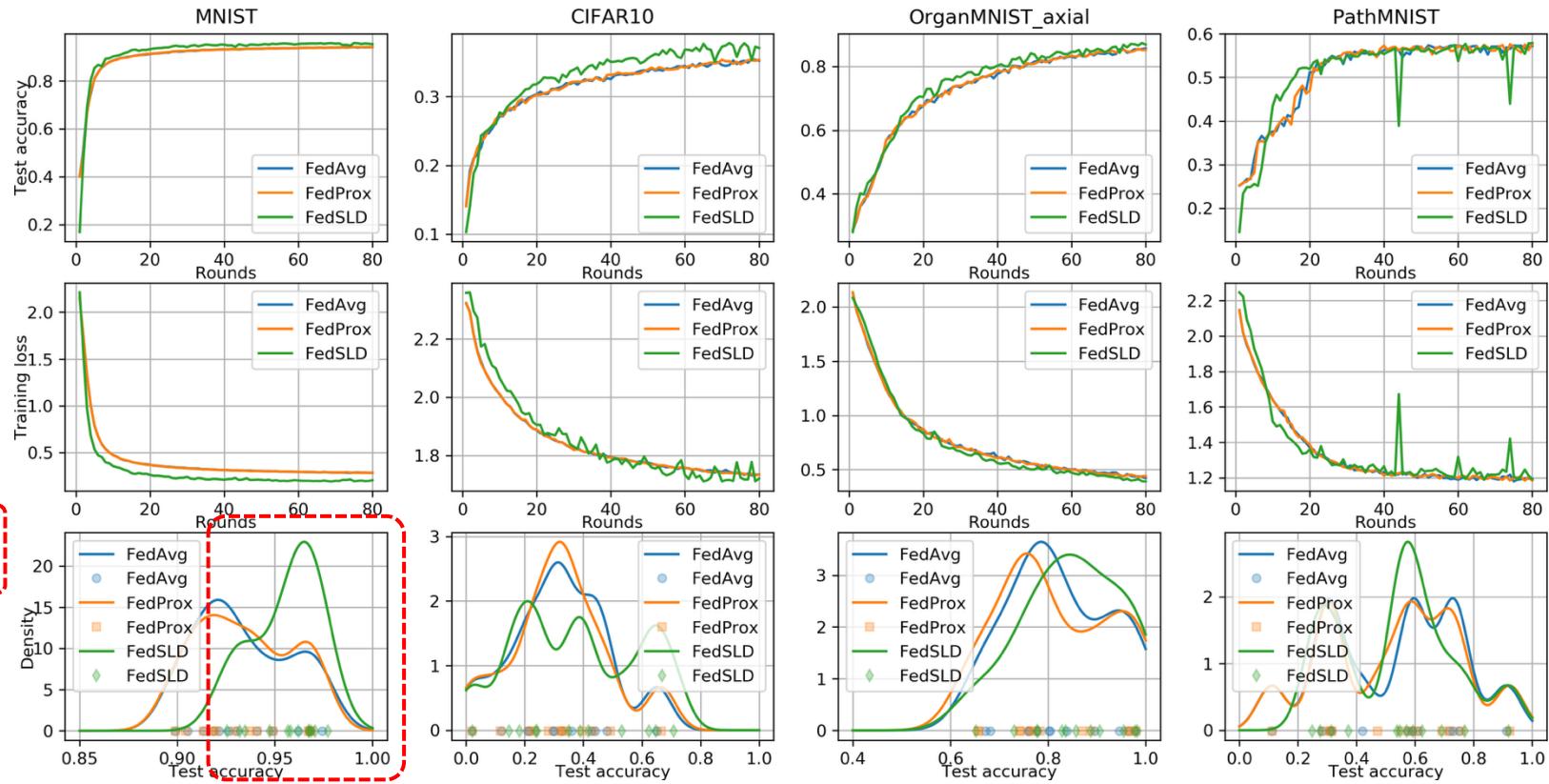
- Two non-IID settings
  - Pathological non-IID (12 clients)
    - Randomly select 2 classes for each client
    - In each class, assign a random number of images
  - Practical non-IID (12 clients)
    - Randomly partition each class of the dataset into 12 shards (10 x 1%, 1 x 10%, 1 x 80%)
    - Randomly assign one shard from each class to each client
    - A simulation that is closer to real-world medical applications
- Compared baselines
  - FedAvg
  - FedProx

# FedSLD – Experiments & results

- Practical non-IID results

Mean personalized acc. / Combined test set (global) acc.

| BMCTA/BTA     | MNIST              | CIFAR10            | Organ-MNIST        | Path-MNIST         |
|---------------|--------------------|--------------------|--------------------|--------------------|
| FedAvg        | 93.41/94.15        | 32.07/35.46        | 82.32/85.69        | 52.70/57.38        |
| FedProx       | 93.45/94.20        | 31.98/35.38        | 81.53/85.54        | 52.77/57.72        |
| FedSLD (Ours) | <b>95.56/95.85</b> | <b>37.48/37.79</b> | <b>84.75/84.75</b> | <b>53.87/57.90</b> |



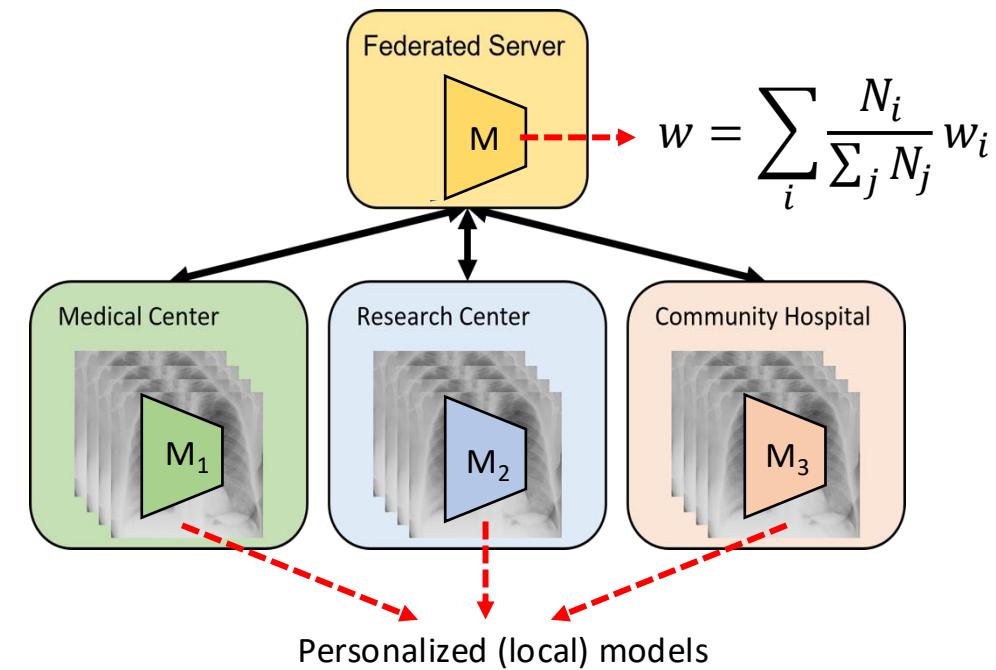
# Overview

- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (FedSLD)
- **Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)**
- PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- Summary

## APPLE – Background and motivation

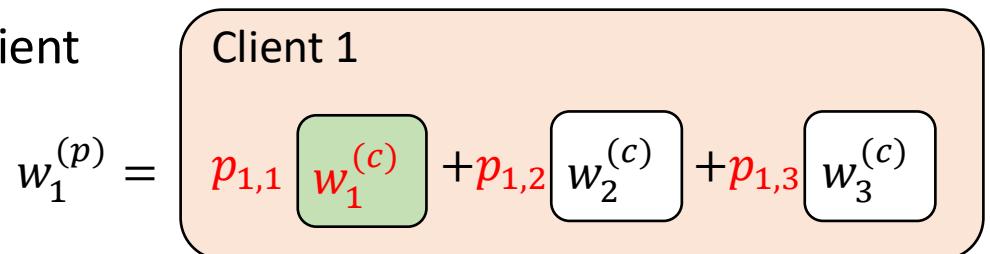
- FedAvg aggregation
  - $w = \sum_i p_i w_i$
  - $p_i = N_i / \sum_j N_j$ , aggregation weights are fixed
- Most existing FL/PFL methods
  - Use FedAvg-like aggregation
  - Training is either global or personalized

**Research question 2:** *How can we develop an **adaptive** aggregation strategy that optimally weighs different clients' contributions for each participant, while maintaining a **flexible balance** between global collaboration and local personalization objectives in cross-silo federated learning?*



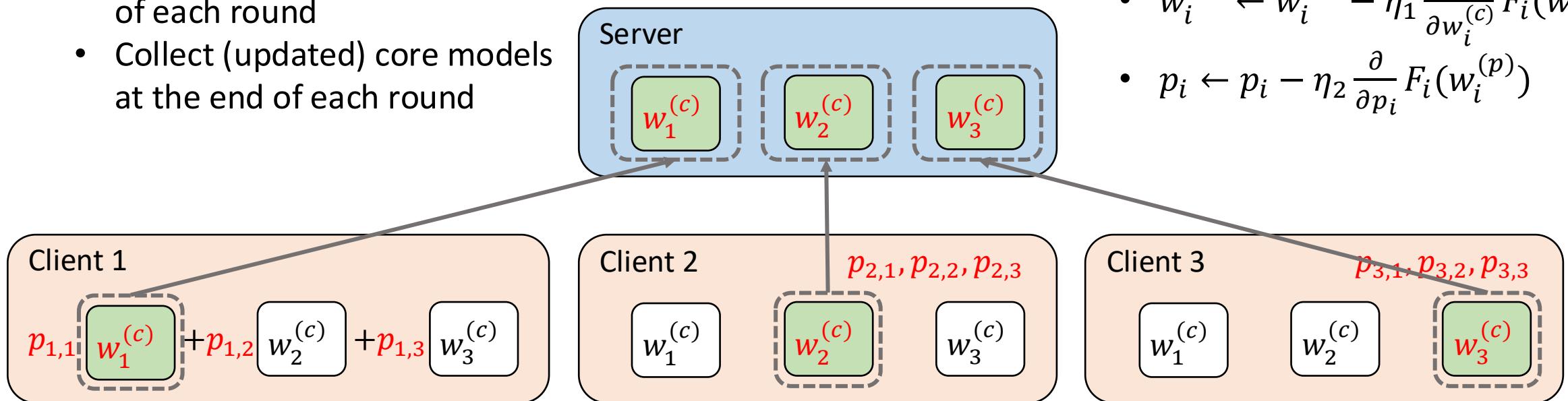
## APPLE – Background and motivation

- *Adaptive Personalized Cross-Silo Federated Learning (APPLE)*
- The model of a client
  - **Personalized model**  $w_i^{(p)}$ : used to do inference on client  $i$
  - **Core model**  $w_i^{(c)}$ : a constructing part of personalized model on client  $i$ , server also maintains core models from every client
  - $w_i^{(p)} = \sum_{j=1}^N p_{i,j} w_j^{(c)}$
  - **Directed relationship (DR) vector**  $p_i$ : learnable weights (coefficients for core models) on client  $i$ , always kept locally



# APPLE – Method

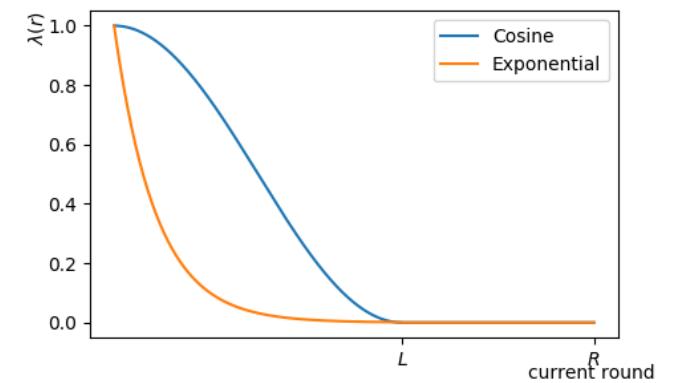
- Server
  - Broadcast core models to each client at the beginning of each round
  - Collect (updated) core models at the end of each round



- Local training
  - Clients' own core models and DR vectors are updated
  - $w_i^{(c)} \leftarrow w_i^{(c)} - \eta_1 \frac{\partial}{\partial w_i^{(c)}} F_i(w_i^{(p)})$
  - $p_i \leftarrow p_i - \eta_2 \frac{\partial}{\partial p_i} F_i(w_i^{(p)})$

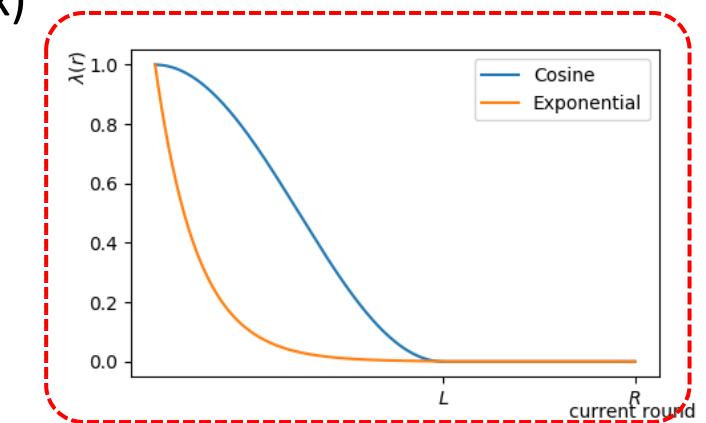
## APPLE – Method

- Proximal Directed Relationships
  - Since downloaded core models are not trained from local empirical risk, training might be drawn to resembling individual learning (DR matrix drawn to identity matrix)
  - Penalize DR vector by a proximal term
  - $$F_i(w_i^{(p)}) = \frac{1}{n_i} \sum_{\xi \in D_i^{tr}} \mathcal{L}(w_i^{(p)}; \xi) + \boxed{\lambda(r) \frac{\mu}{2} \|p_i - p_0\|_2^2}$$
  - Prox-center  $p_0 = [\frac{n_1}{n}, \dots, \frac{n_N}{n}]$
  - Loss scheduler  $\lambda(r) \in [0,1]$ : a decreasing function w.r.t. current round, controls the focus of training;  $\mu$ : the peak value of the proximal term coefficient
  - Proximal term coefficient:  $\infty \rightarrow$  FedAvg; large  $\rightarrow$  facilitate learning global high-level feature; small  $\rightarrow$  concentrate on local empirical risk, learning the personalization



## APPLE – Method

- Proximal Directed Relationships
  - Since downloaded core models are not trained from local empirical risk, training might be drawn to resembling individual learning (DR matrix drawn to identity matrix)
  - Penalize DR vector by a proximal term
  - $$F_i(w_i^{(p)}) = \frac{1}{n_i} \sum_{\xi \in D_i^{tr}} \mathcal{L}(w_i^{(p)}; \xi) + \lambda(r) \frac{\mu}{2} \|p_i - p_0\|_2^2$$
  - Prox-center  $p_0 = [\frac{n_1}{n}, \dots, \frac{n_N}{n}]$
  - Loss scheduler  $\lambda(r) \in [0,1]$ : a decreasing function w.r.t. current round, controls the focus of training;  $\mu$ : the peak value of the proximal term coefficient
  - Proximal term coefficient:  $\infty \rightarrow$  FedAvg; large  $\rightarrow$  facilitate learning global high-level feature; small  $\rightarrow$  concentrate on local empirical risk, learning the personalization



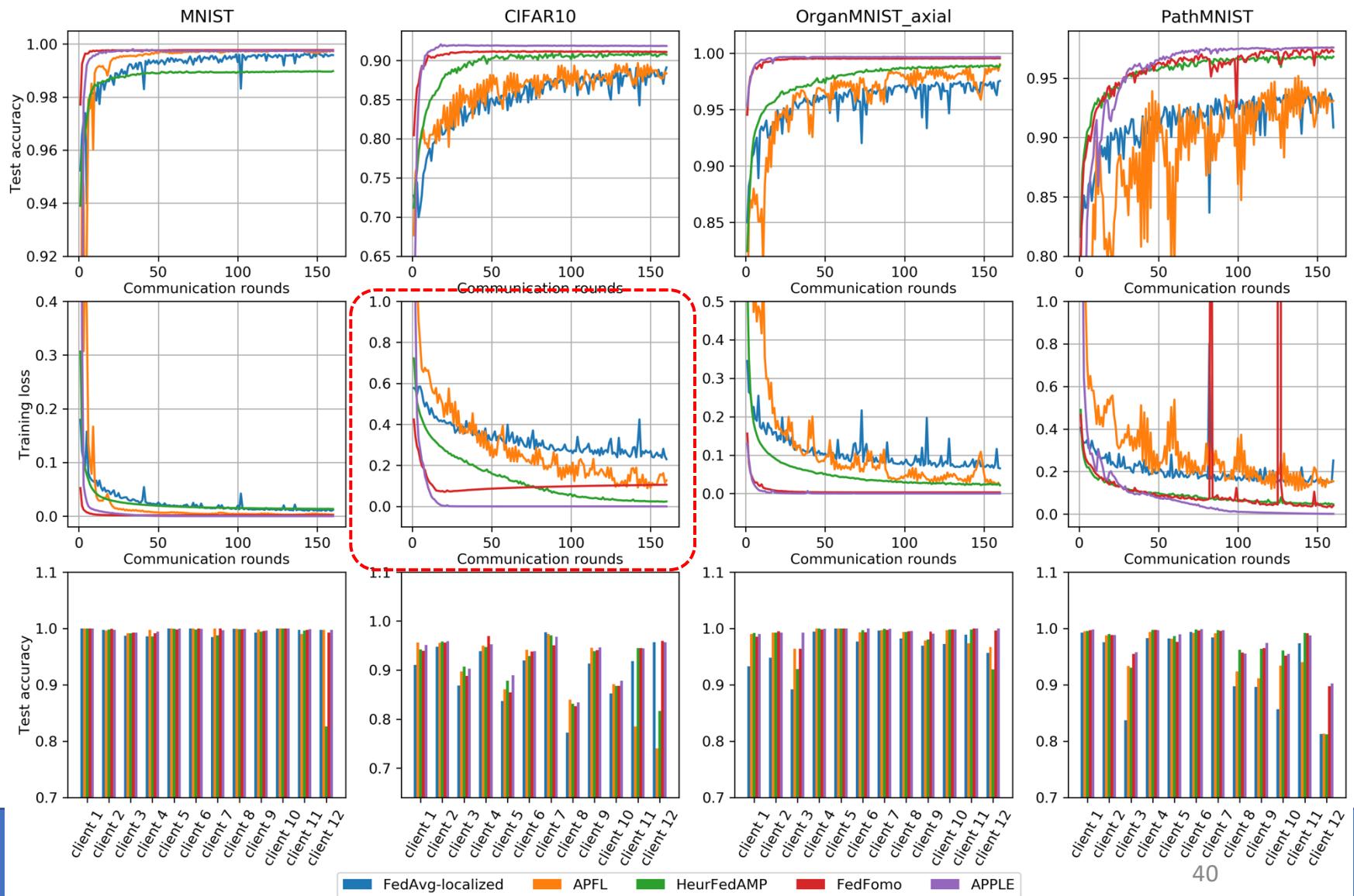
## APPLE – Experiments

- Datasets
  - MNIST
  - CIFAR10
  - OraganMNIST (axial)
  - PathMNIST
- Two non-IID settings (same with FedSLD)
  - Pathological non-IID
  - Practical non-IID
- Compared baselines
  - Separate training
  - FedAvg (McMahan et al., 2017)
  - FedAvg-local
  - FedAvg-FT, FedProx-FT (Wang et al., 2019)
  - APFL (Deng et al., 2020)
  - HeurFedAMP (Huang et al., 2021)
  - FedFomo (Zhang et al., 2021)

