# PGFed: Personalize Each Client's Global Objective for Federated Learning
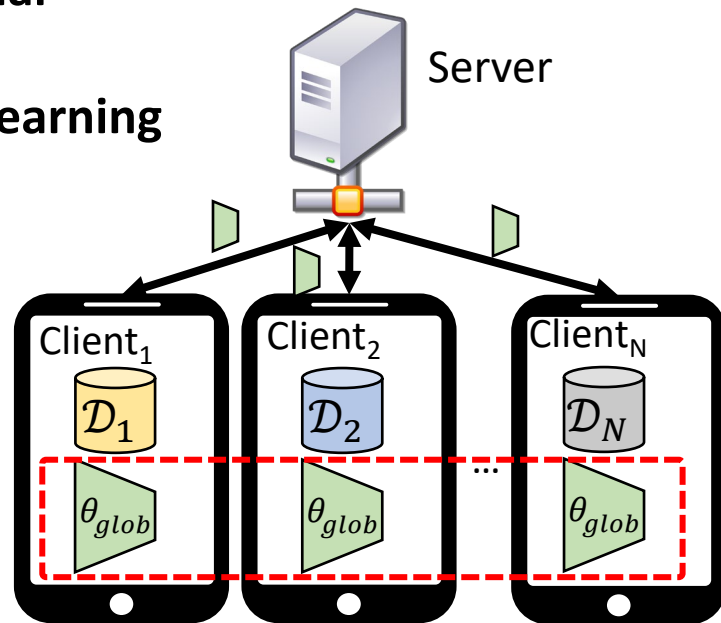
**Jun Luo**[†], Matias Mendieta[‡], Chen Chen[‡], and Shandong Wu[†]

[†] University of Pittsburgh, Pittsburgh, PA, USA
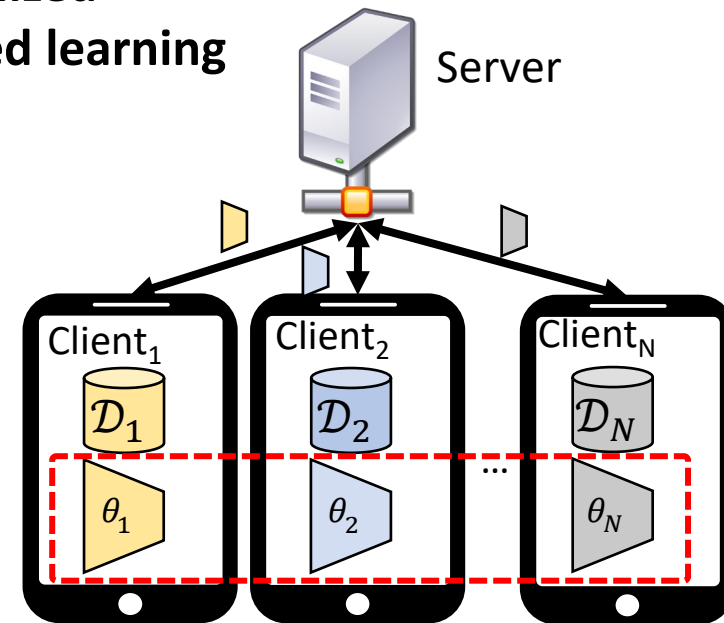[‡] University of Central Florida, Orlando, FL, USA

# Background

**Conventional (Global) federated learning**



$$\min_{\theta_{glob}} F(\theta_{glob}) = \min_{\theta_{glob}} \sum_i p_i F_i(\theta_{glob})$$
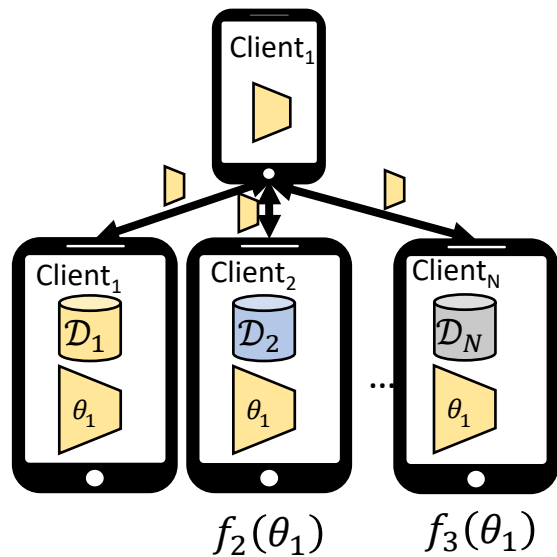
**Personalized federated learning**



$$\min_{\Theta} F(\Theta) = \min_{\theta_1, \theta_2, .., \theta_N} \sum_i p_i F_i(\theta_i)$$

# Background

- In existing personalized FL (pFL) algorithms (with heterogeneous data), *the way in which the **collaborative knowledge** transfers from the server to the clients is **implicit**.*

  - **Collaborative knowledge**: non-local information
    - E.g., Global FL's objective: $F(\theta_{glob}) = \sum_i p_i F_i(\theta_{glob})$

  - **Explicitness** (as opposite of **implicitness**): Direct engagement of multiple clients' empirical risks (explicit since not embed non-local info into model weights or regularization)
    - E.g., Global FL's objective: $F(\theta_{glob}) = \sum_i p_i F_i(\theta_{glob})$ where $F_i(\theta_{glob}) = f_i(\theta_{glob})$

# 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.
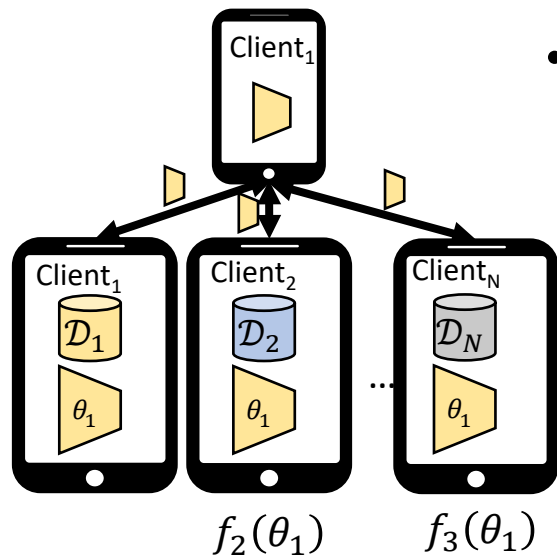
# Motivation

- 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.
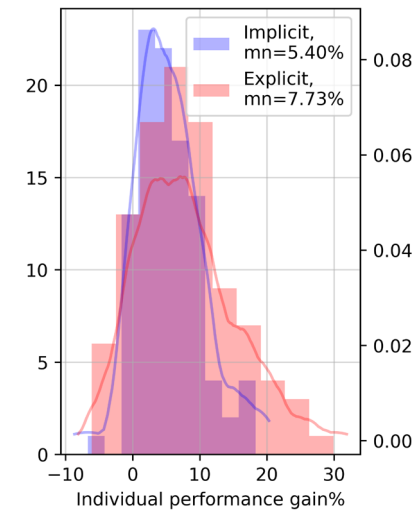


- Toy experiment on exemplar design
  - Cifar10, 100 heterogeneous clients
  - Explicit: $F_i(\theta_i) = f_i(\theta_i) + \frac{\mu}{N-1}\sum_{j \neq i} f_j(\theta_i)$
  - Implicit: $F_i(\theta_i) = f_i(\theta_i)$ (local model of FedAvg)

# Motivation

- Difficulty to achieve explicitness
    - $O(N^2)$ communication overhead
    - Proper coefficient for each non-local risk