# APPLE – Results

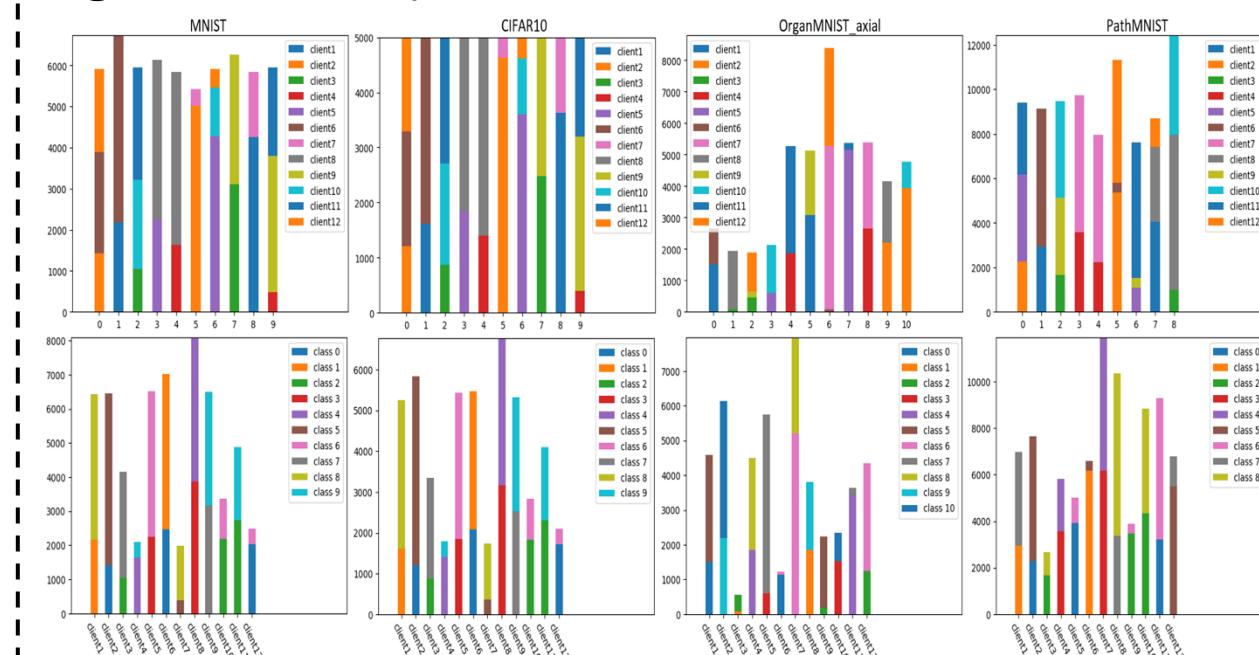
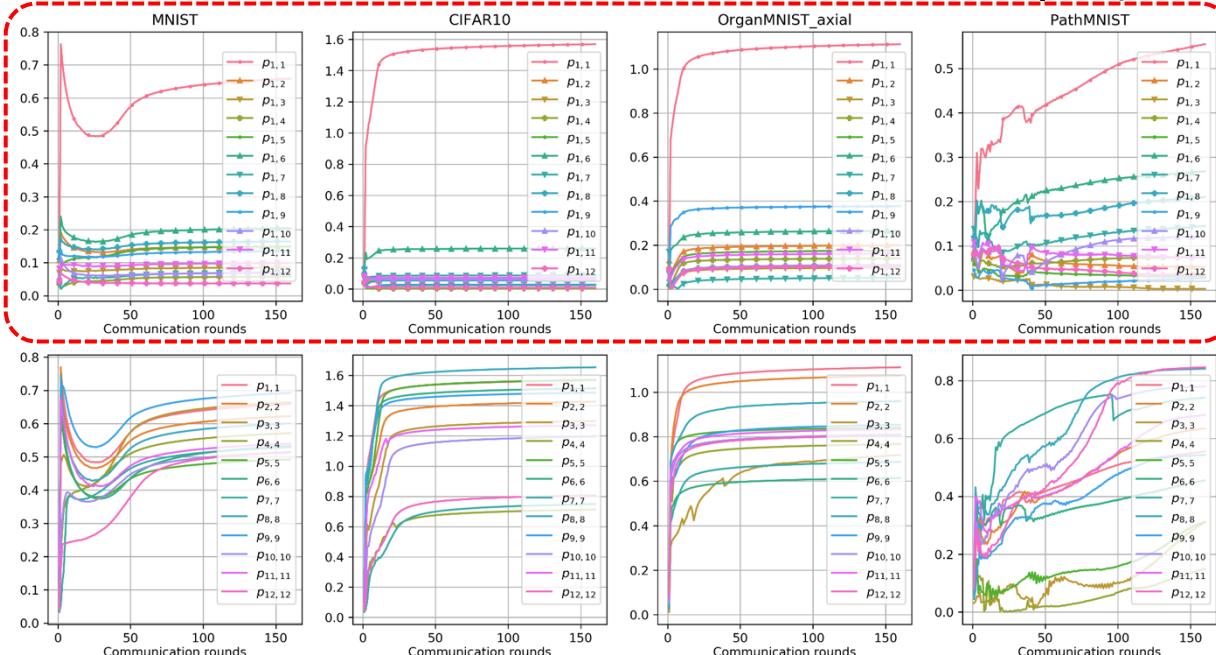
- Pathological non-IID

|                     | Pathological non-IID |              |                     |              |
|---------------------|----------------------|--------------|---------------------|--------------|
|                     | MNIST                | CIFAR10      | Organ-MNIST (axial) | Path-MNIST   |
| Separate            | 97.34                | 74.96        | 93.14               | 87.09        |
| FedAvg              | 95.71                | 51.44        | 59.43               | 56.61        |
| FedAvg-local        | 99.52                | 90.10        | 96.76               | 93.21        |
| FedAvg-FT           | 99.43                | 90.49        | 97.03               | 92.31        |
| FedProx-FT          | 99.43                | 90.49        | 97.03               | 92.38        |
| APFL                | 99.75                | 89.30        | 98.72               | 94.98        |
| HeurFedAMP          | 98.13                | 91.10        | 98.39               | 96.55        |
| FedFomo             | 99.71                | 91.96        | 99.31               | 97.24        |
| APPLE, $\mu = 0$    | 99.73                | 92.22        | <b>99.66</b>        | 96.78        |
| APPLE, $\mu \neq 0$ | <b>99.77</b>         | <b>92.68</b> | 99.61               | <b>97.51</b> |



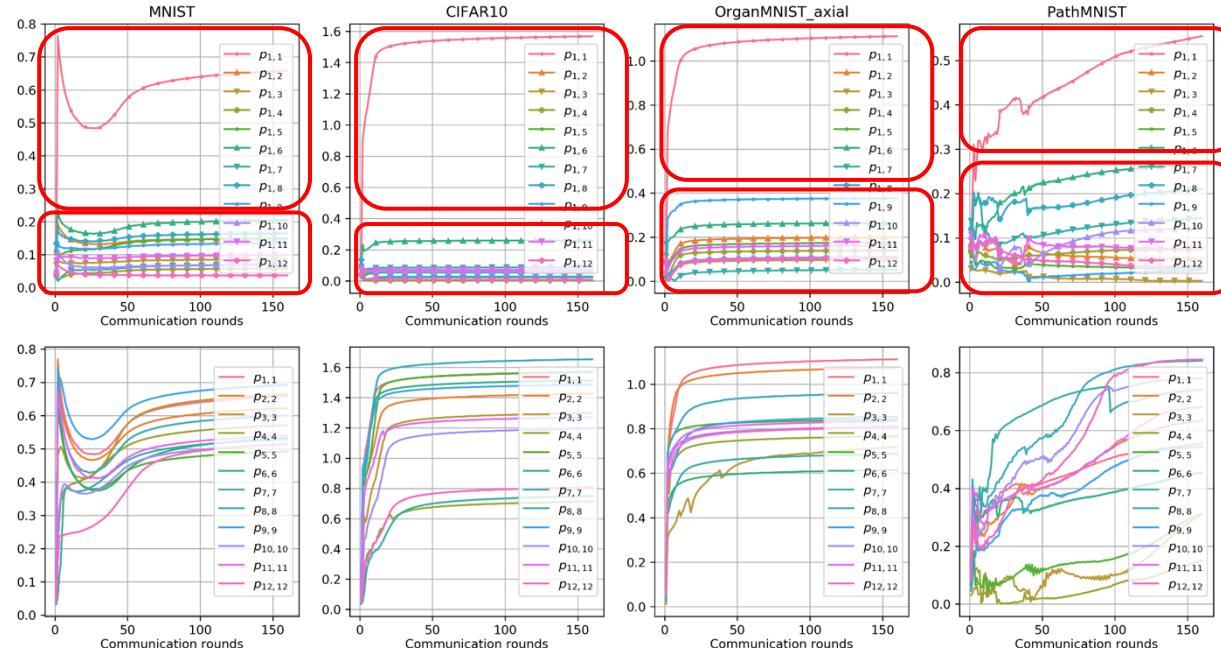
# APPLE – Results

- Visualization of Directed Relationships (Pathological non-IID)

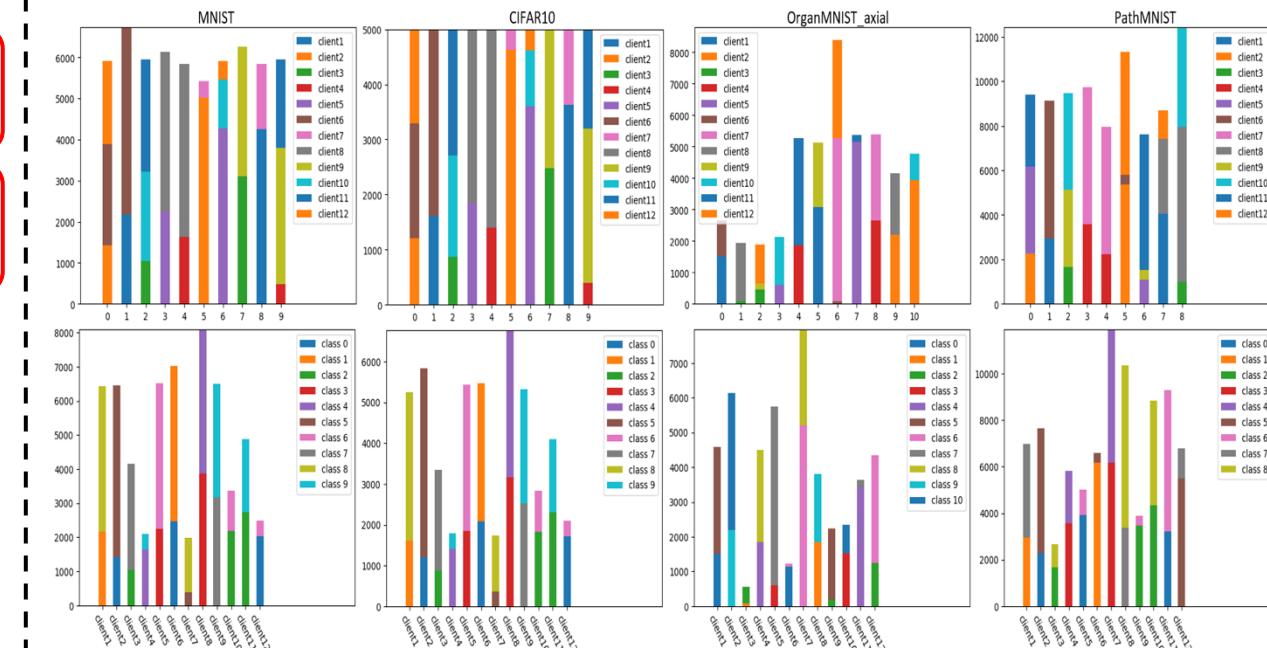


# APPLE – Results

- Visualization of Directed Relationships (Pathological non-IID)



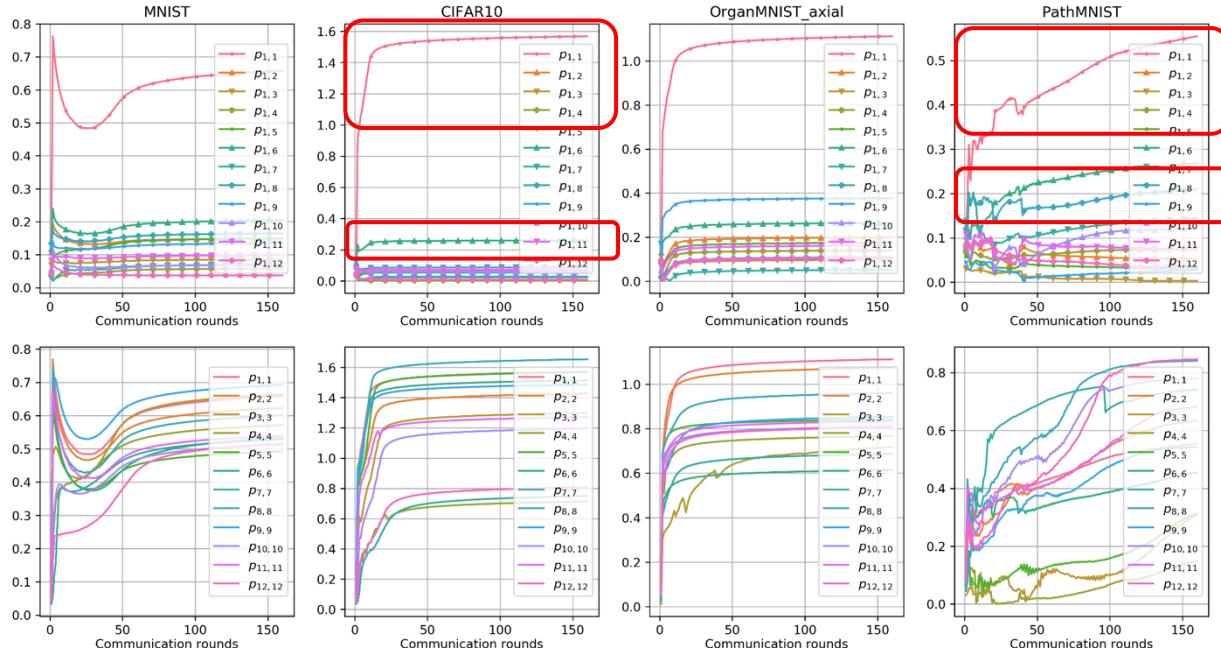
Visualization of DR



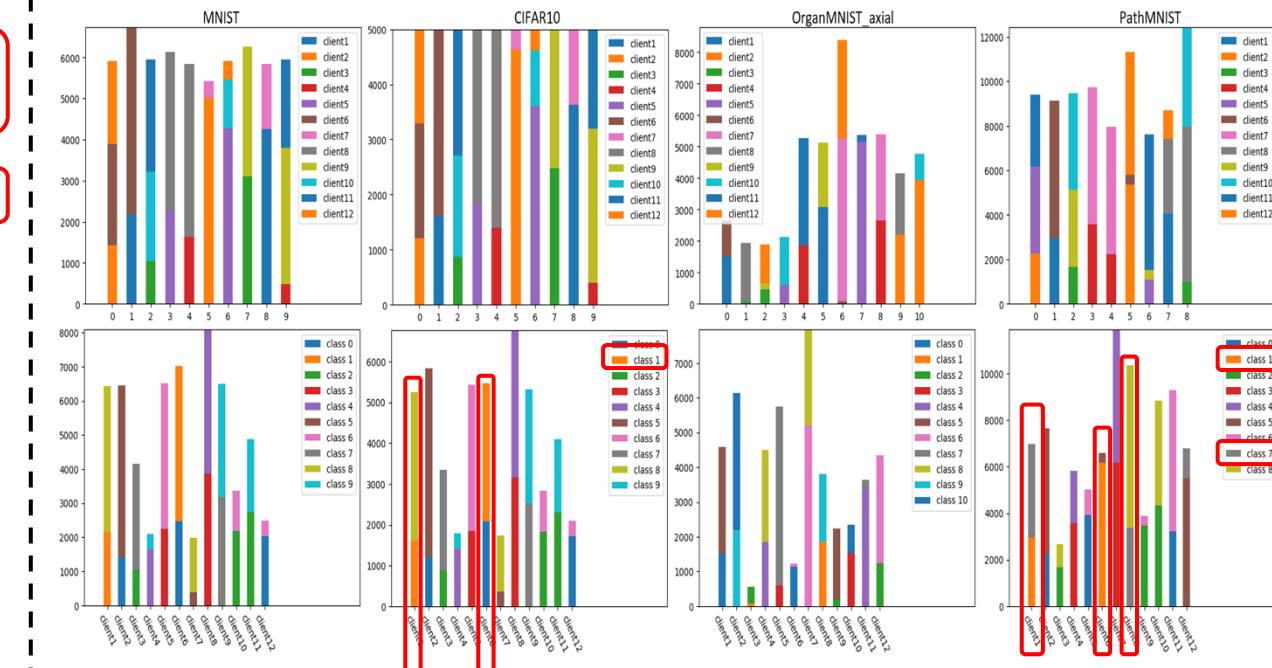
Data distribution

# APPLE – Results

- Visualization of Directed Relationships (Pathological non-IID)



Visualization of DR

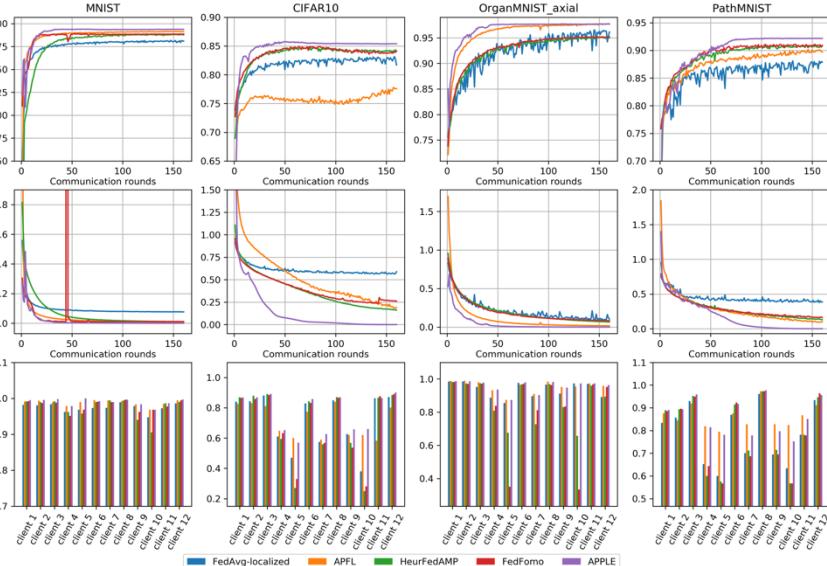


Data distribution

# APPLE – Results

- Practical non-IID

|                     | Practical non-IID |              |                     |              |
|---------------------|-------------------|--------------|---------------------|--------------|
|                     | MNIST             | CIFAR10      | Organ-MNIST (axial) | Path-MNIST   |
| Separate            | 78.20             | 63.06        | 65.21               | 61.36        |
| FedAvg              | 94.00             | 34.32        | 86.56               | 53.83        |
| FedAvg-local        | 97.47             | 71.99        | 93.75               | 78.70        |
| FedAvg-FT           | 97.66             | 72.08        | 94.13               | 78.69        |
| FedProx-FT          | 97.66             | 72.08        | 94.13               | 78.69        |
| APFL                | 98.80             | 71.19        | 95.53               | 86.35        |
| HeurFedAMP          | 97.45             | 69.54        | 86.82               | 79.33        |
| FedFomo             | 98.05             | 70.15        | 82.86               | 79.39        |
| APPLE, $\mu = 0$    | <b>99.00</b>      | 75.62        | <b>95.70</b>        | 84.22        |
| APPLE, $\mu \neq 0$ | 98.97             | <b>77.41</b> | 95.62               | <b>86.39</b> |



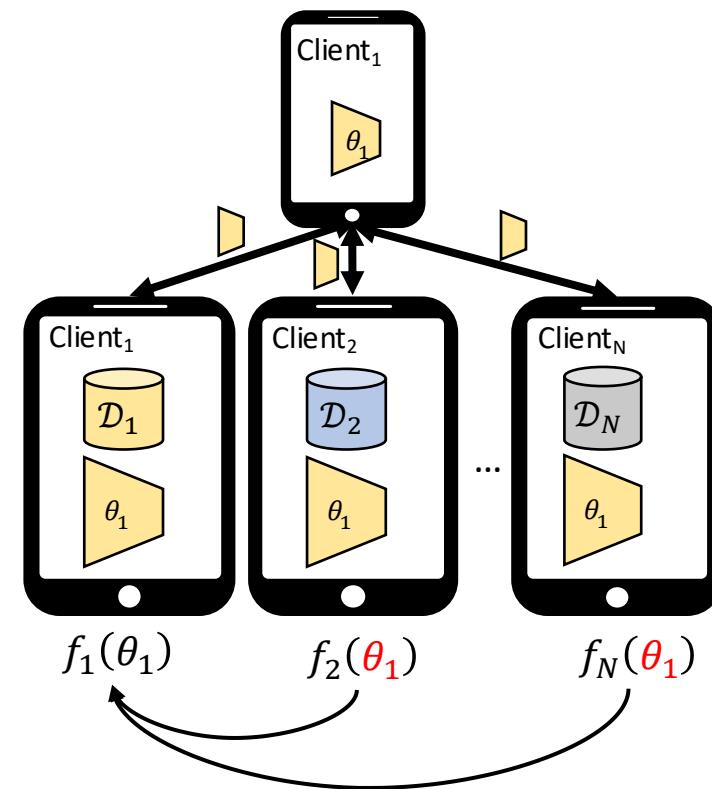
- Under limited bandwidth

|          | Pathological non-IID |         |                     |            | Practical non-IID |         |                     |            |
|----------|----------------------|---------|---------------------|------------|-------------------|---------|---------------------|------------|
|          | MNIST                | CIFAR10 | Organ-MNIST (axial) | Path-MNIST | MNIST             | CIFAR10 | Organ-MNIST (axial) | Path-MNIST |
| $M = 11$ | FedFomo              | 99.71   | 91.96               | 99.31      | 97.24             | 98.05   | 70.15               | 82.86      |
| $M = 11$ | APPLE                | 99.73   | 92.22               | 99.66      | 96.78             | 99.00   | 75.62               | 95.70      |
| $M = 7$  | FedFomo              | 99.71   | 91.95               | 99.31      | 97.33             | 97.65   | 70.24               | 80.88      |
| $M = 7$  | APPLE                | 99.73   | 92.17               | 99.53      | 97.15             | 98.70   | 76.14               | 94.21      |
| $M = 5$  | FedFomo              | 99.71   | 91.94               | 99.31      | 97.40             | 97.47   | 70.44               | 82.83      |
| $M = 5$  | APPLE                | 99.72   | 92.28               | 99.48      | 97.17             | 98.45   | 75.63               | 94.49      |
| $M = 2$  | FedFomo              | 99.71   | 91.98               | 99.31      | 97.25             | 96.51   | 69.87               | 79.53      |
| $M = 2$  | APPLE                | 99.70   | 92.41               | 99.47      | 97.11             | 98.29   | 74.84               | 92.29      |
| $M = 1$  | FedFomo              | 99.71   | 91.95               | 99.31      | 97.15             | 91.54   | 69.93               | 78.37      |
| $M = 1$  | APPLE                | 99.66   | 92.31               | 99.59      | 96.29             | 98.52   | 73.03               | 93.55      |
|          |                      |         |                     |            |                   |         |                     | 83.35      |