# Motivation

- Difficulty to achieve explicitness
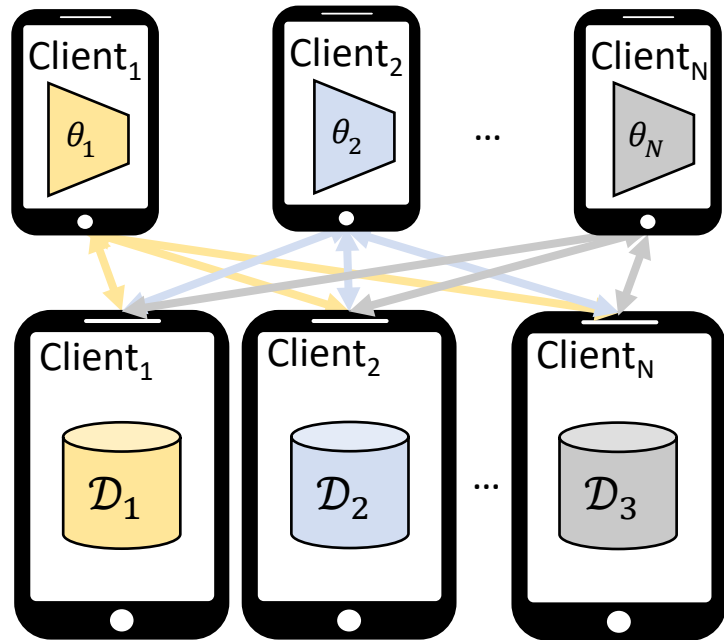  - $O(N^2)$ communication overhead
  - Proper coefficient for each non-local risk

# Method

- Proposed solution: **PGFed**
  - ✓ Estimate $f_j(\theta_i) \approx f_j(\theta_j) + \nabla f_j(\theta_j)^T(\theta_i - \theta_j), O(N^2) \rightarrow O(N)$
  - ✓ Use adaptive coefficient $\alpha_{ij} \forall i, j \in [N]$



Server
$\min_{\Theta} \sum_i \frac{n_i}{n} F_i(\theta_i)$

Client $i$ $\min_{\theta_i, \alpha_i} F_i(\theta_i, \alpha_i) =$

$f_i(\theta_i) + \mu \sum_{j \neq i} \alpha_{ij} \cdot \boxed{f_j(\theta_i)}$

**Explicit knowledge transfer (ours)**

Too much comm. to access $f_j(\cdot), \forall j$?

**Estimate!** $f_j(\theta_i) \approx f_j(\theta_j) + \nabla f_j(\theta_j)^T(\theta_i - \theta_j)$

Client $k$
$\min_{\theta_k} F_k(\theta_k) = f_k(\theta_k)$

**Implicit knowledge transfer (existing)**

# Method

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

  - Global objective:  $$\min_{\boldsymbol{\Theta}, \boldsymbol{A}} F(\boldsymbol{\Theta}, \boldsymbol{A}) = \min_{\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_N, \boldsymbol{\alpha}_1, \ldots, \boldsymbol{\alpha}_N} \sum_{i=1}^{N} p_i F_i(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i)$$

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

  - 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

$$F_i(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) \approx f_i(\boldsymbol{\theta}_i) + \mathcal{R}_{aux}^{[N]}(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i)$$

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

# Method

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

- **Gradient-based update**
  - W.r.t $\theta_i$:
    $$\nabla_{\boldsymbol{\theta}_i} F_i(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) = \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i) + \nabla_{\boldsymbol{\theta}_i} \mathcal{R}_{aux}^{[N]}(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i)$$
    $$= \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i) + \mu \underbrace{\sum_{j \in [N]} \alpha_{ij} \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)}_{\tilde{\boldsymbol{g}}_{[N]}} .$$

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

# Method

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

- **Gradient-based update**
  - W.r.t $\alpha_{ij}$:
  
$$\nabla_{\alpha_{ij}} F_i(\boldsymbol{\theta}_i, \boldsymbol{\alpha}_i) = \mu \left( f_j(\boldsymbol{\theta}_j) + \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j) \right)$$

$$= \underbrace{\mu \left( f_j(\boldsymbol{\theta}_j) - \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T \boldsymbol{\theta}_j \right)}_{g_\alpha^{(1)}} + \underbrace{\mu \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j)^T \boldsymbol{\theta}_i}_{g_\alpha^{(2)}} .$$

  - $g_\alpha^{(1)}$ (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: $g_\alpha^{(2)} \approx \bar{g}_{[N]}^T \boldsymbol{\theta}_i = \dfrac{\mu}{N} \left( \sum_{j \in [N]} \nabla_{\boldsymbol{\theta}_j} f_j(\boldsymbol{\theta}_j) \right)^T \boldsymbol{\theta}_i$

      - Client $j$ uploading local gradient

# Method

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

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

- 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 $\theta_1^T, ..., \theta_N^T$.

**ServerExecute:**
1: Initialize $\alpha_{ij} \leftarrow 1/M \ \forall i, j \in [N]$, global model $\theta_{glob}^0$
2: $\mathbf{A}[i] \leftarrow \boldsymbol{\alpha}_i \ \forall i \in [N]$
3: **for** $t \leftarrow 1, 2, ..., T$ **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:              $\theta_i^t, g_\alpha^{(1)}, \nabla f(\theta_i^t), \alpha_i \leftarrow$ **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:              $\theta_i^t, g_\alpha^{(1)}, \nabla f(\theta_i^t), \alpha_i \leftarrow$ **ClientUpdate**$(\theta_{glob}^{t-1}, t,$
             $\tilde{g}_{S_{t-1}}, \bar{g}_{S_{t-1}}, g_{t-1}^{(1)})$
13:          **end if**
14:          // the next line records the values for next round
15:          $\mathbf{A}[i] \leftarrow \boldsymbol{\alpha}_i; g_t^{(1)}[i] \leftarrow g_\alpha^{(1)}; \nabla_t[i] \leftarrow \nabla f(\theta_i^t)$
16:          $\theta_{glob}^t \leftarrow \sum_{i \in S_t} p_i \theta_i^t$
17:      **end for**
18:      **for** $i \in ([N] - S_t)$ **in parallel do**
19:          $\theta_i^t \leftarrow \theta_i^{t-1}; \tilde{g}_i^t \leftarrow \tilde{g}_i^{t-1}$
20:      **end for**
21: **end for**
22: **return** $\theta_1^T, ..., \theta_N^T$

---

**ClientUpdate**$(\theta_{global}^{t-1}, t \ (, \tilde{g}, \bar{g}, g_{t-1}^{(1)}))$:
1: **if** t=1 **then**
2:      $\theta_i^t \leftarrow$ **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:          $\theta_i^t \leftarrow \theta_i^t - \eta_1(\nabla f(\theta_i^t, \mathcal{B}) + \tilde{g}_t^i)$
9:          $g^{(2)} = \bar{g}^T \theta_i$
10:          $\forall j \in g_{t-1}^{(1)}: \ \alpha_{ij} \leftarrow \alpha_{ij} - \eta_2(g_{t-1}^{(1)}[j] + g^{(2)})$
11:      **end for**
12: **end if**
13: $g_\alpha^{(1)} \leftarrow \mu \left( f(\theta_i^t) - \nabla f(\theta_i^t)^T \theta_i^t \right)$ // for next round
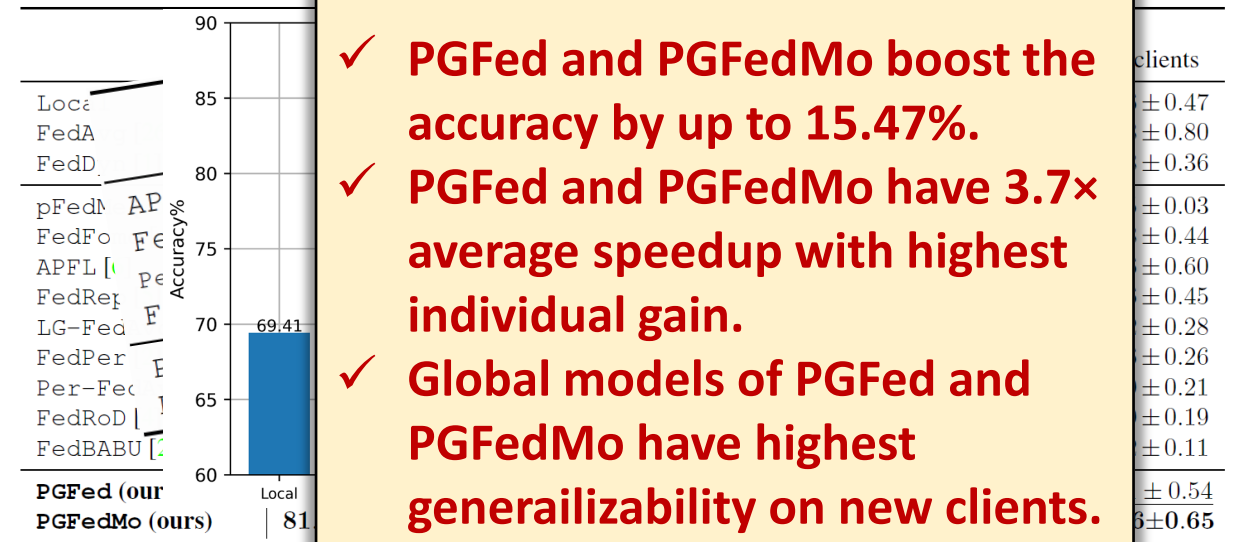14: **return** $\theta_i^t, g_\alpha^{(1)}, \nabla f(\theta_i^t), \alpha_i$