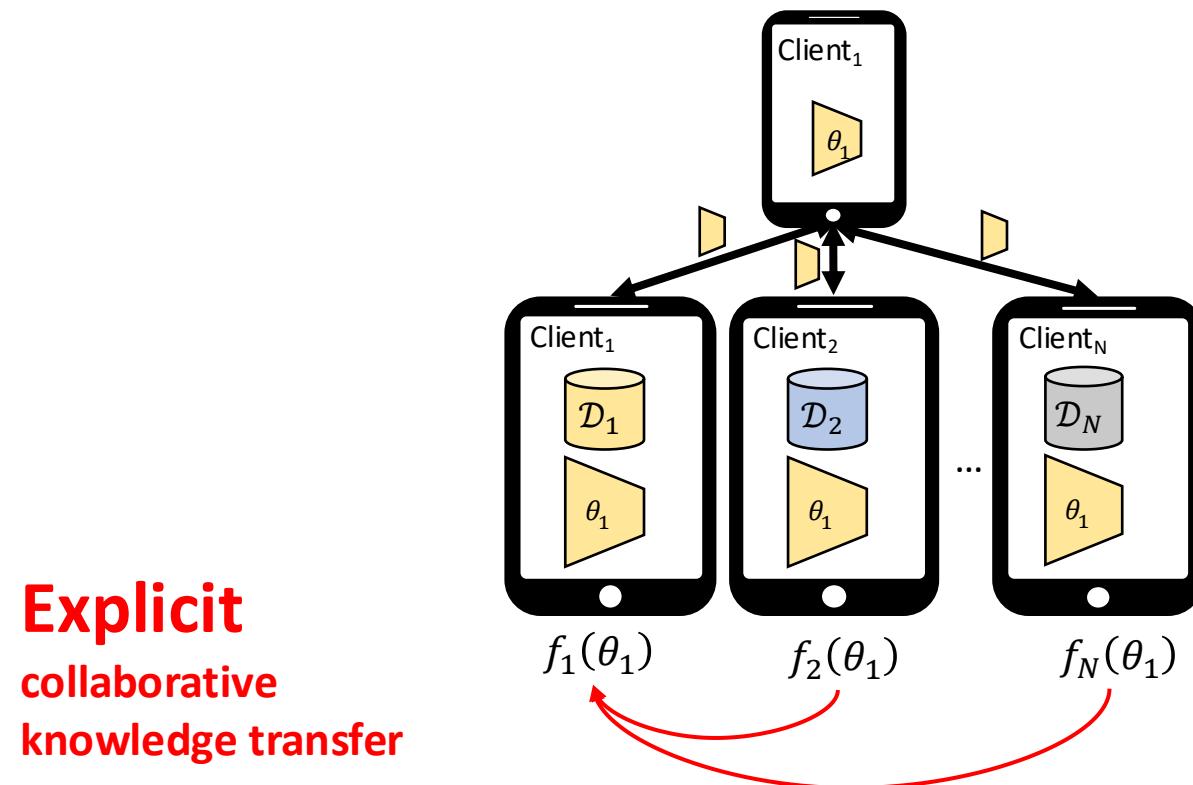
# Overview

- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (FedSLD)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)
- **PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)**
  - Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)
  - Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- Summary

# PGFed – Background and motivation

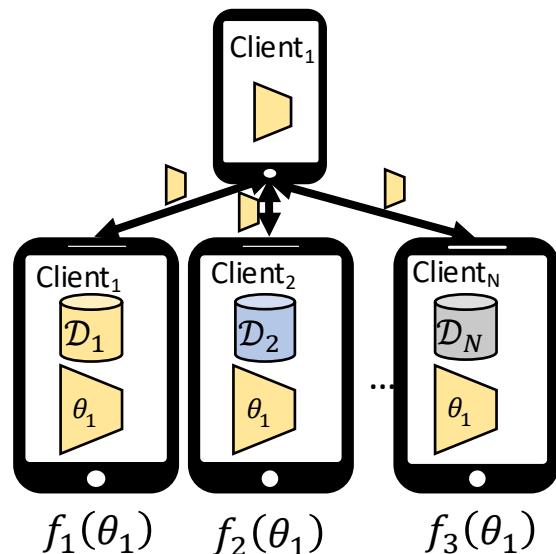


# PGFed – Background and motivation



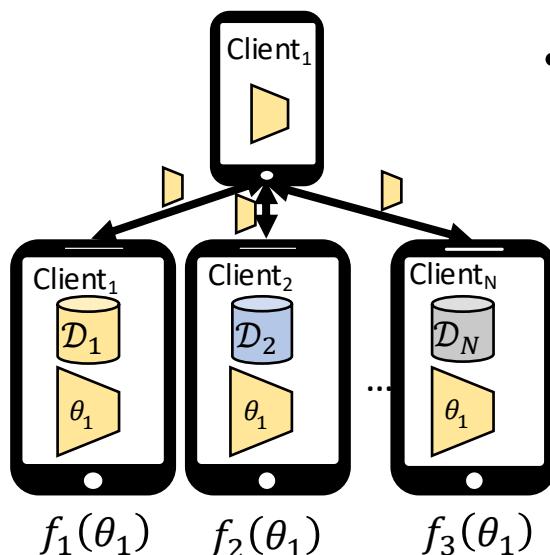
## PGFed – Background and motivation

- Why explicit (especially for personalized model update)?
  - **(Explicitness:** Direct engagement of multiple clients' empirical risks)
  - Intuition/motivation: facilitate the generalizability of  $\theta_i$  directly by penalizing its performance over other clients' empirical risks.

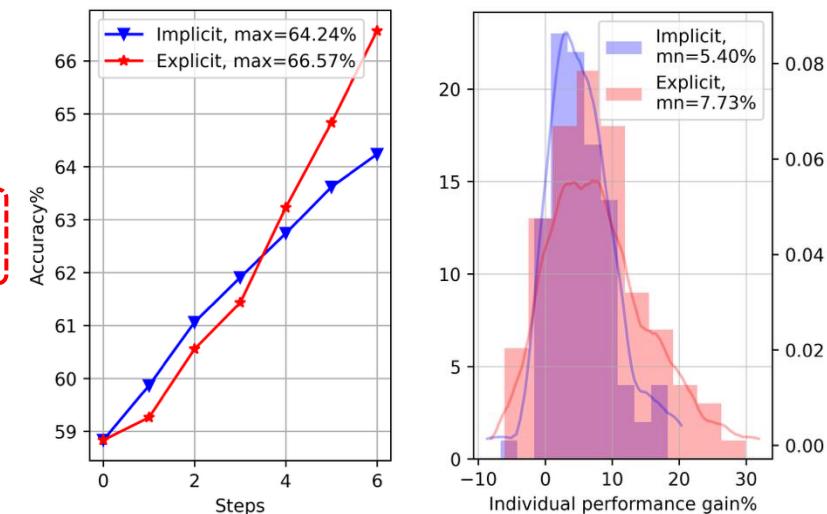


# PGFed – Background and motivation

- Why explicit (especially for personalized model update)?
  - **(Explicitness:** Direct engagement of multiple clients' empirical risks)
  - Intuition/motivation: facilitate the generalizability of  $\theta_i$  directly by penalizing its performance over other clients' empirical risks.

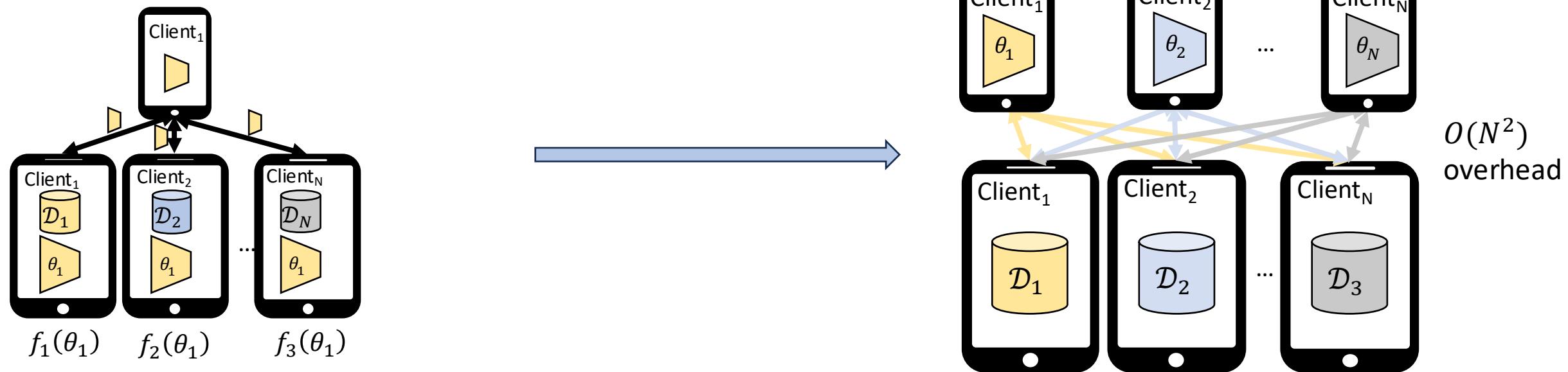


- Toy experiment on exemplar design
  - Cifar10, 100 heterogeneous clients
  - Explicit: 
$$F_i(\theta_i) = f_i(\theta_i) + \frac{\mu}{N-1} \sum_{j \neq i} f_j(\theta_i)$$
  - Implicit: 
$$F_i(\theta_i) = f_i(\theta_i)$$
 (local model of FedAvg)



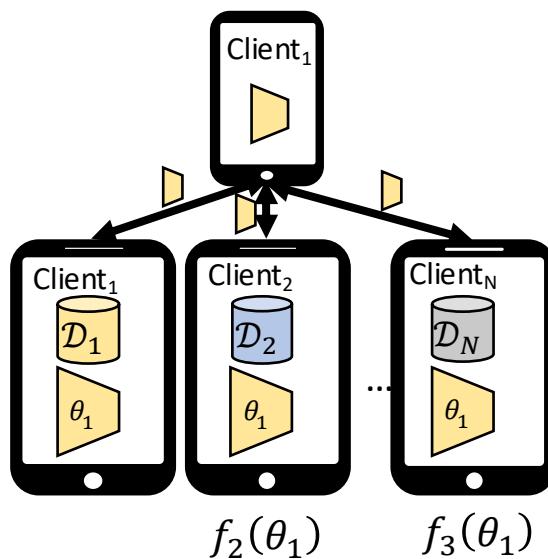
## PGFed – Background and motivation

- Why explicit (especially for personalized model update)?
  - **(Explicitness:** Direct engagement of multiple clients' empirical risks)
  - Intuition/motivation: facilitate the generalizability of  $\theta_i$  directly by penalizing its performance over other clients' empirical risks.

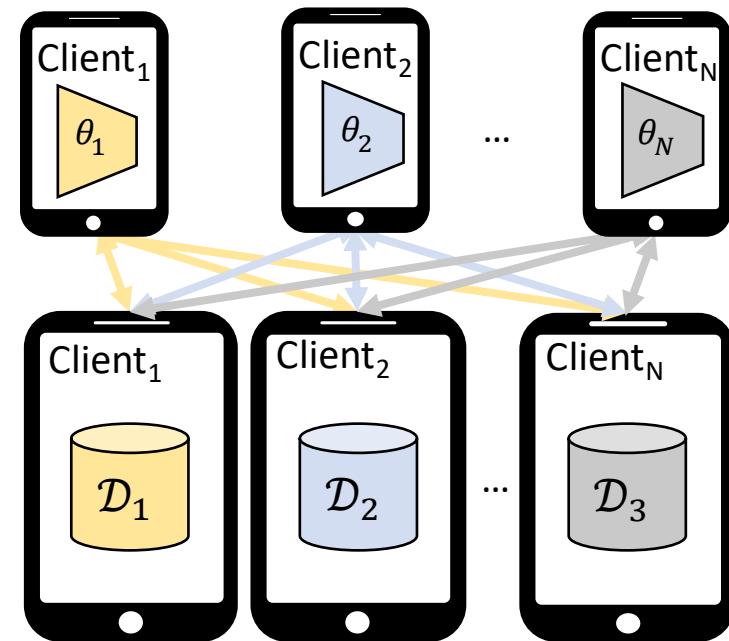


## PGFed – Background and motivation

- Why explicit (especially for personalized model update)?
  - **(Explicitness:** Direct engagement of multiple clients' empirical risks)
  - Intuition/motivation: facilitate the generalizability of  $\theta_i$  directly by penalizing its performance over other clients' empirical risks.

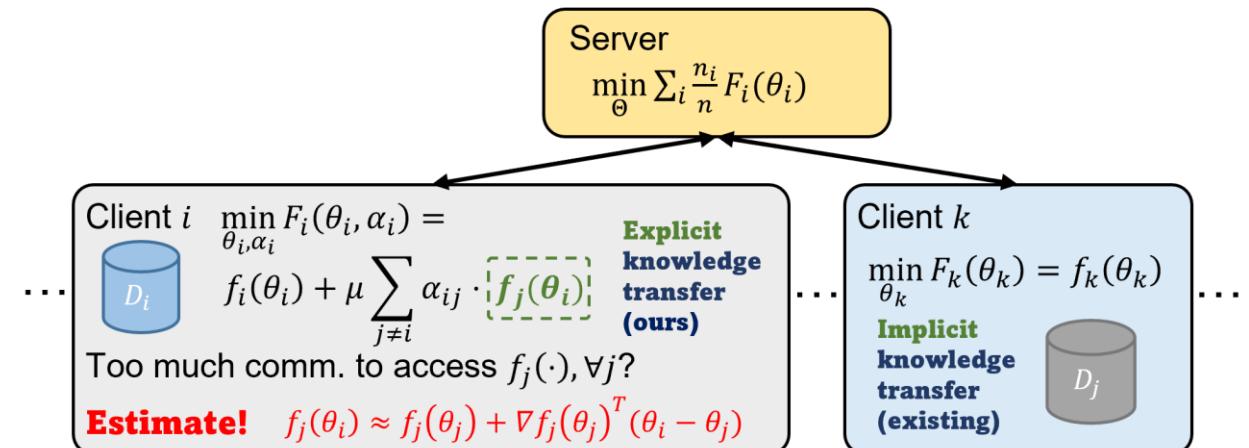
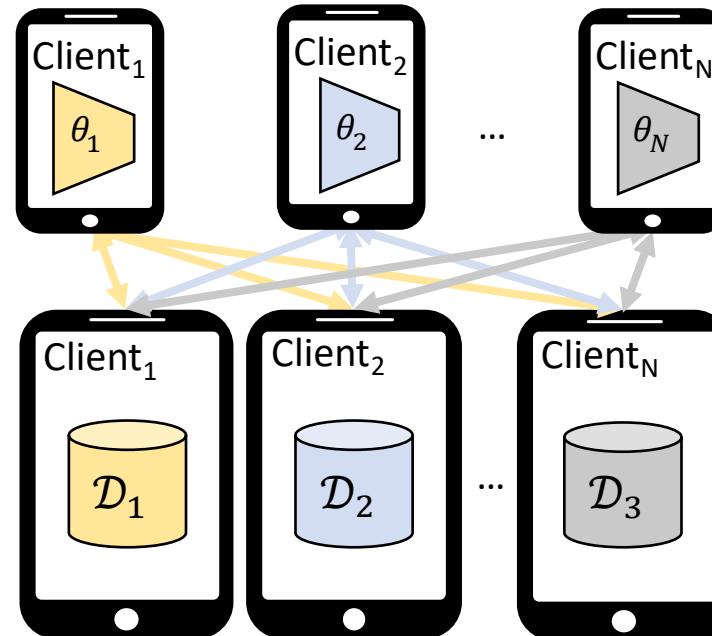


**Research question 3:** How can we design an **explicit** PFL framework to further boost the model performance with **linear communication** complexity that remains practical for both cross-silo and cross-device federated learning scenarios?



## PGFed – Background and motivation

- Difficulty to achieve explicitness
  - $O(N^2)$  communication overhead  $\longrightarrow \checkmark$  Estimate  $f_i(\theta_i) \approx f_i(\theta_i) + \nabla f_i(\theta_i)^T (\theta_i - \theta_j)$ ,  $O(N^2) \rightarrow O(N)$
  - Proper coefficient for each non-local risk  $\longrightarrow \checkmark$  Use adaptive coefficient  $\alpha_{ij} \forall i, j \in [N]$
- Proposed solution: **Personalized Global FL (PGFed)**



## PGFed – Method

- Objectives of Personalized Global Federated Learning (**PGFed**)

- Global objective:  $\min_{\Theta, A} F(\Theta, A) = \min_{\theta_1, \dots, \theta_N, \alpha_1, \dots, \alpha_N} \sum_{i=1}^N p_i F_i(\theta_i, \alpha_i)$

- Local objective:  $F_i(\theta_i, \alpha_i) = f_i(\theta_i) + \mu \sum_{j \in [N]} \alpha_{ij} f_j(\theta_i)$

- Plugging  $f_j(\theta_i) \approx f_j(\theta_j) + \nabla f_j(\theta_j)^T (\theta_i - \theta_j)$  into Local objective, we have

$$F_i(\theta_i, \alpha_i) \approx f_i(\theta_i) + \mathcal{R}_{aug}^{[N]}(\theta_i, \alpha_i)$$

$$\mathcal{R}_{aug}^{[N]}(\theta_i, \alpha_i) = \mu \sum_{j \in [N]} \alpha_{ij} (f_j(\theta_j) + \nabla_{\theta_j} f_j(\theta_j)^T (\theta_i - \theta_j))$$

$$\mathcal{R}_{aug}^{[N]}(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) = \mu \sum_{j \in [N]} \alpha_{ij} (f_j(\boldsymbol{\theta}_j) + \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j))$$

## PGFed – Method

- Gradient-based update

- W.r.t  $\theta_i$ : 
$$\begin{aligned} \nabla_{\boldsymbol{\theta}_i} F_i(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) &= \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i) + \nabla_{\boldsymbol{\theta}_i} \mathcal{R}_{aug}^{[N]}(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) \\ &= \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i) + \underbrace{\mu \sum_{j \in [N]} \alpha_{ij} \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)}_{\tilde{\mathbf{g}}_{[N]}}. \end{aligned}$$

- $\tilde{\mathbf{g}}_{[N]}$  can be computed by the server with:
  - Client  $i$  uploading  $\alpha_i$
  - Client  $j$  uploading local gradient

$$\mathcal{R}_{aug}^{[N]}(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) = \mu \sum_{j \in [N]} \alpha_{ij} (f_j(\boldsymbol{\theta}_j) + \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j))$$

## PGFed – Method

- Gradient-based update
  - W.r.t  $\alpha_{ij}$ : 
$$\begin{aligned} \nabla_{\alpha_{ij}} F_i(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) &= \mu (f_j(\boldsymbol{\theta}_j) + \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)) \\ &= \underbrace{\mu (f_j(\boldsymbol{\theta}_j) - \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T \boldsymbol{\theta}_j)}_{g_{\alpha}^{(1)}} + \underbrace{\mu \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T \boldsymbol{\theta}_i}_{g_{\alpha}^{(2)}}. \end{aligned}$$
  - $g_{\alpha}^{(1)}$  (a scalar) can be computed and uploaded by the client  $j$
  - $g_{\alpha}^{(2)}$  (exact value needs to transmit all gradients to client  $i$  (takes  $O(N^2)$  comm.))
    - Estimate: 
$$g_{\alpha}^{(2)} \approx \bar{\mathbf{g}}_{[N]}^T \boldsymbol{\theta}_i = \frac{\mu}{N} \left( \sum_{j \in [N]} \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j) \right)^T \boldsymbol{\theta}_i$$
    - Client  $j$  uploading local gradient

# PGFed – Method

- To accommodate to  $M$  selected clients per round:  $[N] \rightarrow S_t$  (selected set of clients in round  $t$ )

$$\tilde{g}_{S_t} = \mu \sum_{j \in S_t} \alpha_{ij} \nabla_{\theta_j} f_j(\theta_j) \quad \bar{g}_{S_t} = \frac{\mu}{M} \left( \sum_{j \in S_t} \nabla_{\theta_j} f_j(\theta_j) \right)$$

- To keep information from clients selected in previous round, use momentum (PGFedMo)

$$\tilde{g}_{S_t}^i = (1 - \beta) \tilde{g}_{S_t}^i \text{ (downloaded)} + \beta \tilde{g}_{S_t}^i \text{ (previous)}$$

---

**Algorithm 1** PGFed and PGFedMo
 

---

**Input:**  $N$  clients, learning rates  $\eta_1, \eta_2$ , number of rounds

$T$ , coefficient  $\mu$ , momentum  $\beta$  for PGFedMo

**Output:** Personalized models  $\theta_1^T, \dots, \theta_N^T$ .

**ServerExecute:**

```

1: Initialize  $\alpha_{ij} \leftarrow 1/M \forall i, j \in [N]$ , global model  $\theta_{glob}^0$ 
2:  $\mathbf{A}[i] \leftarrow \alpha_i \forall i \in [N]$ 
3: for  $t \leftarrow 1, 2, \dots, T$  do
4:   Select a subset of  $M$  clients,  $S_t$ 
5:    $g_t^{(1)} \leftarrow \{\}; \nabla_t \leftarrow \{\}$  // built for next round
6:   for  $i \in S_t$  in parallel do
7:     if  $t=1$  then
8:        $\theta_i^t, g_\alpha^{(1)}, \nabla f(\theta_i^t), \alpha_i \leftarrow \text{ClientUpdate}(\theta_{glob}^{t-1}, t)$ 
9:     else
10:       $\tilde{g}_{S_{t-1}} \leftarrow \mu \sum_{j \in S_{t-1}} \alpha_{ij} \nabla_{t-1}[j]$ 
11:       $\bar{g}_{S_{t-1}} \leftarrow \frac{\mu}{M} \sum_{j \in S_{t-1}} \nabla_{t-1}[j]$ 
12:       $\theta_i^t, g_\alpha^{(1)}, \nabla f(\theta_i^t), \alpha_i \leftarrow \text{ClientUpdate}(\theta_{glob}^{t-1}, t,$ 
13:       $\tilde{g}_{S_{t-1}}, \bar{g}_{S_{t-1}}, g_t^{(1)})$ 
14:    end if
15:    // the next line records the values for next round
16:     $\mathbf{A}[i] \leftarrow \alpha_i; g_t^{(1)}[i] \leftarrow g_\alpha^{(1)}; \nabla_t[i] \leftarrow \nabla f(\theta_i^t)$ 
17:     $\theta_{glob}^t \leftarrow \sum_{i \in S_t} p_i \theta_i^t$ 
18:  end for
19:  for  $i \in ([N] - S_t)$  in parallel do
20:     $\theta_i^t \leftarrow \theta_i^{t-1}; \tilde{g}_i^t \leftarrow \tilde{g}_i^{t-1}$ 
21:  end for
22: return  $\theta_1^T, \dots, \theta_N^T$ 

```

---

**ClientUpdate**( $\theta_{global}^{t-1}, t$  ,  $\tilde{g}, \bar{g}, g_t^{(1)}$ ):

```

1: if  $t=1$  then
2:    $\theta_i^t \leftarrow \text{ClientUpdate}(\theta_{global}^{t-1}, \eta_1)$  as in FedAvg
3: else
4:    $\theta_i^t \leftarrow \theta_{global}^{t-1}$ 
5:    $\tilde{g}_i^t \leftarrow \tilde{g}$  // without momentum
6:    $\tilde{g}_i^t \leftarrow (1 - \beta) \tilde{g} + \beta \tilde{g}_i^{t-1}$  // with momentum
7:   for Batch of data  $\mathcal{B} \in \mathcal{D}_i$  do
8:      $\theta_i^t \leftarrow \theta_i^t - \eta_1 (\nabla f(\theta_i^t, \mathcal{B}) + \tilde{g}_i^t)$ 
9:      $g^{(2)} = \bar{g}^T \theta_i$ 
10:     $\forall j \in g_t^{(1)} : \alpha_{ij} \leftarrow \alpha_{ij} - \eta_2 (g_{t-1}^{(1)}[j] + g^{(2)})$ 
11:  end for
12: end if
13:  $g_\alpha^{(1)} \leftarrow \mu (f(\theta_i^t) - \nabla f(\theta_i^t)^T \theta_i^t)$  // for next round
14: return  $\theta_i^t, g_\alpha^{(1)}, \nabla f(\theta_i^t), \alpha_i$ 

```

# PGFed – Experiments & results

- Settings
  - Datasets: CIFAR10, CIFAR100, OrganMNIST, Office-home
  - Partition
    - CIFAR10/100:  $\text{Dir}(\alpha = 0.3)$ , 25, 50, 100 clients, 25% sample rate
    - OrganMNIST: 25 clients,  $\text{Dir}(\alpha = 1.0)$ , 50% sample rate  
50, 100 clients,  $\text{Dir}(\alpha = 0.3)$ , 25% sample rate
    - Office-home: 5 clients/domain x 4 domains,  $\text{Dir}(\alpha = 0.3)$ , 25% sample rate
  - Metric: mean personalized test accuracy
  - Compared methods
 

|           |              |
|-----------|--------------|
| • Local   | • FedReP     |
| • FedAvg  | • LG-FedAvg  |
| • FedDyn  | • FedPer     |
| • pFedMe  | • Per-FedAvg |
| • FedFomo | • FedRoD     |
| • APFL    | • FedBABU    |

Heterogeneous partition of a dataset based on  
Dirichlet distribution:

- $\alpha = \infty \rightarrow$  homogeneous
- $\alpha = 0.3/0.5/1.0 \rightarrow$  very heterogeneous, with  
1.0 slightly balanced (tend to have lower acc.)
- $\alpha = 0 \rightarrow$  one class per client

# PGFed – Experiments & results

- Performance on CIFAR10 & CIFAR100

|                       | CIFAR10                            |                                    |                                    | CIFAR100                           |                                    |                                    |
|-----------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
|                       | 25 clients                         | 50 clients                         | 100 clients                        | 25 clients                         | 50 clients                         | 100 clients                        |
| Local                 | 72.40 $\pm$ 0.45                   | 70.28 $\pm$ 0.38                   | 67.39 $\pm$ 0.20                   | 32.74 $\pm$ 0.08                   | 26.05 $\pm$ 0.34                   | 23.06 $\pm$ 0.47                   |
| FedAvg [26]           | 65.07 $\pm$ 0.25                   | 64.41 $\pm$ 0.66                   | 63.19 $\pm$ 0.46                   | 28.48 $\pm$ 0.59                   | 26.06 $\pm$ 0.65                   | 25.58 $\pm$ 0.80                   |
| FedDyn [1]            | 67.31 $\pm$ 0.36                   | 65.02 $\pm$ 0.91                   | 62.49 $\pm$ 0.06                   | 34.17 $\pm$ 0.43                   | 27.06 $\pm$ 0.18                   | 23.88 $\pm$ 0.36                   |
| pFedMe [35]           | 70.60 $\pm$ 0.23                   | 68.92 $\pm$ 0.35                   | 66.40 $\pm$ 0.04                   | 27.97 $\pm$ 0.24                   | 23.82 $\pm$ 0.06                   | 22.35 $\pm$ 0.03                   |
| FedFomo [43]          | 72.33 $\pm$ 0.03                   | 72.17 $\pm$ 0.48                   | 70.86 $\pm$ 0.27                   | 32.15 $\pm$ 0.61                   | 25.90 $\pm$ 1.17                   | 24.48 $\pm$ 0.44                   |
| APFL [6]              | 77.03 $\pm$ 0.26                   | 77.36 $\pm$ 0.18                   | 76.29 $\pm$ 0.13                   | 39.16 $\pm$ 0.93                   | 35.15 $\pm$ 0.65                   | 33.86 $\pm$ 0.60                   |
| FedRep [5]            | 76.85 $\pm$ 0.44                   | 76.03 $\pm$ 0.17                   | 72.30 $\pm$ 0.52                   | 33.43 $\pm$ 0.80                   | 26.86 $\pm$ 0.39                   | 22.76 $\pm$ 0.45                   |
| LG-FedAvg [23]        | 72.83 $\pm$ 0.28                   | 70.44 $\pm$ 0.31                   | 67.55 $\pm$ 0.09                   | 33.65 $\pm$ 0.19                   | 27.13 $\pm$ 0.37                   | 24.82 $\pm$ 0.28                   |
| FedPer [2]            | 77.84 $\pm$ 0.18                   | 77.76 $\pm$ 0.22                   | 75.01 $\pm$ 0.20                   | 35.22 $\pm$ 0.67                   | 28.63 $\pm$ 0.70                   | 25.56 $\pm$ 0.26                   |
| Per-FedAvg [7]        | 75.49 $\pm$ 0.74                   | 76.27 $\pm$ 0.50                   | 75.41 $\pm$ 0.35                   | 32.89 $\pm$ 0.43                   | 32.24 $\pm$ 0.75                   | 32.59 $\pm$ 0.21                   |
| FedRoD [4]            | 79.73 $\pm$ 0.68                   | 79.61 $\pm$ 0.22                   | 77.76 $\pm$ 0.32                   | 39.55 $\pm$ 0.58                   | 33.87 $\pm$ 2.42                   | 31.49 $\pm$ 0.19                   |
| FedBABU [28]          | 78.92 $\pm$ 0.36                   | 79.35 $\pm$ 0.84                   | 76.34 $\pm$ 0.22                   | 32.71 $\pm$ 0.23                   | 29.66 $\pm$ 0.64                   | 27.72 $\pm$ 0.11                   |
| <b>PGFed (ours)</b>   | <u>81.02 <math>\pm</math> 0.41</u> | <u>81.42 <math>\pm</math> 0.31</u> | <u>78.56 <math>\pm</math> 0.35</u> | <u>43.12 <math>\pm</math> 0.03</u> | <u>38.45 <math>\pm</math> 0.44</u> | <u>35.71 <math>\pm</math> 0.54</u> |
| <b>PGFedMo (ours)</b> | <u>81.20 <math>\pm</math> 0.08</u> | <u>81.48 <math>\pm</math> 0.32</u> | <u>78.74 <math>\pm</math> 0.22</u> | <u>43.44 <math>\pm</math> 0.14</u> | <u>38.50 <math>\pm</math> 0.45</u> | <u>35.76 <math>\pm</math> 0.65</u> |

✓ **PGFed and PGFedMo boost the accuracy by up to 15.47%.**

# PGFed – Experiments & results

- Convergence speed

- CIFAR10

|            | 25 clients |              | 50 clients |              | 100 clients |              |
|------------|------------|--------------|------------|--------------|-------------|--------------|
|            | #round     | speedup      | #round     | speedup      | #round      | speedup      |
| APFL       | 31         | 1.0 $\times$ | 28         | 1.7 $\times$ | 24          | 2.6 $\times$ |
| FedPer     | 8          | 3.9 $\times$ | 6          | 7.8 $\times$ | 8           | 7.9 $\times$ |
| Per-FedAvg | 31         | 1.0 $\times$ | 47         | 1.0 $\times$ | 63          | 1.0 $\times$ |
| FedRoD     | 26         | 1.2 $\times$ | 35         | 1.3 $\times$ | 10          | 6.3 $\times$ |
| PGFed      | 9          | 3.4 $\times$ | 14         | 3.4 $\times$ | 15          | 4.2 $\times$ |
| PGFedMo    | 9          | 3.4 $\times$ | 14         | 3.4 $\times$ | 15          | 4.2 $\times$ |

- Mean individual gain over Local

- CIFAR10

|            | 25 clients        | 50 clients        | 100 clients       |
|------------|-------------------|-------------------|-------------------|
| FedAvg     | $-8.99 \pm 10.36$ | $-8.90 \pm 15.48$ | $-5.02 \pm 14.30$ |
| APFL       | $2.79 \pm 8.07$   | $5.73 \pm 8.43$   | $8.37 \pm 6.91$   |
| FedPer     | $5.31 \pm 2.56$   | $8.31 \pm 6.00$   | $8.63 \pm 5.26$   |
| Per-FedAvg | $0.72 \pm 6.22$   | $5.02 \pm 7.39$   | $8.09 \pm 7.00$   |
| FedRoD     | $7.80 \pm 3.68$   | $8.84 \pm 6.29$   | $10.68 \pm 6.14$  |
| PGFed      | $8.49 \pm 4.67$   | $10.78 \pm 5.88$  | $11.15 \pm 5.06$  |
| PGFedMo    | $8.61 \pm 3.59$   | $10.90 \pm 6.11$  | $11.16 \pm 5.44$  |

- CIFAR100

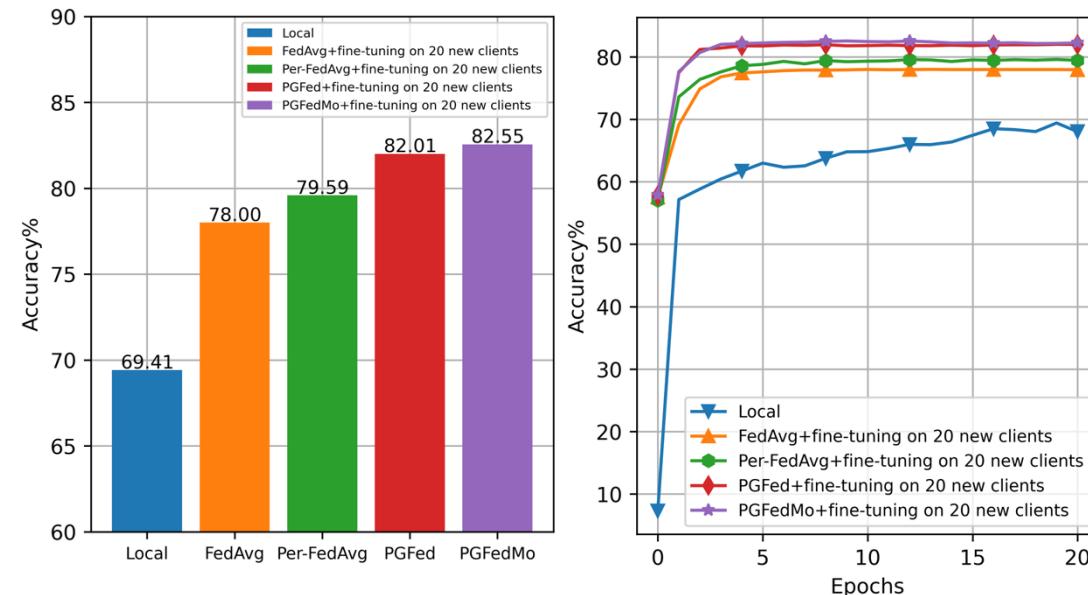
|            | 25 clients       | 50 clients      | 100 clients      |
|------------|------------------|-----------------|------------------|
| FedAvg     | $-3.29 \pm 4.22$ | $0.02 \pm 4.63$ | $1.77 \pm 6.38$  |
| APFL       | $6.48 \pm 2.93$  | $8.70 \pm 3.37$ | $9.31 \pm 4.55$  |
| FedPer     | $3.43 \pm 1.80$  | $2.16 \pm 2.45$ | $2.31 \pm 3.54$  |
| Per-FedAvg | $0.07 \pm 3.71$  | $5.47 \pm 3.86$ | $7.49 \pm 5.73$  |
| FedRoD     | $7.32 \pm 2.68$  | $6.59 \pm 3.17$ | $7.47 \pm 3.69$  |
| PGFed      | $9.34 \pm 1.71$  | $9.01 \pm 2.97$ | $12.05 \pm 3.93$ |
| PGFedMo    | $9.40 \pm 1.87$  | $8.99 \pm 2.76$ | $12.07 \pm 3.97$ |



PGFed and PGFedMo have 3.7 $\times$  average speedup with highest individual gain.

## PGFed – Experiments & results

- Adaptive ability on new clients
  - CIFAR10 & CIFAR100
  - FL on 80 clients, fine-tune global model for 20 epochs on 20 new clients
  - Mean personalized acc. on 20 new clients



**Global models of PGFed and PGFedMo have highest generalizability**

# PGFed – More experiments & results

Details in paper

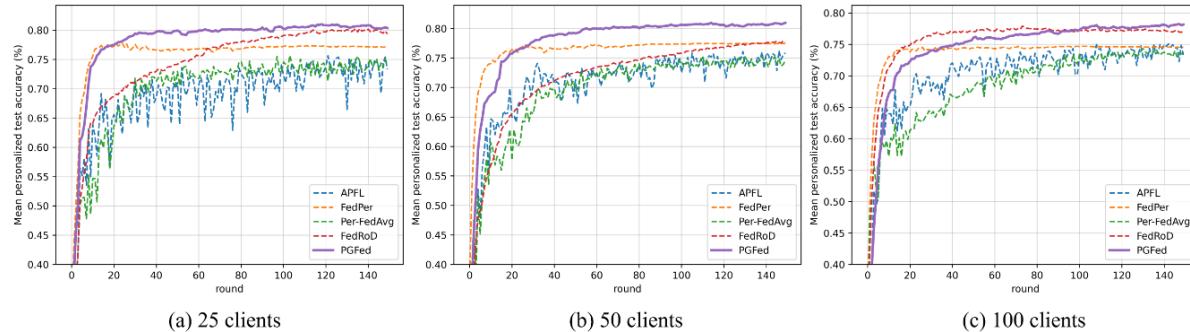


Figure 1. Convergence behavior of the personalized FL approaches with top performance on CIFAR10. While achieving the highest accuracy performance, PGFed is also able to consistently converge faster than several of the baselines that reach high accuracies.

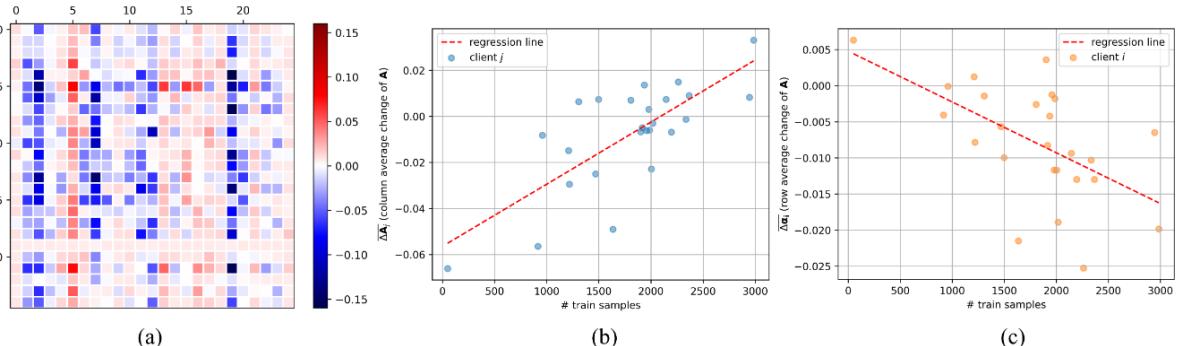


Figure 2. Visualization of the change in  $\mathbf{A}$ . Figure (a) is a heatmap of the change in  $\mathbf{A}$ . For Figure (b) and (c), the Y-axis of Figure (b) represents the column average of the change in  $\mathbf{A}$  (the average change of weights of client  $j$ 's empirical risk on other clients). The Y-axis of Figure (c) is the row average of the change in  $\mathbf{A}$  (the average change of weights of the auxiliary risk on client  $i$ ). Through the regression line, we verify the positive correlation between  $\Delta A_j$  and  $n_j$  in Figure (b), and the negative correlation between  $\Delta \alpha_i$  and  $n_i$  in Figure (c).

|            | Art                                | Clipart                            | Product                            | Real World                         | Mean                               |
|------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| Local      | $17.16 \pm 0.85$                   | $37.65 \pm 0.47$                   | $43.83 \pm 0.40$                   | $24.50 \pm 0.21$                   | $30.79 \pm 0.23$                   |
| FedAvg     | $11.68 \pm 1.26$                   | $41.29 \pm 0.85$                   | $42.49 \pm 1.28$                   | $19.14 \pm 0.89$                   | $28.65 \pm 0.49$                   |
| APFL       | $19.11 \pm 1.55$                   | $44.67 \pm 0.61$                   | <b><math>50.40 \pm 0.56</math></b> | $25.85 \pm 0.88$                   | $35.00 \pm 0.41$                   |
| FedRep     | $20.24 \pm 1.45$                   | $38.43 \pm 1.02$                   | $43.70 \pm 1.04$                   | $24.02 \pm 0.81$                   | $31.60 \pm 0.05$                   |
| LGFedAvg   | $17.54 \pm 0.45$                   | $38.75 \pm 0.13$                   | $44.59 \pm 0.62$                   | $25.79 \pm 0.61$                   | $31.67 \pm 0.21$                   |
| FedPer     | $17.83 \pm 1.07$                   | $38.97 \pm 0.35$                   | $45.87 \pm 0.13$                   | $25.01 \pm 0.52$                   | $31.92 \pm 0.24$                   |
| Per-FedAvg | $14.62 \pm 0.40$                   | $39.94 \pm 1.29$                   | $44.40 \pm 1.32$                   | $21.58 \pm 0.65$                   | $30.13 \pm 0.07$                   |
| FedRoD     | $19.67 \pm 1.23$                   | $42.44 \pm 0.77$                   | $44.34 \pm 2.07$                   | $24.28 \pm 1.69$                   | $32.68 \pm 0.69$                   |
| FedBABU    | $18.18 \pm 3.54$                   | $42.10 \pm 2.31$                   | $43.51 \pm 0.91$                   | <b><math>26.81 \pm 1.86</math></b> | $33.38 \pm 0.29$                   |
| PGFed      | <b><math>22.40 \pm 0.26</math></b> | <b><math>46.48 \pm 1.00</math></b> | <b><math>49.86 \pm 2.14</math></b> | $26.04 \pm 0.80$                   | <b><math>36.19 \pm 0.92</math></b> |
| PGFedMo    | $22.16 \pm 0.45$                   | $45.88 \pm 0.83$                   | $49.45 \pm 0.19$                   | $26.60 \pm 0.99$                   | $36.02 \pm 0.20$                   |

Table 2. Mean and standard deviation over three trials of the mean personalized accuracy% of the four domains (5 clients/domain) and the average performance on Office-home dataset. The highest and second-highest accuracies under each setting are in **bold** and underlined, respectively.