# Experiments & results

- Settings
  - Datasets: CIFAR10, CIFAR100, OrganAMNIST, Office-home
  - Heterogeneity: Dir $\alpha = 0.3, 1.0$
  - Number of clients: 20, 25, 50, 100 clients
  - Metrics
    - Mean local test accuracy
    - Mean individual gain over `Local`
    - #Rounds to reach 70% acc. & speedup
    - Accuracy of fine-tuning resulting global model on new clients
    - Throughput
    - Etc. (see full paper)



## Main takeaways:

✓ **PGFed and PGFedMo boost the accuracy by up to 15.47%.**
✓ **PGFed and PGFedMo have 3.7× average speedup with highest individual gain.**
✓ **Global models of PGFed and PGFedMo have highest generailizability on new clients.**

# More experiments & results



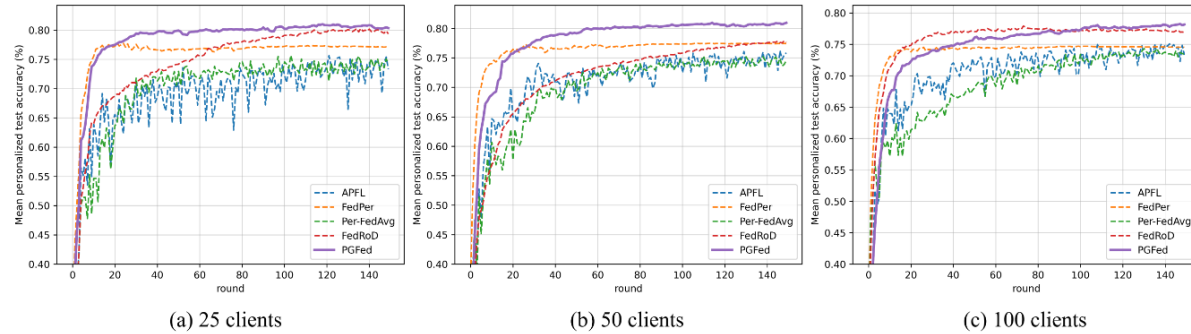(a) 25 clients    (b) 50 clients    (c) 100 clients

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.
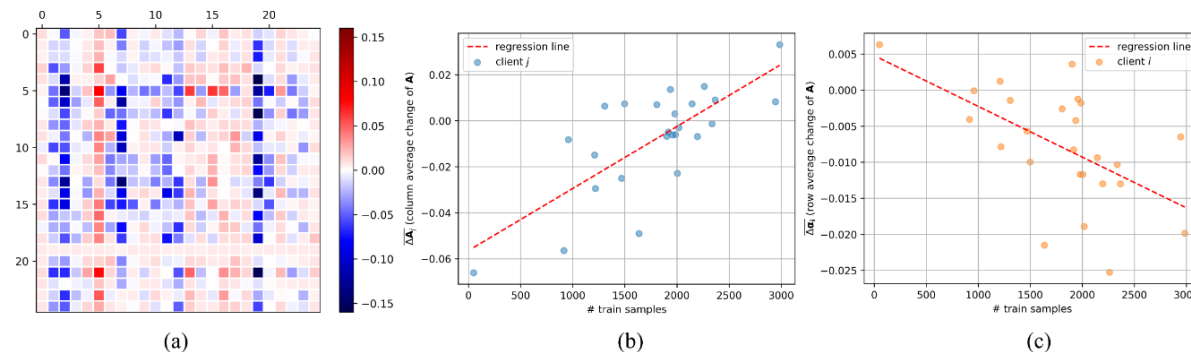