|            | 25 clients<br>sample 50%<br>Dir(1.0) | 50 clients<br>sample 25%<br>Dir(0.3) | 100 clients<br>sample 25%<br>Dir(0.3) |
|------------|--------------------------------------|--------------------------------------|---------------------------------------|
| Local      | $90.45 \pm 0.19$                     | $90.63 \pm 0.07$                     | $87.14 \pm 0.10$                      |
| FedAvg     | $99.11 \pm 0.03$                     | $98.74 \pm 0.04$                     | $98.47 \pm 0.08$                      |
| APFL       | $97.49 \pm 0.05$                     | $97.53 \pm 0.06$                     | $96.19 \pm 0.11$                      |
| FedRep     | $95.06 \pm 0.16$                     | $94.86 \pm 0.07$                     | $92.47 \pm 0.04$                      |
| LGFedAvg   | $90.47 \pm 0.18$                     | $90.99 \pm 0.08$                     | $87.52 \pm 0.22$                      |
| FedPer     | $97.89 \pm 0.06$                     | $97.55 \pm 0.08$                     | $95.56 \pm 0.33$                      |
| Per-FedAvg | $98.40 \pm 0.02$                     | $96.80 \pm 0.04$                     | $95.09 \pm 0.07$                      |
| FedRoD     | $98.61 \pm 0.05$                     | $98.14 \pm 0.09$                     | $97.05 \pm 0.06$                      |
| FedBABU    | $96.49 \pm 0.28$                     | $94.33 \pm 0.13$                     | $91.07 \pm 0.23$                      |
| PGFed      | $99.20 \pm 0.04$                     | <b><math>99.17 \pm 0.05</math></b>   | <b><math>98.94 \pm 0.02</math></b>    |
| PGFedMo    | <b><math>99.21 \pm 0.04</math></b>   | $99.17 \pm 0.07$                     | $98.86 \pm 0.06$                      |

|            | Images/s | Relative speed | Accuracy         |
|------------|----------|----------------|------------------|
| FedAvg     | 6917.1   | 100.00%        | $64.41 \pm 0.66$ |
| APFL       | 3389.8   | 48.99%         | $77.36 \pm 0.18$ |
| Per-FedAvg | 3464.5   | 50.09%         | $76.27 \pm 0.50$ |
| FedRoD     | 6682.4   | 96.61%         | $79.61 \pm 0.22$ |
| PGFed      | 6120.0   | 88.48%         | $81.42 \pm 0.31$ |
| PGFedMo    | 6032.8   | 87.22%         | $81.48 \pm 0.32$ |
| PGFed-CE*  | 6175.5   | 89.28%         | $81.16 \pm 0.56$ |

\* A more communication-efficient variation of PGFed, introduced in Appendix D  
Table 3. Computational speed (in terms of “images/s”) and accuracy on CIFAR10 with 50 clients

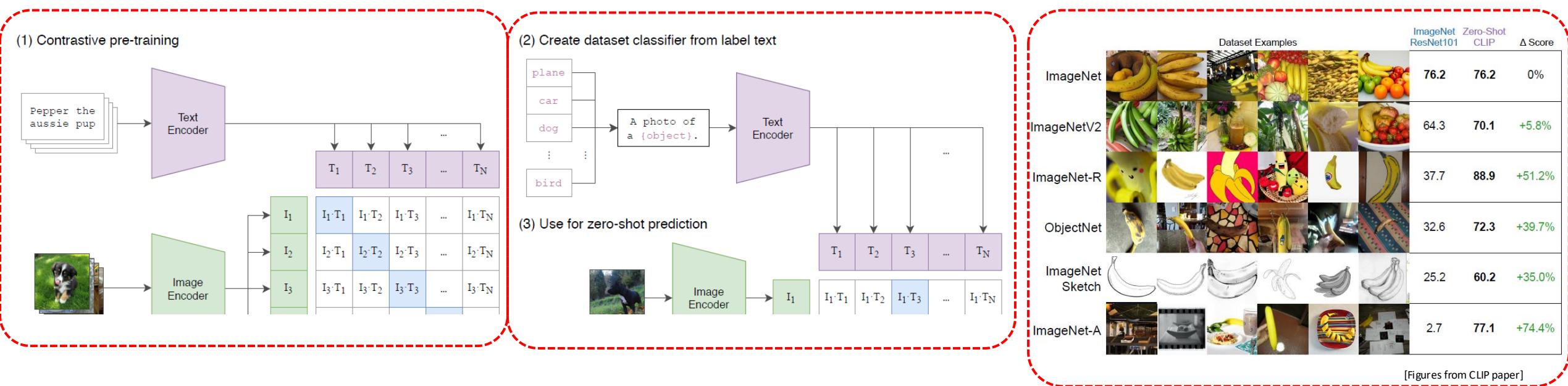
Table 1. Mean and standard deviation over three trials of the mean personalized test accuracy (%) on OrganAMNIST

# Overview

- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (FedSLD)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)
- **Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)**
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- Summary

# pFedMoAP – Background and motivation

- Vision-Language Models (VLMs) like CLIP with their robust representation learning capabilities, show promise for addressing data heterogeneity in federated learning.



# pFedMoAP – Background and motivation

- Vision-Language Models (VLMs) like CLIP with their robust representation learning capabilities, show promise for addressing data heterogeneity in federated learning.
- Traditional fine-tuning of VLMs in federated settings is challenging due to high communication overhead, leading researchers to explore prompt learning as a more efficient adaptation technique.

| Caltech101   | Prompt  | Accuracy     |
|--|---|--------------|
|  | a [CLASS].  | 82.68        |
|  | a photo of [CLASS].   | 80.81        |
|  | a photo of a [CLASS].   | 86.29        |
|  | [V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS]. | <b>91.83</b> |

(a)

| Flowers102  | Prompt  | Accuracy     |
|---|---|--------------|
|  | a photo of a [CLASS].   | 60.86        |
|   | a flower photo of a [CLASS].                                    | 65.81        |
|   | a photo of a [CLASS], a type of flower.                         | 66.14        |
|   | [V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS]. | <b>94.51</b> |

(b)

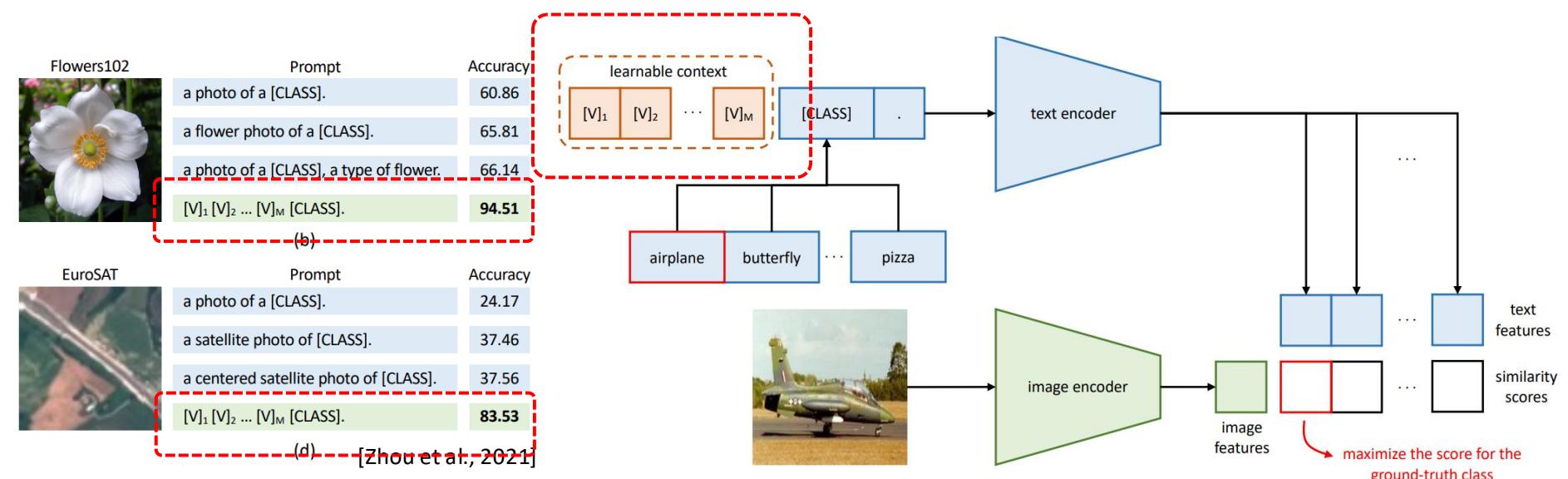
| Describable Textures (DTD)  | Prompt  | Accuracy     |
|---|---|--------------|
|  | a photo of a [CLASS].   | 39.83        |
|   | a photo of a [CLASS] texture.                                   | 40.25        |
|   | [CLASS] texture.  | 42.32        |
|   | [V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS]. | <b>63.58</b> |

(c)

| EuroSAT  | Prompt  | Accuracy     |
|--|---|--------------|
|  | a photo of a [CLASS].   | 24.17        |
|  | a satellite photo of [CLASS].                                   | 37.46        |
|  | a centered satellite photo of [CLASS].                          | 37.56        |
|  | [V] <sub>1</sub> [V] <sub>2</sub> ... [V] <sub>M</sub> [CLASS]. | <b>83.53</b> |

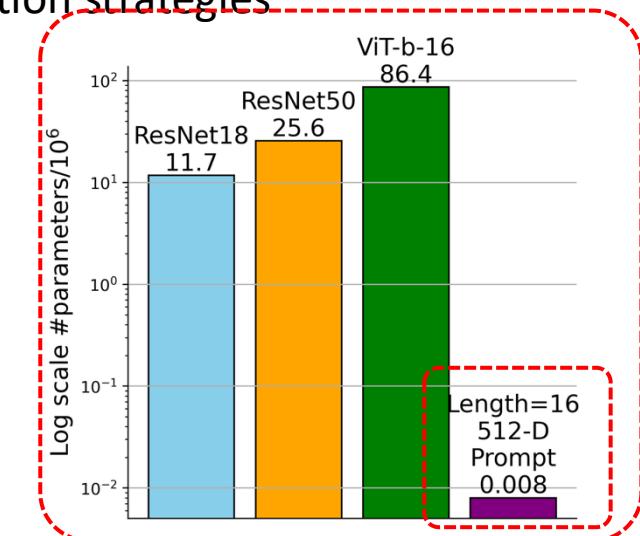
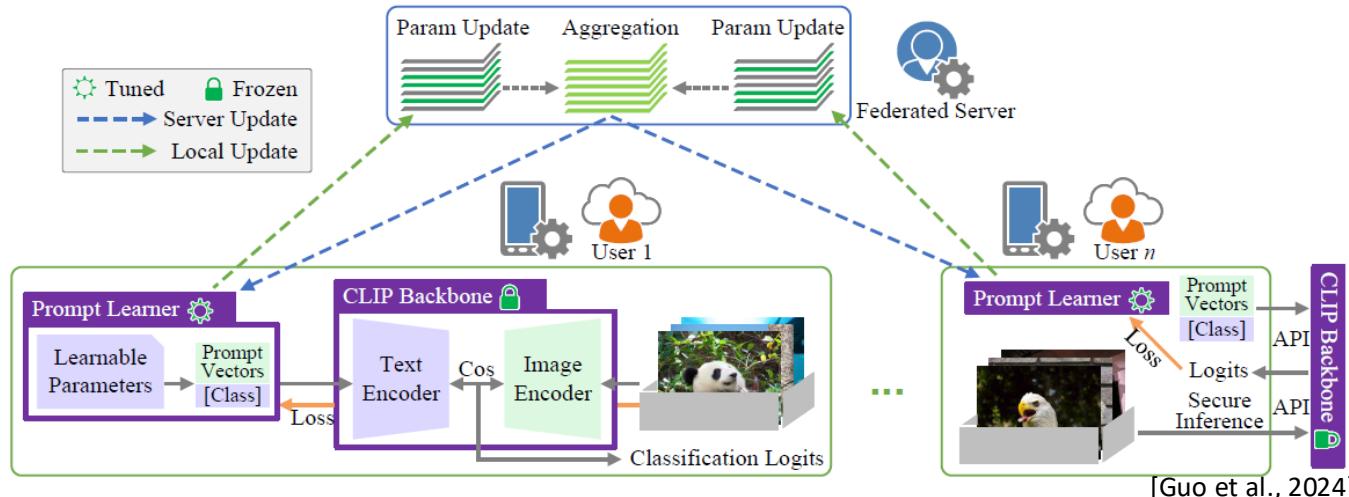
(d) [Zhou et al., 2021]



[Zhou et al., 2021]

# pFedMoAP – Background and motivation

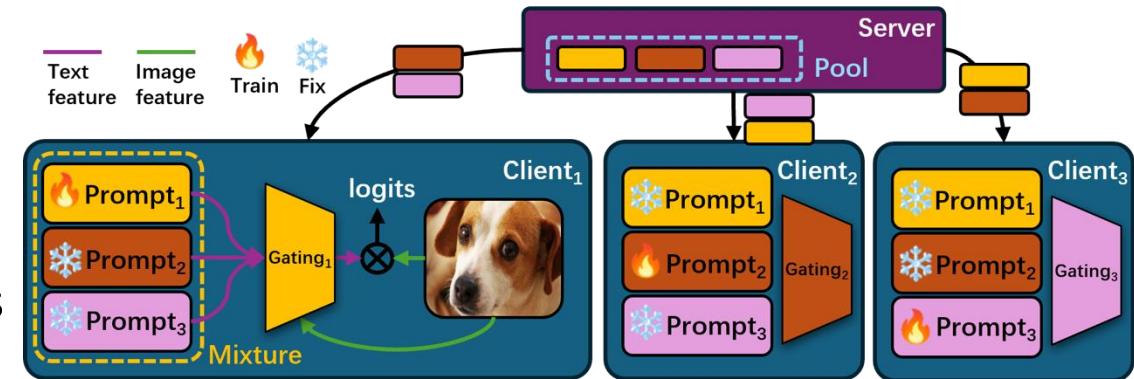
- Vision-Language Models (VLMs) like CLIP with their robust representation learning capabilities, show promise for addressing data heterogeneity in federated learning.
- Traditional fine-tuning of VLMs in federated settings is challenging due to high communication overhead, leading researchers to explore prompt learning as a more efficient adaptation technique.
- Existing federated prompt learning works
  - Habitually fall into traditional FL paradigm where clients are restricted to downloading only a single globally aggregated model – not fully leveraging the prompt's lightweight nature
  - Struggling to handle extreme data heterogeneity, lacking personalization strategies



# pFedMoAP – Background and motivation

**Research question 4:** *How can we devise a personalized federated learning framework, tailored for prompt learning in CLIP-like VLMs, while fully exploiting the lightweight nature of the prompts?*

- Personalized **Federated Mixture of Adaptive Prompts (pFedMoAP)**
  - Allows download of multiple pre-aggregated prompts
  - Uses a Mixture of Experts approach to treat locally updated prompts as specialized experts
  - Implements a client-specific, attention-based gating network to generate enhanced text features



# pFedMoAP – Method

- Formulations for existing paradigms

- Global objective of PFL

$$\min_{\theta_1, \dots, \theta_N} F(\theta_1, \dots, \theta_N) = \min_{\theta_1, \dots, \theta_N} \sum_{i=1}^N p_i F_i(\theta_i)$$

- Prompt learning for CLIP-like VLMs

- Learnable prompt  $P = \{p_1, \dots, p_l\} \in \mathbb{R}^{l \times d}$

- Full prompt  $P^{(c)}$  of class  $c$  is  $P$  with embedding of label  $c$

- Classification

$$\text{logit}^{(c)} = \text{sim}(f(\mathbf{x}), g(\mathbf{P}^{(c)}))$$

$$p(\hat{y} = c | \mathbf{x}) = \frac{\exp(\text{logit}^{(c)}) / \tau}{\sum_{k=1}^C \exp(\text{logit}^{(k)} / \tau)}$$

$\mathbf{x}$ : image

$f(\cdot)$ : CLIP's image encoder

$g(\cdot)$ : CLIP's text encoder

$\tau$ : temperature

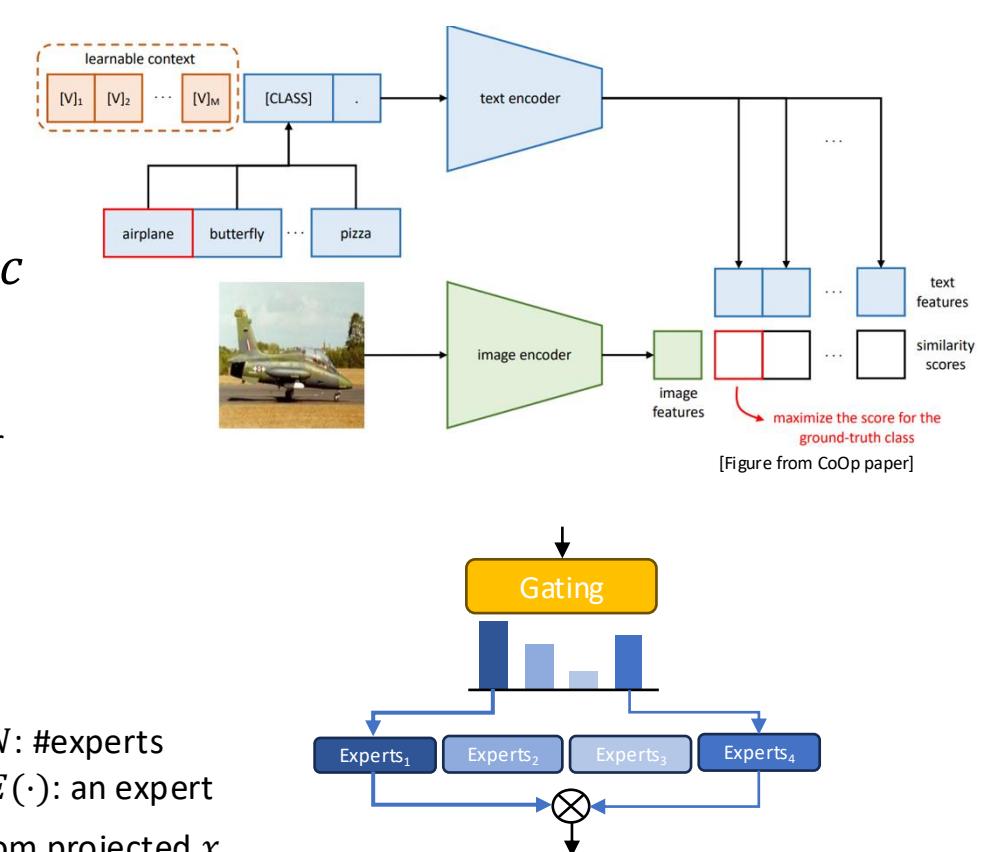
- In FL, aggregated global prompt

$$\mathbf{P}_g^t = \sum_{i \in \mathcal{S}_t} \frac{n_i}{\sum_{k \in \mathcal{S}_t} n_k} \mathbf{P}_i^t$$

- Mixture of Experts (MoE) output

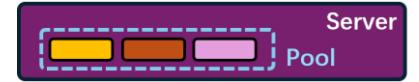
$$\text{MoE}(\mathbf{x}) = \sum_{i=1}^N G(\mathbf{x})_i \cdot E_i(\mathbf{x})$$

$G(\cdot)$ : gating, usually softmax of TopK/N from projected  $\mathbf{x}$



# pFedMoAP – Method

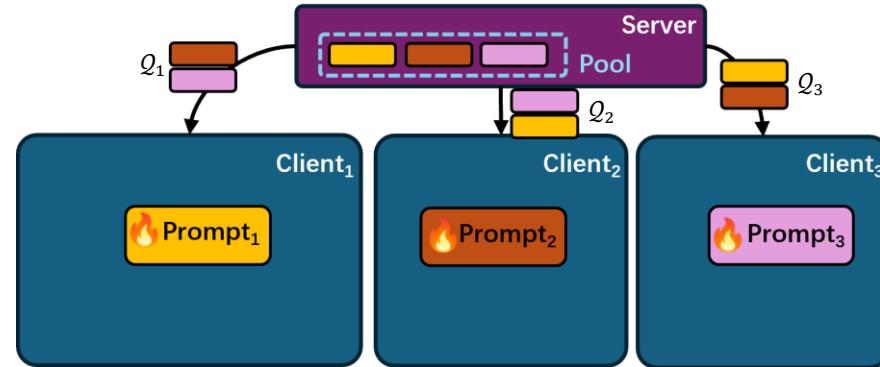
- Workflow
  - Server maintains a pool of prompts  $\mathcal{P}_t = \mathcal{P}_{t-1} - \{P_i^{t-1}\}_{i \in \mathcal{P}_{t-1} \cap \mathcal{S}_t} + \{P_j^t\}_{j \in \mathcal{S}_t}$



# pFedMoAP – Method

- Workflow

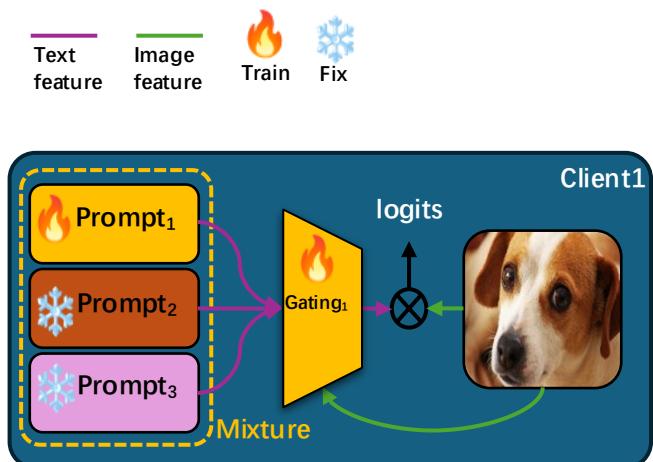
- Server maintains a pool of prompts  $\mathcal{P}_t = \mathcal{P}_{t-1} - \{P_i^{t-1}\}_{i \in \mathcal{P}_{t-1} \cap \mathcal{S}_t} + \{P_j^t\}_{j \in \mathcal{S}_t}$
- Each client  $i \in \mathcal{S}_t$  download  $K$  pre-aggregated (non-local) prompt
  - K-Nearest Neighbors (KNN) since most likely to have similar distribution
  - $\mathcal{Q}_i = \{NL_j\}_{j=1}^K$  : set of clients assigned to client  $i$ , with prompts  $P_{NL_j}$  ( $NL$ = abbr. for non-local)



# pFedMoAP – Method

- Workflow

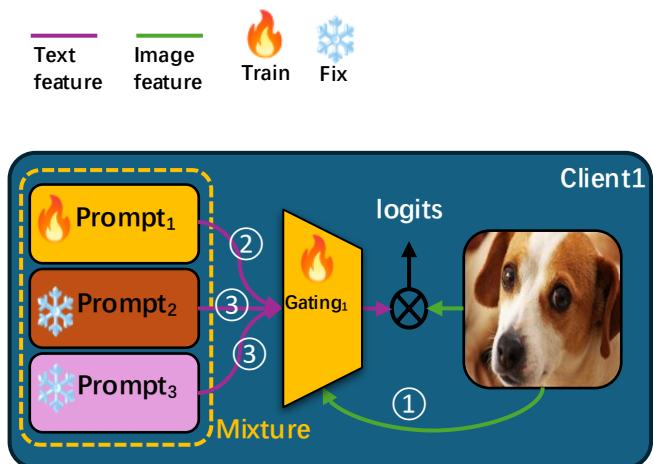
- Server maintains a pool of prompts  $\mathcal{P}_t = \mathcal{P}_{t-1} - \{P_i^{t-1}\}_{i \in \mathcal{P}_{t-1} \cap \mathcal{S}_t} + \{P_j^t\}_{j \in \mathcal{S}_t}$
- Each client  $i \in \mathcal{S}_t$  download  $K$  pre-aggregated (non-local) prompt
  - K-Nearest Neighbors (KNN) since most likely to have similar distribution
  - $\mathcal{Q}_i = \{NL_j\}_{j=1}^K$  : set of clients assigned to client  $i$ , with prompts  $P_{NL_j}$  ( $NL$ = abbr. for non-local)
- Before local training, for once, client compute (fixed) text feature from non-local prompts  $\forall c \in [C], T_{NL}^{(c)} \triangleq \{T_{NL_j}^{(c)} | T_{NL_j}^{(c)} = g(P_{NL_j}^{(c)}), \forall NL_j \in \mathcal{Q}_i\}$



# pFedMoAP – Method

- Workflow

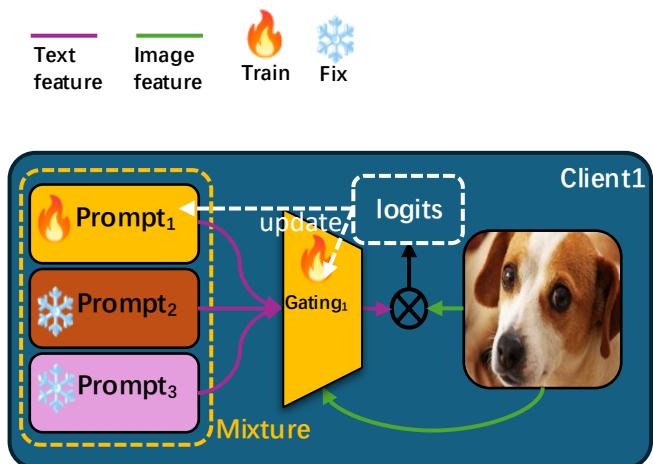
- Server maintains a pool of prompts  $\mathcal{P}_t = \mathcal{P}_{t-1} - \{P_i^{t-1}\}_{i \in \mathcal{P}_{t-1} \cap \mathcal{S}_t} + \{P_j^t\}_{j \in \mathcal{S}_t}$
- Each client  $i \in \mathcal{S}_t$  download  $K$  pre-aggregated (non-local) prompt
  - K-Nearest Neighbors (KNN) since most likely to have similar distribution
  - $\mathcal{Q}_i = \{NL_j\}_{j=1}^K$  : set of clients assigned to client  $i$ , with prompts  $P_{NL_j}$  ( $NL$ = abbr. for non-local)
- Before local training, for once, client compute (fixed) text feature from non-local prompts  $\forall c \in [C], T_{NL}^{(c)} \triangleq \{T_{NL_j}^{(c)} | T_{NL_j}^{(c)} = g(P_{NL_j}^{(c)}), \forall NL_j \in \mathcal{Q}_i\}$
- Gating (detailed in following slides)
  - Input type ①: image feature  $I_k = f(x_k)$
  - Input type ②: text feature from local prompt  $T_L^{(c)} = g(P_i^{(c)})$
  - Input type ③: text features from non-local prompts  $T_{NL}^{(c)}$
  - Output: MoE text feature  $\forall c \in [C], T_{MoE}^{(c)} \triangleq G(I_k, T_L^{(c)}, T_{NL}^{(c)} | \theta_i)$



# pFedMoAP – Method

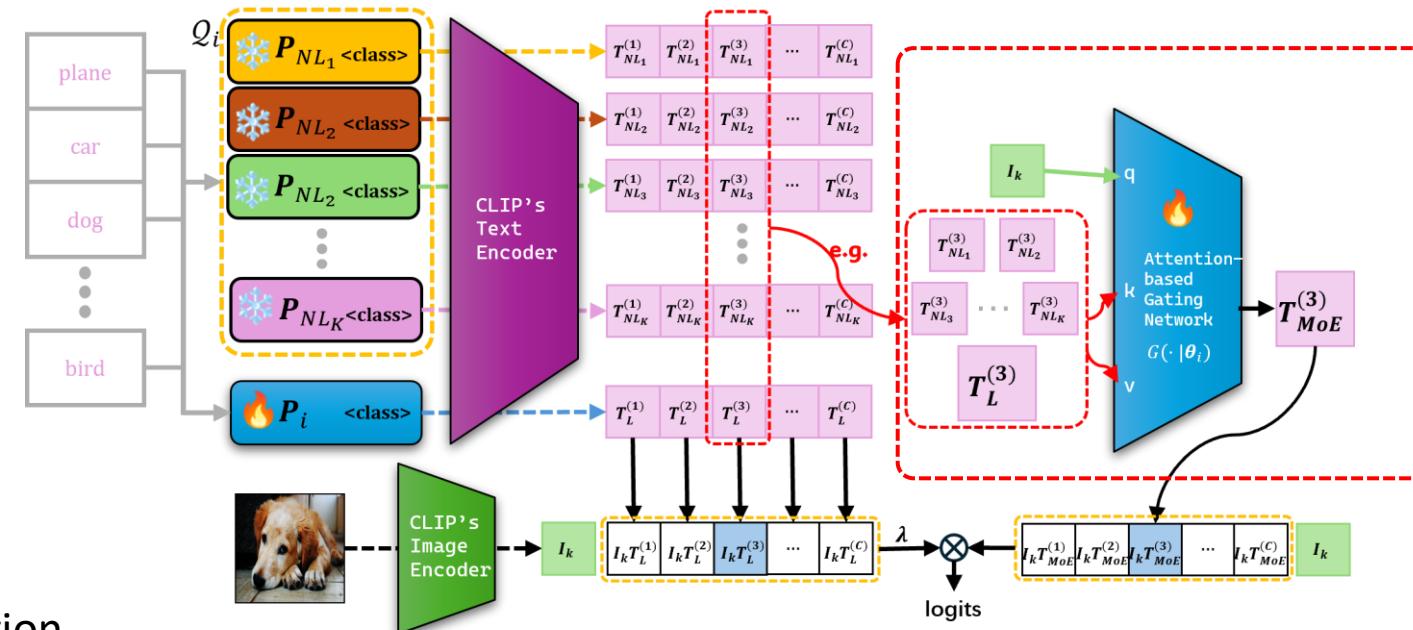
- Workflow

- Server maintains a pool of prompts  $\mathcal{P}_t = \mathcal{P}_{t-1} - \{P_i^{t-1}\}_{i \in \mathcal{P}_{t-1} \cap \mathcal{S}_t} + \{P_j^t\}_{j \in \mathcal{S}_t}$
- Each client  $i \in \mathcal{S}_t$  download  $K$  pre-aggregated (non-local) prompt
  - K-Nearest Neighbors (KNN) since most likely to have similar distribution
  - $\mathcal{Q}_i = \{NL_j\}_{j=1}^K$  : set of clients assigned to client  $i$ , with prompts  $P_{NL_j}$  ( $NL$ = abbr. for non-local)
- Before local training, for once, client compute (fixed) text feature from non-local prompts  $\forall c \in [C], T_{NL}^{(c)} \triangleq \{T_{NL_j}^{(c)} | T_{NL_j}^{(c)} = g(P_{NL_j}^{(c)}), \forall NL_j \in \mathcal{Q}_i\}$
- Gating (detailed in following slides)
  - Input type ①: image feature  $I_k = f(x_k)$
  - Input type ②: text feature from local prompt  $T_L^{(c)} = g(P_i^{(c)})$
  - Input type ③: text features from non-local prompts  $T_{NL}^{(c)}$
  - Output: MoE text feature  $\forall c \in [C], T_{MoE}^{(c)} \triangleq G(I_k, T_L^{(c)}, T_{NL}^{(c)} | \theta_i)$
- Final step: compute logits, manually address local prompt since it is the only locally learnable prompt  $\forall c \in [C], \text{logit}^{(c)} = \text{sim}(I_k, T_{MoE}^{(c)}) + \lambda \cdot \text{sim}(I_k, T_L^{(c)})$



## pFedMoAP – Method

- Attention-based gating network: mechanism



- Multi-head attention
- Pooling on features to reduce the size of gating from 1024 to 128
- $Q = \text{Pooling}(I_k)$ ,  $K = V = \text{Pooling}\{T_L^{(c)}, T_{NL_1}^{(c)}, T_{NL_2}^{(c)}, \dots, T_{NL_K}^{(c)}\}$
- MoE text feature:  $T_{MoE}^{(c)} = G(I_k, T_L^{(c)}, T_{NL}^{(c)} | \theta_i) = \text{MHA}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$     $\text{head}_q = \text{Attention}(QW_q^Q, KW_q^K, VW_q^V)$

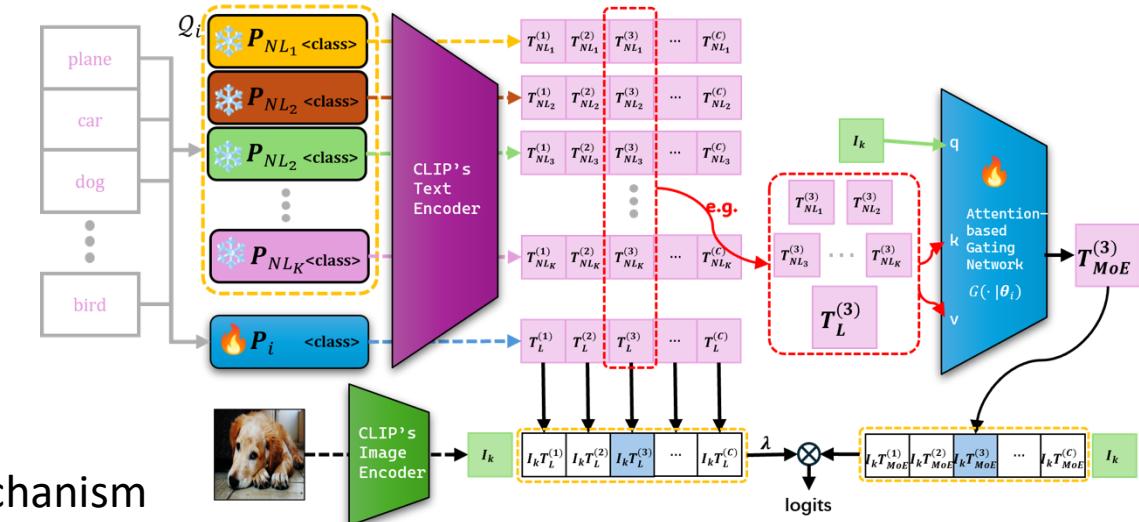
## pFedMoAP – Method

- Attention-based gating network: design rationale against traditional projection-based gating network

- Projection-based gating network  $G_{\text{proj}}(\mathbf{x}_k) \in \mathbb{R}^{K+1}$

$$MoE(\mathbf{x}) = \sum_{i=1}^N G(\mathbf{x})_i \cdot E_i(\mathbf{x})$$

- Attention-based gating against projection-based gating
  - is more robust to adaptive experts
  - serves as linear probing with more capacity
  - leverages CLIP's feature alignment with attention mechanism
  - is agnostic to experts' order



# pFedMoAP – Experiments & results

- Datasets

| Dataset          | Training Set Size | Test Set Size | Number of Classes | Number of Clients | Sample Rate | Data Heterogeneity   |
|------------------|-------------------|---------------|-------------------|-------------------|-------------|----------------------|
| Flowers102       | 4,093             | 2,463         | 102               | 10                | 100%        | Pathological non-IID |
| OxfordPets       | 2,944             | 3,669         | 37                | 10                | 100%        | Pathological non-IID |
| Food101          | 50,500            | 30,300        | 101               | 10                | 100%        | Pathological non-IID |
| Caltech101       | 4,128             | 2,465         | 100               | 10                | 100%        | Pathological non-IID |
| DTD              | 2,820             | 1,692         | 47                | 10                | 100%        | Pathological non-IID |
| Office-Caltech10 | 2,025             | 508           | 10                | 20                | 50%         | Dir(0.3)             |
| DomainNet        | 18,278            | 4,573         | 10                | 30                | 25%         | Dir(0.3)             |
| CIFAR10          | 50,000            | 10,000        | 10                | 100               | 10%         | Dir(0.5)             |
| CIFAR100         | 50,000            | 10,000        | 100               | 100               | 10%         | Dir(0.5)             |


 CLIP datasets, pathological label shift  

 Domain adaptation datasets, feature + label shift  

 CIFAR 10/100, Practical label shift

# pFedMoAP – Experiments & results

- Datasets

| Dataset          | Training Set Size | Test Set Size | Number of Classes | Number of Clients | Sample Rate | Data Heterogeneity   |
|------------------|-------------------|---------------|-------------------|-------------------|-------------|----------------------|
| Flowers102       | 4,093             | 2,463         | 102               | 10                | 100%        | Pathological non-IID |
| OxfordPets       | 2,944             | 3,669         | 37                | 10                | 100%        | Pathological non-IID |
| Food101          | 50,500            | 30,300        | 101               | 10                | 100%        | Pathological non-IID |
| Caltech101       | 4,128             | 2,465         | 100               | 10                | 100%        | Pathological non-IID |
| DTD              | 2,820             | 1,692         | 47                | 10                | 100%        | Pathological non-IID |
| Office-Caltech10 | 2,025             | 508           | 10                | 20                | 50%         | Dir(0.3)             |
| DomainNet        | 18,278            | 4,573         | 10                | 30                | 25%         | Dir(0.3)             |
| CIFAR10          | 50,000            | 10,000        | 10                | 100               | 10%         | Dir(0.5)             |
| CIFAR100         | 50,000            | 10,000        | 100               | 100               | 10%         | Dir(0.5)             |

- Compared methods
  - Local methods
    - Zero-shot CLIP
    - CoOp (prompt learning)
  - Federated prompt learning + FL/PFL
    - PromptFL
    - PromptFL + FedProx
    - PromptFL + FT
    - PromptFL + FedAMP
    - PromptFL + FedPer
  - Personalization designed for federated prompt learning
    - pFedPrompt
    - FedOTP

# pFedMoAP – Experiments & results

- Pathological label shift on CLIP datasets

|                            | Flowers102                         | OxfordPets                         | Food101                            | Caltech101                         | DTD                                |
|----------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| ZS-CLIP [71]               | $62.17 \pm 0.12$                   | $84.47 \pm 0.01$                   | $75.27 \pm 0.05$                   | $85.14 \pm 0.24$                   | $40.21 \pm 0.12$                   |
| CoOp [100]                 | $70.14 \pm 0.76$                   | $83.21 \pm 1.30$                   | $70.43 \pm 2.42$                   | $87.37 \pm 0.44$                   | $44.23 \pm 0.63$                   |
| PromptFL [31]              | $72.80 \pm 1.14$                   | $90.79 \pm 0.61$                   | $77.31 \pm 1.64$                   | $89.70 \pm 1.99$                   | $54.11 \pm 0.22$                   |
| PromptFL+FT [12]           | $72.31 \pm 0.91$                   | $91.23 \pm 0.50$                   | $77.16 \pm 1.56$                   | $89.70 \pm 0.25$                   | $53.74 \pm 1.36$                   |
| PromptFL+FedPer [5]        | $72.11 \pm 1.35$                   | $89.50 \pm 1.62$                   | $71.29 \pm 1.87$                   | $86.72 \pm 1.45$                   | $50.23 \pm 0.82$                   |
| PromptFL+FedProx [50]      | $66.40 \pm 0.29$                   | $89.24 \pm 0.41$                   | $76.24 \pm 1.94$                   | $89.41 \pm 0.55$                   | $44.26 \pm 1.11$                   |
| PromptFL+FedAMP [37]       | $69.10 \pm 0.13$                   | $80.21 \pm 0.44$                   | $74.48 \pm 1.71$                   | $87.31 \pm 1.60$                   | $47.16 \pm 0.92$                   |
| pFedPrompt [30]            | $86.46 \pm 0.15$                   | $91.84 \pm 0.41$                   | $92.26 \pm 1.34$                   | $96.54 \pm 1.31$                   | $77.14 \pm 0.09$                   |
| FedOTP [48]                | $96.23 \pm 0.44$                   | $98.82 \pm 0.11$                   | $92.73 \pm 0.15$                   | $97.02 \pm 0.36$                   | $87.64 \pm 0.70$                   |
| pFedMoAP ( $\lambda=0.0$ ) | $97.61 \pm 0.11$                   | $94.83 \pm 0.65$                   | $86.71 \pm 0.15$                   | $95.71 \pm 0.37$                   | $85.64 \pm 0.34$                   |
| pFedMoAP ( $\lambda=0.5$ ) | <b><math>98.41 \pm 0.04</math></b> | <b><math>99.06 \pm 0.09</math></b> | <b><math>93.39 \pm 0.09</math></b> | <b><math>97.95 \pm 0.07</math></b> | <b><math>89.13 \pm 0.54</math></b> |

# pFedMoAP – Experiments & results

- Practical label shift on CIFAR datasets
  - $\text{Dir}(\alpha = 0.5)$ , 100 clients, 10% sample rate, 120 rounds
  - CLIP backbone: ResNet50
- Feature + label shift on domain adaptation datasets
  - 5 clients/domain,  $\text{Dir}(\alpha = 0.3)$
  - DomainNet = 30 clients, 25% sample rate, 25 rounds
  - Office-Caltech10 = 20 clients, 50% sample rate
  - CLIP backbone: ViT-b-16

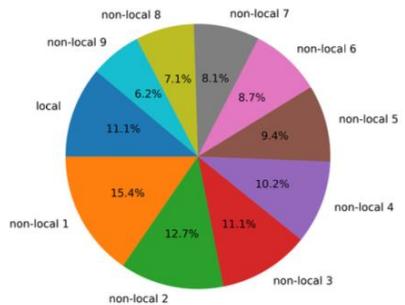
|                     | CIFAR10          | CIFAR10          |
|---------------------|------------------|------------------|
| ZS CLIP [79]        | $53.46 \pm 0.21$ | $32.68 \pm 0.00$ |
| CoOp [115]          | $80.84 \pm 0.39$ | $48.74 \pm 0.17$ |
| PromptFL [36]       | $73.29 \pm 0.37$ | $45.00 \pm 0.62$ |
| Prompt+FedProx [57] | $73.32 \pm 0.34$ | $45.63 \pm 0.75$ |
| pFedMoAP            | $83.46 \pm 0.53$ | $53.42 \pm 0.22$ |

| DomainNet      | Clipart           | Infograph         | Painting         | Quickdraw         | Real              | Sketch           | Average           |
|----------------|-------------------|-------------------|------------------|-------------------|-------------------|------------------|-------------------|
| ZS CLIP        | $9.18 \pm 0.62$   | $10.03 \pm 0.16$  | $9.93 \pm 0.51$  | $10.25 \pm 0.40$  | $9.90 \pm 1.30$   | $9.54 \pm 1.13$  | $9.81 \pm 0.30$   |
| CoOp           | $43.84 \pm 3.51$  | $45.72 \pm 0.85$  | $29.94 \pm 0.46$ | $36.83 \pm 1.17$  | $31.64 \pm 0.49$  | $33.97 \pm 0.78$ | $36.99 \pm 0.79$  |
| PromptFL       | $27.63 \pm 16.41$ | $27.69 \pm 18.07$ | $21.62 \pm 8.34$ | $23.45 \pm 13.49$ | $20.62 \pm 11.03$ | $25.90 \pm 8.10$ | $24.48 \pm 12.52$ |
| Prompt+FedProx | $22.23 \pm 15.42$ | $21.75 \pm 17.00$ | $18.58 \pm 8.15$ | $19.40 \pm 12.59$ | $17.17 \pm 10.25$ | $22.49 \pm 8.44$ | $20.27 \pm 11.83$ |
| pFedMoAP       | $47.49 \pm 0.64$  | $46.73 \pm 0.71$  | $32.74 \pm 0.84$ | $37.16 \pm 0.34$  | $31.02 \pm 0.59$  | $37.67 \pm 0.72$ | $38.80 \pm 0.11$  |