(a)    (b)    (c)

Figure 2. Visualization of the change in $A$. Figure (a) is a heat map of the change in $A$. For Figure (b) and (c), the Y-axis of Figure (b) represents the column average of the change in $A$ (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 $A$ (the average change of weights of the auxiliary risk on client $i$). Through the regression line, we verify the positive correlation between $\overline{\Delta A_j}$ and $n_j$ in Figure (b), and the negative correlation between $\overline{\Delta \alpha_i}$ 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 | $19.11 \pm 1.55$ | $44.67 \pm 0.61$ | $\mathbf{50.40 \pm 0.56}$ | $25.85 \pm 0.88$ | $35.00 \pm 0.41$ |
| 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 | $18.18 \pm 3.54$ | $42.10 \pm 2.31$ | $43.51 \pm 0.91$ | $\mathbf{26.81 \pm 1.86}$ | $33.38 \pm 0.29$ |
| PGFed | $\mathbf{22.40 \pm 0.26}$ | $\mathbf{46.48 \pm 1.00}$ | $49.86 \pm 2.14$ | $26.04 \pm 0.80$ | $\mathbf{36.19 \pm 0.92}$ |
| PGFedMo | $\underline{22.16 \pm 0.45}$ | $\underline{45.88 \pm 0.83}$ | $\underline{49.45 \pm 0.19}$ | $\underline{26.60 \pm 0.99}$ | $\underline{36.02 \pm 0.20}$ |

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 sample 50% Dir(1.0) | 50 clients sample 25% Dir(0.3) | 100 clients sample 25% 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 | $99.20 \pm 0.04$ | $\mathbf{99.17 \pm 0.05}$ | $\mathbf{98.94 \pm 0.02}$ |
| PGFedMo | $\mathbf{99.21 \pm 0.04}$ | $99.17 \pm 0.07$ | $98.86 \pm 0.06$ |

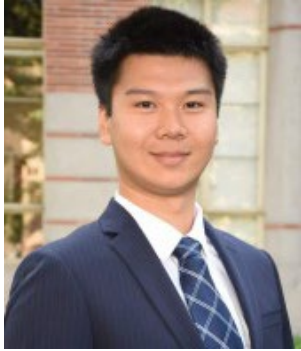Table 1. Mean and standard deviation over three trials of the mean personalized test accuracy (%) on OrganAMNIST

| | 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
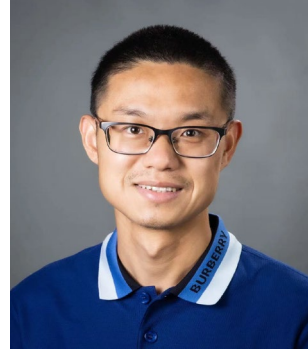
***More details in full paper…***
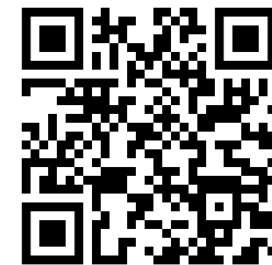
# Acknowledgements

**Jun Luo**  Matias Mendieta  Dr. Chen Chen  Dr. Shandong Wu

Intelligent Computing for Clinical Imaging (ICCI) Lab, University of Pittsburgh

Jun Luo
jul117@pitt.edu

# Thank you!



Full paper



Code

Also check out our poster
Oct. 4th (Wed.)
02:30 PM-04:30 PM

Jun Luo
University of Pittsburgh
jul117@pitt.edu