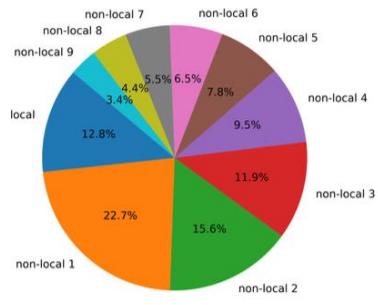
| Office-Caltech10    | Amazon           | Caltech           | DSLR             | Webcam           | Average          |
|---------------------|------------------|-------------------|------------------|------------------|------------------|
| ZS-CLIP [79]        | $9.83 \pm 1.63$  | $10.67 \pm 0.89$  | $10.89 \pm 1.40$ | $6.20 \pm 3.84$  | $9.40 \pm 0.77$  |
| CoOp [115]          | $30.29 \pm 3.64$ | $35.88 \pm 1.30$  | $29.89 \pm 5.15$ | $33.43 \pm 2.25$ | $32.37 \pm 1.81$ |
| PromptFL [36]       | $21.08 \pm 9.60$ | $23.72 \pm 12.21$ | $22.94 \pm 7.96$ | $25.88 \pm 7.72$ | $23.41 \pm 9.06$ |
| Prompt+FedProx [57] | $18.64 \pm 8.58$ | $19.56 \pm 11.59$ | $20.89 \pm 7.38$ | $22.96 \pm 7.56$ | $20.51 \pm 8.48$ |
| pFedMoAP            | $35.47 \pm 1.37$ | $37.45 \pm 1.33$  | $45.11 \pm 3.14$ | $35.22 \pm 1.04$ | $38.31 \pm 1.21$ |

# pFedMoAP – Experiments & results

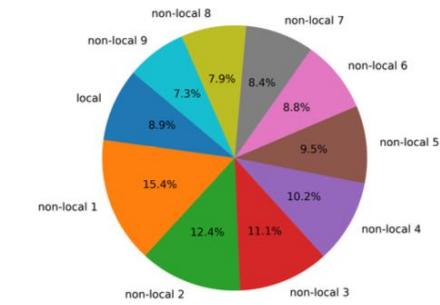
- Contributions of experts



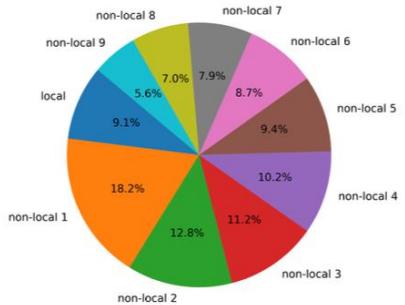
(a) Caltech101, 10 experts



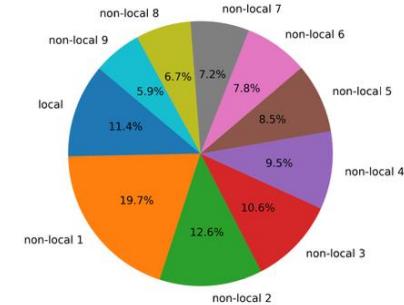
(b) DTD, 10 experts



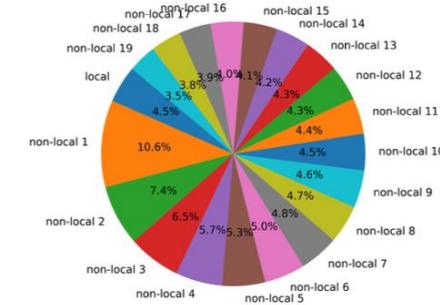
(c) Food101, 10 experts



(d) Flowers102, 10 experts



(e) OxfordPets, 10 experts



(f) DomainNet, 20 experts

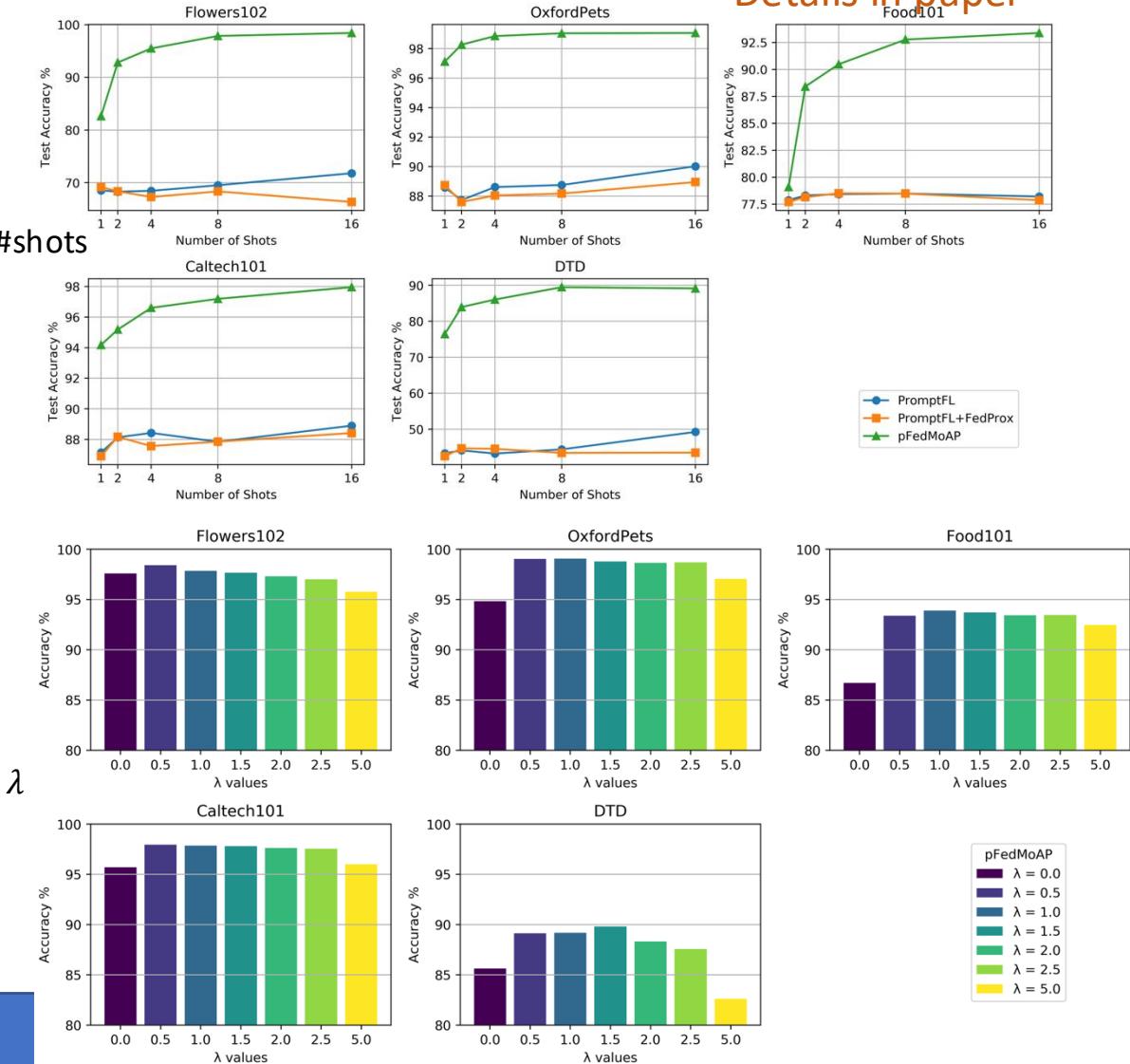
## pFedMoAP – Experiments & results

- Attention-based vs. linear projection-base gating network

|   | Flowers102       | OxfordPets       | Food101          | Caltech101       | DTD              |
|---|------------------|------------------|------------------|------------------|------------------|
| Linear projection-based (3 experts)         | 86.92 $\pm$ 1.84 | 90.54 $\pm$ 1.33 | 78.19 $\pm$ 3.07 | 89.59 $\pm$ 1.46 | 61.42 $\pm$ 5.43 |
| Linear projection-based (10 experts)        | 69.64 $\pm$ 4.57 | 52.78 $\pm$ 6.88 | 77.39 $\pm$ 3.29 | 86.57 $\pm$ 1.96 | 30.42 $\pm$ 7.14 |
| Attention-based, with aggregation           | 97.56 $\pm$ 0.07 | 98.24 $\pm$ 0.12 | 91.89 $\pm$ 0.19 | 96.17 $\pm$ 0.18 | 87.52 $\pm$ 0.69 |
| Attention-based, without aggregation (ours) | 98.41 $\pm$ 0.04 | 99.06 $\pm$ 0.09 | 93.39 $\pm$ 0.09 | 97.95 $\pm$ 0.07 | 89.13 $\pm$ 0.54 |

# pFedMoAP – Ablation studies

Details in paper



## Differential privacy

|   | Flowers102          | OxfordPets | Food101    | Caltech101 | DTD        |            |
|---|---------------------|------------|------------|------------|------------|------------|
| <b>Without differential privacy (from Tab. 11)</b>            |                     |            |            |            |            |            |
| PromptFL [36]   | 72.80±1.14          | 90.79±0.61 | 77.31±1.64 | 89.70±1.99 | 54.11±0.22 |            |
| PromptFL+FedProx [57]   | 66.40±0.29          | 89.24±0.41 | 76.24±1.94 | 89.41±0.55 | 44.26±1.11 |            |
| pFedMoAP(ours)  | 98.41±0.04          | 99.06±0.09 | 93.39±0.09 | 97.95±0.07 | 89.13±0.54 |            |
| <b>With differential privacy (<math>\epsilon = 50</math>)</b> |                     |            |            |            |            |            |
| PromptFL [36]   | 67.07±0.60          | 88.05±0.32 | 77.41±0.60 | 84.83±0.42 | 38.39±1.25 |            |
| PromptFL+FedProx [57]   | 66.22±0.63          | 87.78±0.61 | 77.27±0.59 | 84.68±0.64 | 39.43±1.11 |            |
| pFedMoAP(ours)  | 98.34±0.06          | 99.08±0.02 | 93.36±0.04 | 97.90±0.08 | 89.99±0.49 |            |
| <b>With differential privacy (<math>\epsilon = 25</math>)</b> |                     |            |            |            |            |            |
| PromptFL [36]   | 64.25±1.10          | 86.26±1.07 | 76.84±0.66 | 85.00±1.59 | 38.19±0.66 |            |
| PromptFL+FedProx [57]   | 62.87±0.99          | 86.82±0.47 | 76.21±0.64 | 84.51±1.52 | 37.82±0.52 |            |
| pFedMoAP(ours)  | 98.36±0.12          | 99.02±0.04 | 93.41±0.13 | 97.99±0.06 | 89.11±0.28 |            |
| <b>Feature dimension</b>                                      |                     |            |            |            |            |            |
|   | Gating network size | Flowers102 | OxfordPets | Food101    | Caltech101 | DTD        |
| $d_{\text{feature}} = 32$                                     | 4.2K                | 97.28±0.18 | 98.75±0.32 | 93.42±0.08 | 97.37±0.08 | 88.61±0.89 |
| $d_{\text{feature}} = 64$                                     | 16.6K               | 98.55±0.10 | 98.91±0.23 | 93.89±0.12 | 97.75±0.12 | 89.96±0.09 |
| $d_{\text{feature}} = 128$                                    | 66.0K               | 98.41±0.04 | 99.06±0.09 | 93.39±0.09 | 97.95±0.07 | 89.13±0.54 |
| $d_{\text{feature}} = 256$                                    | 263.2K              | 99.01±0.05 | 98.88±0.21 | 92.49±0.20 | 97.93±0.07 | 90.88±0.16 |
| $d_{\text{feature}} = 512$                                    | 1.1M                | 98.18±0.38 | 96.85±0.22 | 90.34±0.31 | 96.99±0.11 | 89.65±0.10 |
| $d_{\text{feature}} = 1024$                                   | 4.2M                | 98.11±0.33 | 95.81±0.84 | 89.20±0.37 | 96.82±0.26 | 89.03±0.14 |

# Overview

- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (FedSLD)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)

RSNA 2025 • **Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection**

- Summary

## Case study: PFL for real-world breast cancer detection

- Challenges in FL for Medical Imaging
  - Limited sample size
    - Due to high costs of medical imaging and labeling
    - Leads to more severely inconsistent local objectives

## Case study: PFL for real-world breast cancer detection

- Challenges in FL for Medical Imaging
  - Limited sample size
  - Data distribution bias
    - Local demographics can hardly represent large population
    - Global distribution often presents extreme imbalance

## Case study: PFL for real-world breast cancer detection

- Challenges in FL for Medical Imaging
  - Limited sample size
  - Data distribution bias
  - Both feature and label shift
    - Disease prevalence Geographical regions and demographics
    - Institutional specialization
    - Equipment variations
    - Clinical protocol differences
    - Healthcare access across demographics

## Case study: PFL for real-world breast cancer detection

- Challenges in FL for Medical Imaging
  - Limited sample size
  - Data distribution bias
  - Both feature and label shift
  - **Uncertainty in labels**
    - Subjective nature of medical image interpretation causes challenges in terms of label quality and consistency

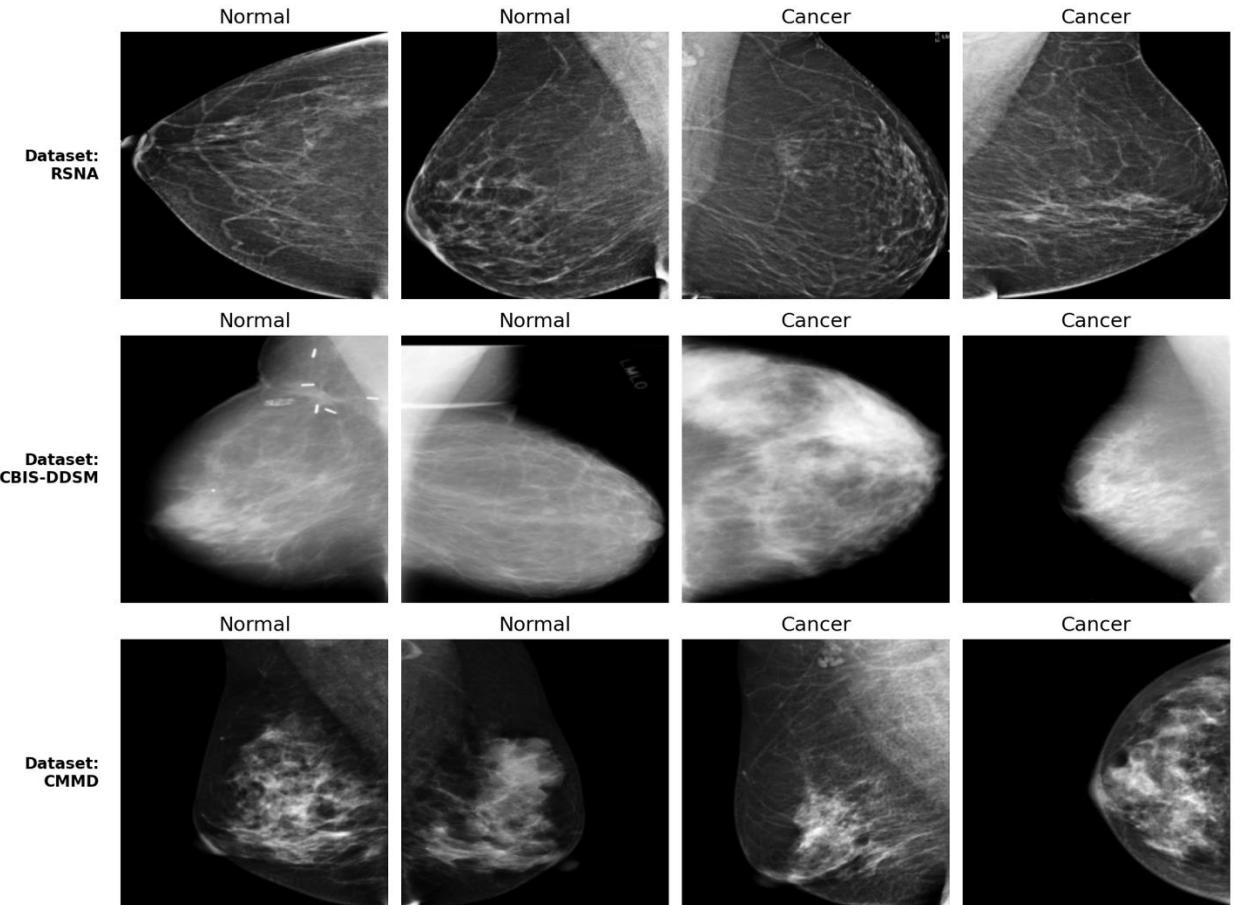
## Case study: PFL for real-world breast cancer detection

- Challenges in FL for Medical Imaging
  - Limited sample size
  - Data distribution bias
  - Both feature and label shift
  - Uncertainty in labels

**Research question 5:** *How can we carefully implement, train, and evaluate existing FL/PFL algorithms and potentially design novel federated frameworks to address the challenges in medical imaging applications?*

# Case study: PFL for real-world breast cancer detection

- Datasets
  - RSNA breast cancer detection: 54K mammograms with metadata
  - CBIS-DDSM: 3K annotated mammograms with pixel-level lesion masks.
  - CMMMD: 5K studies from two Chinese hospitals for cross-domain evaluation.



# Case study: PFL for real-world breast cancer detection

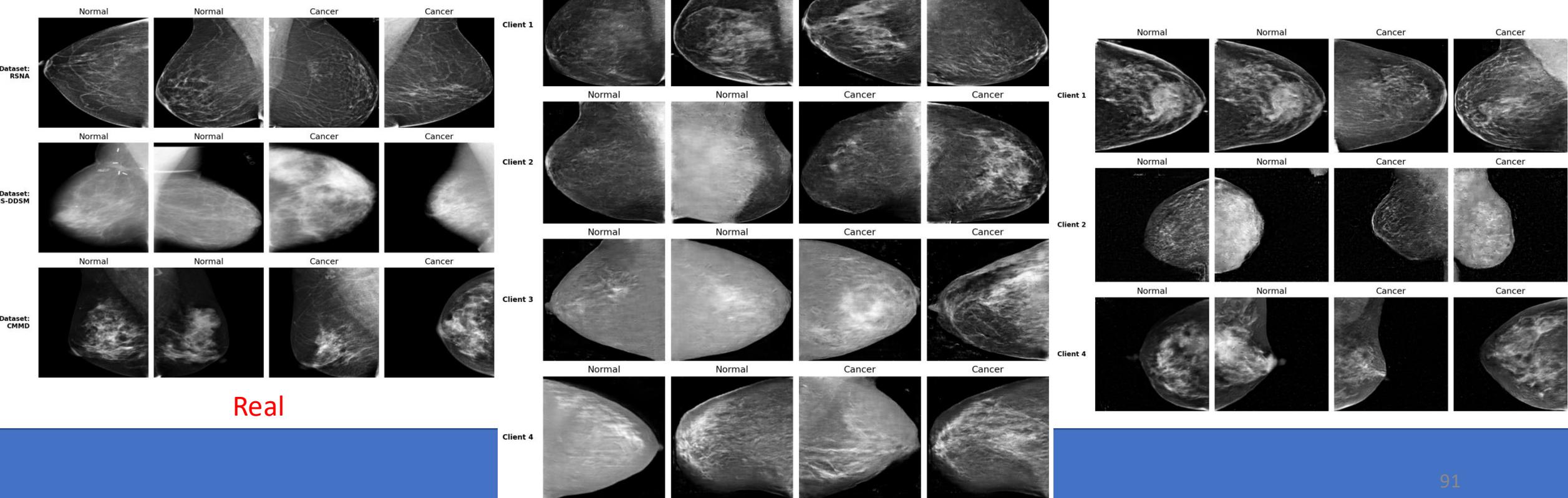
- Data partition for FL

- **Setting 1:**
  - RSNA only
  - Partitioned by machine ID
- **Setting 2:**
  - RSNA, CBIS-DDSM, CMMMD
  - Partitioned by datasets
- 85% training, 15% validation

| Setting   | Client   | Dataset Source  | Cancer Images | Normal Images |
|-----------|----------|-----------------|---------------|---------------|
| Setting 1 | Client 1 | RSNA Machine 49 | 628           | 2,512         |
|           | Client 2 | RSNA Machine 48 | 187           | 1,122         |
|           | Client 3 | RSNA Machine 29 | 184           | 1,104         |
|           | Client 4 | RSNA Machine 21 | 159           | 1,272         |
| Setting 2 | Client 1 | RSNA Site A     | 664           | 2,656         |
|           | Client 2 | RSNA Site B     | 494           | 2,964         |
|           | Client 3 | CBIS-DDSM       | 1,350         | 1,753         |
|           | Client 4 | CMMMD           | 4,094         | 1,108         |

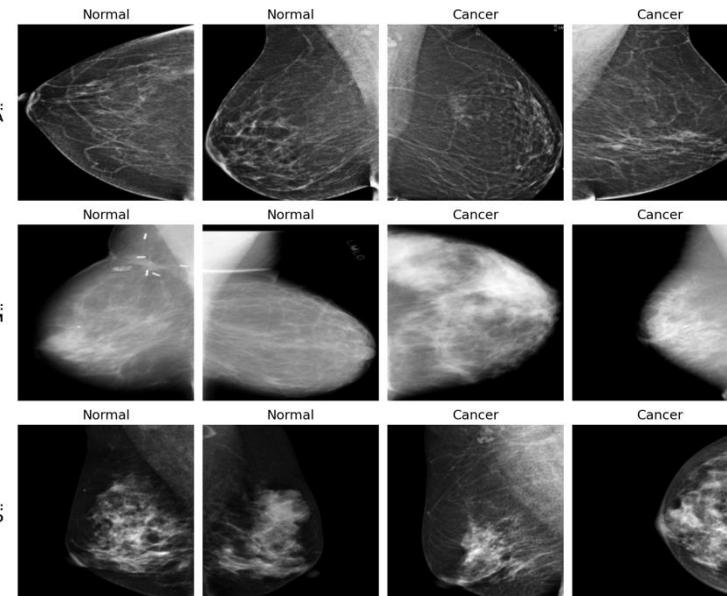
# Case study: PFL for real-world breast cancer detection

- Methods
  - Local
  - FedAvg
  - APPLE
  - PGFed
- Generative data augmentation
  - Mitigating imbalance
  - Mitigating heterogeneity

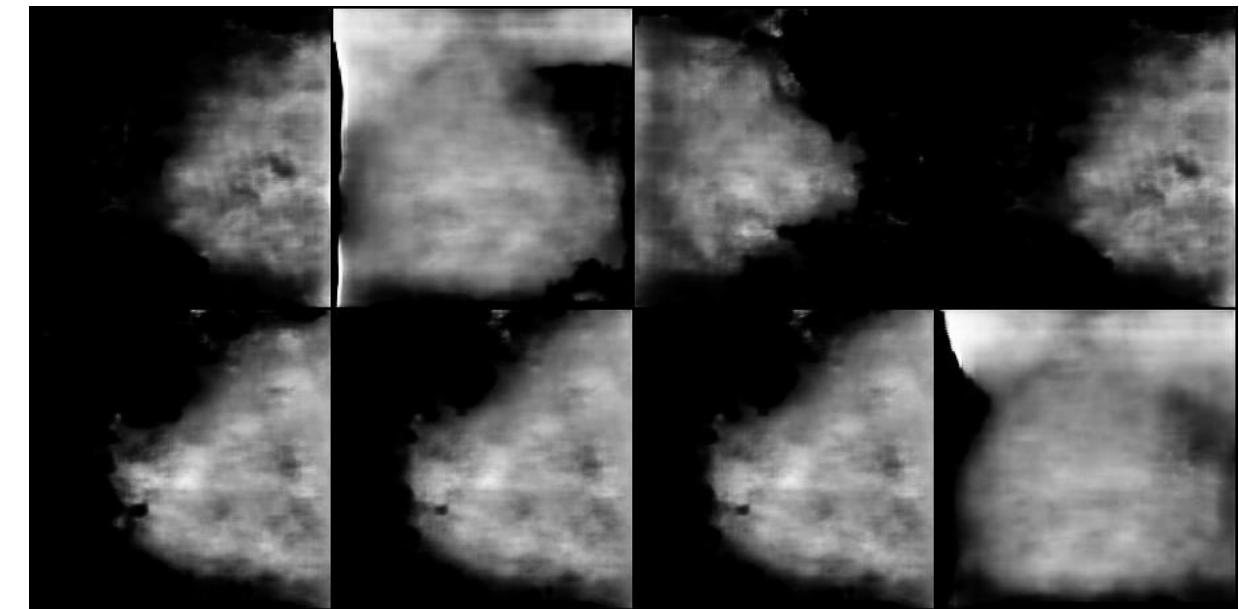


# Case study: PFL for real-world breast cancer detection

- Methods
  - Local
  - FedAvg
  - APPLE
  - PGFed
- Generative data augmentation
  - Mitigating imbalance
  - Mitigating heterogeneity



Real



Failed cases: Setting 2 client 2 synthesized images

# Case study: PFL for real-world breast cancer detection

- Results

| Algorithm             | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std |
|-----------------------|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|
| <i>Real data only</i> |          |           |                     |                     |                     |                     |                |          |           |                     |                     |                     |                     |                |
| Local                 | 66.59    | 82.29     | 64.85 (0.00)        | 73.37 (0.00)        | 60.93 (0.00)        | 67.21 (0.00)        | 0.00 ± 0.00    | 62.37    | 73.86     | 60.96 (0.00)        | 62.47 (0.00)        | 64.74 (0.00)        | 61.32 (0.00)        | 0.00 ± 0.00    |
| FedAvg                | 63.09    | 84.26     | 74.24 (9.39)        | 60.90 (-12.47)      | 57.29 (-3.64)       | 66.99 (-0.22)       | -1.73 ± 7.82   | 64.48    | 75.43     | 69.55 (8.59)        | 65.50 (3.04)        | 64.49 (-0.25)       | 62.13 (0.81)        | 3.05 ± 3.41    |
| APPLE                 | 69.01    | 85.24     | 71.92 (7.07)        | 70.94 (-2.43)       | 64.48 (3.55)        | 72.32 (5.10)        | 3.32 ± 3.55    | 65.89    | 76.13     | 71.51 (10.54)       | 64.22 (1.75)        | 69.26 (4.51)        | 63.28 (1.96)        | 4.69 ± 3.55    |
| PGFed                 | 69.32    | 84.76     | 74.54 (9.69)        | 72.25 (-1.12)       | 66.72 (5.79)        | 73.95 (6.74)        | 5.27 ± 3.96    | 66.40    | 77.13     | 70.21 (9.24)        | 69.16 (6.69)        | 68.33 (3.58)        | 64.98 (3.66)        | 5.79 ± 2.35    |

Setting 1: RSNA only

Setting 2: Multiple datasets

# Case study: PFL for real-world breast cancer detection

- Results

| Algorithm             | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std |
|-----------------------|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|
| <i>Real data only</i> |          |           |                     |                     |                     |                     |                |          |           |                     |                     |                     |                     |                |
| Local                 | 66.59    | 82.29     | 64.85 (0.00)        | 73.37 (0.00)        | 60.93 (0.00)        | 67.21 (0.00)        | 0.00 ± 0.00    | 62.37    | 73.86     | 60.96 (0.00)        | 62.47 (0.00)        | 64.74 (0.00)        | 61.32 (0.00)        | 0.00 ± 0.00    |
| FedAvg                | 63.09    | 84.26     | 74.24 (9.39)        | 60.90 (-12.47)      | 57.29 (-3.64)       | 66.99 (-0.22)       | -1.73 ± 7.82   | 64.48    | 75.43     | 69.55 (8.59)        | 65.50 (3.04)        | 64.49 (-0.25)       | 62.13 (0.81)        | 3.05 ± 3.41    |
| APPLE                 | 69.01    | 85.24     | 71.92 (7.07)        | 70.94 (-2.43)       | 64.48 (3.55)        | 72.32 (5.10)        | 3.32 ± 3.55    | 65.89    | 76.13     | 71.51 (10.54)       | 64.22 (1.75)        | 69.26 (4.51)        | 63.28 (1.96)        | 4.69 ± 3.55    |
| PGFed                 | 69.32    | 84.76     | 74.54 (9.69)        | 72.25 (-1.12)       | 66.72 (5.79)        | 73.95 (6.74)        | 5.27 ± 3.96    | 66.40    | 77.13     | 70.21 (9.24)        | 69.16 (6.69)        | 68.33 (3.58)        | 64.98 (3.66)        | 5.79 ± 2.35    |

Setting 1: RSNA only

Setting 2: Multiple datasets

# Case study: PFL for real-world breast cancer detection

- Results

| Algorithm  | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std | Mean AUC  | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std |
|--|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|---|-----------|---------------------|---------------------|---------------------|---------------------|----------------|
| <i>Real data only</i>  |          |           |                     |                     |                     |                     |                |   |           |                     |                     |                     |                     |                |
| Local  | 66.59    | 82.29     | 64.85 (0.00)        | 73.37 (0.00)        | 60.93 (0.00)        | 67.21 (0.00)        | 0.00 ± 0.00    | 62.37   | 73.86     | 60.96 (0.00)        | 62.47 (0.00)        | 64.74 (0.00)        | 61.32 (0.00)        | 0.00 ± 0.00    |
| FedAvg   | 63.09    | 84.26     | 74.24 (9.39)        | 60.90 (-12.47)      | 57.29 (-3.64)       | 66.99 (-0.22)       | -1.73 ± 7.82   | 64.48   | 75.43     | 69.55 (8.59)        | 65.50 (3.04)        | 64.49 (-0.25)       | 62.13 (0.81)        | 3.05 ± 3.41    |
| APPLE  | 69.01    | 85.24     | 71.92 (7.07)        | 70.94 (-2.43)       | 64.48 (3.55)        | 72.32 (5.10)        | 3.32 ± 3.55    | 65.89   | 76.13     | 71.51 (10.54)       | 64.22 (1.75)        | 69.26 (4.51)        | 63.28 (1.96)        | 4.69 ± 3.55    |
| PGFed  | 69.32    | 84.76     | 74.54 (9.69)        | 72.25 (-1.12)       | 66.72 (5.79)        | 73.95 (6.74)        | 5.27 ± 3.96    | 66.40   | 77.13     | 70.21 (9.24)        | 69.16 (6.69)        | 68.33 (3.58)        | 64.98 (3.66)        | 5.79 ± 2.35    |
| <i>Adding synthesized data; FID score for client 1: 125.61, client 2: 220.52, client 3: 194.73, client 4: 184.41</i> |          |           |                     |                     |                     |                     |                | <i>sized data; FID score for client 1: 120.23, client 2: 200.91, client 3: N/A, client 4: 58.86</i> |           |                     |                     |                     |                     |                |

Setting 1: RSNA only

Setting 2: Multiple datasets

# Case study: PFL for real-world breast cancer detection

- Results

| Algorithm  | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std | Mean AUC | Mean Acc. | Client 1 AUC (gain) | Client 2 AUC (gain) | Client 3 AUC (gain) | Client 4 AUC (gain) | Mean gain ±std |
|--|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|----------|-----------|---------------------|---------------------|---------------------|---------------------|----------------|
| <i>Real data only</i>  |          |           |                     |                     |                     |                     |                |          |           |                     |                     |                     |                     |                |
| Local  | 66.59    | 82.29     | 64.85 (0.00)        | 73.37 (0.00)        | 60.93 (0.00)        | 67.21 (0.00)        | 0.00 ± 0.00    | 62.37    | 73.86     | 60.96 (0.00)        | 62.47 (0.00)        | 64.74 (0.00)        | 61.32 (0.00)        | 0.00 ± 0.00    |
| FedAvg   | 63.09    | 84.26     | 74.24 (9.39)        | 60.90 (-12.47)      | 57.29 (-3.64)       | 66.99 (-0.22)       | -1.73 ± 7.82   | 64.48    | 75.43     | 69.55 (8.59)        | 65.50 (3.04)        | 64.49 (-0.25)       | 62.13 (0.81)        | 3.05 ± 3.41    |
| APPLE  | 69.01    | 85.24     | 71.92 (7.07)        | 70.94 (-2.43)       | 64.48 (3.55)        | 72.32 (5.10)        | 3.32 ± 3.55    | 65.89    | 76.13     | 71.51 (10.54)       | 64.22 (1.75)        | 69.26 (4.51)        | 63.28 (1.96)        | 4.69 ± 3.55    |
| PGFed  | 69.32    | 84.76     | 74.54 (9.69)        | 72.25 (-1.12)       | 66.72 (5.79)        | 73.95 (6.74)        | 5.27 ± 3.96    | 66.40    | 77.13     | 70.21 (9.24)        | 69.16 (6.69)        | 68.33 (3.58)        | 64.98 (3.66)        | 5.79 ± 2.35    |
| <i>Adding synthesized data; FID score for client 1: 125.61, client 2: 220.52, client 3: 194.73, client 4: 184.41</i> |          |           |                     |                     |                     |                     |                |          |           |                     |                     |                     |                     |                |
| Local  | 67.06    | 84.73     | 69.42 (0.00)        | 75.31 (0.00)        | 56.49 (0.00)        | 63.36 (0.00)        | 0.00 ± 0.00    | 63.01    | 72.89     | 64.89 (0.00)        | 60.90 (0.00)        | 64.74 (0.00)        | 61.52 (0.00)        | 0.00 ± 0.00    |
| FedAvg   | 64.51    | 84.49     | 70.47 (1.05)        | 55.28 (-20.03)      | 62.69 (6.20)        | 71.60 (8.24)        | -1.14 ± 11.22  | 64.20    | 74.87     | 68.23 (3.34)        | 64.80 (3.90)        | 65.84 (1.10)        | 63.26 (1.74)        | 2.52 ± 1.14    |
| APPLE  | 67.40    | 84.89     | 69.01 (-0.41)       | 63.69 (-11.62)      | 66.51 (10.02)       | 68.94 (5.58)        | 0.89 ± 8.12    | 66.38    | 76.74     | 67.59 (2.70)        | 63.52 (2.62)        | 72.72 (7.98)        | 65.77 (4.25)        | 4.39 ± 2.17    |
| PGFed  | 70.41    | 84.76     | 72.61 (3.19)        | 68.56 (-6.75)       | 67.77 (11.28)       | 74.32 (10.96)       | 4.67 ± 7.35    | 67.41    | 77.03     | 71.16 (6.27)        | 71.72 (10.82)       | 70.38 (5.64)        | 64.42 (2.90)        | 6.41 ± 2.85    |

Setting 1: RSNA only

Setting 2: Multiple datasets

## Case study: PFL for real-world breast cancer detection

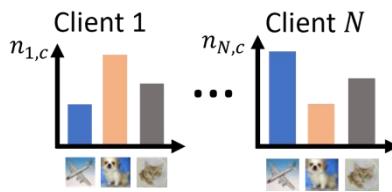
- Summary
  - Traditional FL (e.g. FedAvg) < Local with severe heterogeneity
  - Proposed PFL (APPLE, PGFed) > FedAvg and Local
  - Generated data: higher perceptive quality usually translate to higher performance
- For a FL system deployed in real-world
  - Can use more advanced model with better pretraining.
  - Practical hyperparameter tuning is hard, alg's with more hyperparameters is harder to deploy
    - Research: same values for all clients, global metrics for selection.
    - Real-world: clients can use different values, local metrics for selection.
    - Example, hyperparameters:  $a$  2 values,  $b$  3 values, research has  $2 \times 3 = 6$  combinations, real-world  $a_1, b_1, a_2, b_2, a_3, b_3, a_4, b_4$ :  $(2 \times 3)^4 = 1296$  combinations

# Overview

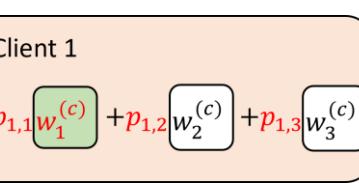
- Federated learning: introduction
- Federated Learning with Shared Label Distribution for Medical Image Classification (FedSLD)
- Adapt to Adaptation: Learning Personalization for Cross-Silo Federated Learning (APPLE)
- PGFed: Personalize Each Client's Global Objective for Federated Learning (PGFed)
- Mixture of Experts Made Personalized: Federated Prompt Learning for Vision-Language Models (pFedMoAP)
- Case Study: Personalized, Real-World, and Cross-Silo Federated Learning for Breast Cancer Detection
- **Summary**

# Summary of the four FL/PFL algorithms

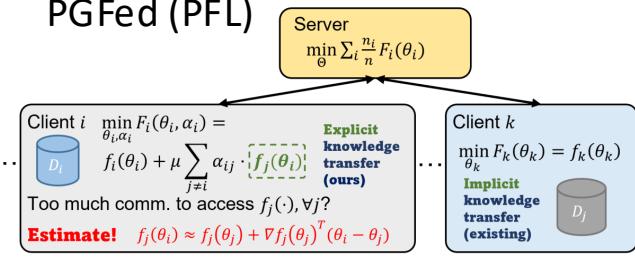
FedSLD (Global FL)



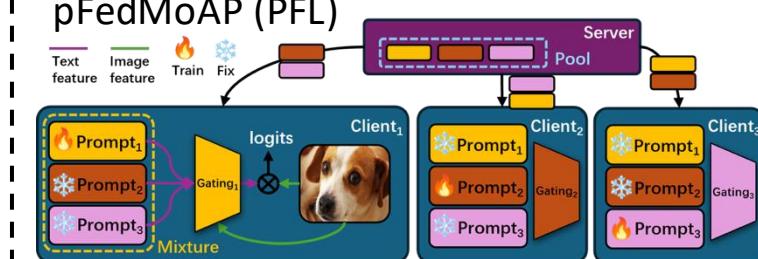
APPLE (PFL)



PGFed (PFL)



pFedMoAP (PFL)



Estimate prior and reweighting PFL with adaptive aggregation

+ Sharable label dist.  
+ Reweights sample loss

- Limited performance gain  
- Only considers label shift

Explicit and efficient personalized global objective with first order approximation

+ More generalized personalized models  
+ Explicitness with  $O(N)$  communication

- Slightly larger communication than FedAvg  
- Requires extra server computation/storage

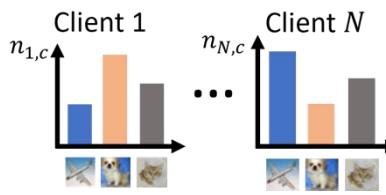
PFL prompt learning for CLIP with attention-based gating network in a MoE structure

+ Pre-aggregated prompts sharing allows MoE  
+ Flexible and robust attention-based gating

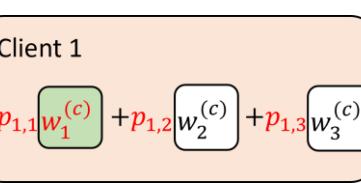
- Clients needs to be able to run CLIP  
- Computation slightly  $\uparrow$  as #experts  $\uparrow$

# Summary of the four FL/PFL algorithms

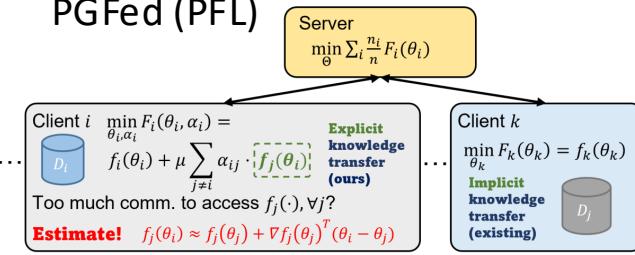
FedSLD (Global FL)



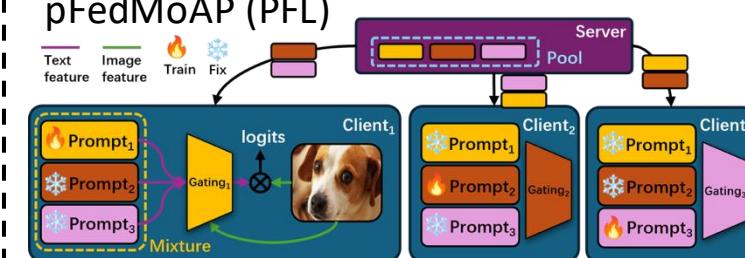
APPLE (PFL)



PGFed (PFL)



pFedMoAP (PFL)



Estimate prior and reweighting

PFL with adaptive aggregation

Explicit and efficient personalized global objective with first order approximation

PFL prompt learning for CLIP with attention-based gating network in a MoE structure

## Summary of FL case study: breast cancer detection

- Traditional FL (e.g. FedAvg) < Local with sever heterogeneity
- Proposed PFL (APPLE, PGFed) > FedAvg and Local
- Generated data: higher perceptive quality usually translate to higher performance
- Deploying FL system in real-world will face more challenges (e.g. hyperparameter tuning across heterogeneous clients).

## Future directions

- Federated learning with large foundation models.
- Synthetic data generation and augmentation via foundation models.
- Enhancing privacy, security, and trustworthiness in FL.
- Developing multimodal federated systems across diverse domains.

## Acknowledgements

- My advisor Dr. Wu and my committee Dr. Zhou, Dr. Jia, and Dr. Tang
- My wife Tianling
- Collaborators: Dr. Chen, Matias and Guangyu from UCF
- ICCI: Oliver, Chang, Zhengbo, Dooman, Jiren, Zhiwei, Giacomo
- Internship: at  , and **Sony AI**

# Intelligent Computing for Clinical Imaging (ICCI) Lab

## -Interfacing computational and medical sciences



### Selected funding

- ❖ NIH OT (#1OT2OD037972-01)
- ❖ NIH/NCI R01 (#CA193603)
- ❖ NIH/NCI R01 (#CA218405)
- ❖ NSF/NIH joint R01(#EB032896)
- ❖ NSF (CICI: SIVD: #2115082)
- ❖ NSF (CBET #2229156)
- ❖ NIH OT #1OT2OD032701-01
- ❖ NIH R01 Supplement (#CA193603-S; #EB032896-03S1)
- ❖ UPMC Enterprise (Early Commercialization)
- ❖ RSNA Research Scholar Grant (#RSCH1530)
- ❖ PA Breast Cancer Coalition
- ❖ Jewish Healthcare Foundation
- ❖ Pittsburgh Foundation
- ❖ Amazon AWS Machine Learning Research Award
- ❖ Stanly Marks Research Foundation
- ❖ Pitt Momentum Funds Scaling Grant

# Thank you!

Jun Luo  
[jul117@pitt.edu](mailto:jul117@pitt.edu)