



DM-PFL: Hitchhiking Generic Federated Learning for Efficient Shift-Robust Personalization

Wenhai Zhang¹, Zimu Zhou², Yansheng Wang¹, Yongxin Tong¹

¹State Key Laboratory of Software Development Environment, Beijing Advanced Innovation Center for Future Blockchain and Privacy Computing, Beihang University, Beijing, China

²City University of Hong Kong, Hong Kong SAR, China

¹{garyzhang, arthur_wang, yxtong}@buaa.edu.cn

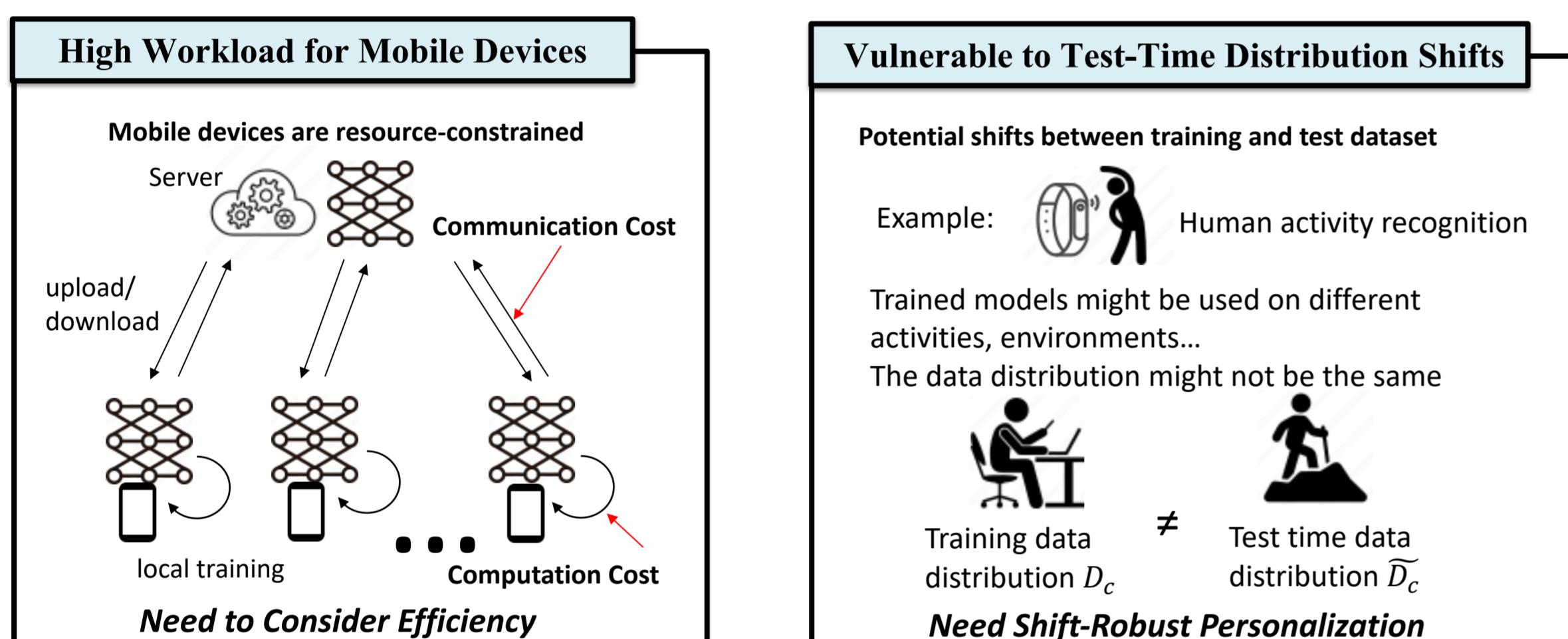
²zimuzhou@cityu.edu.hk

Introduction

- Unlike the generic federated learning that learns a global model for all clients, personalized federated learning trains client-specific models collaboratively, which holds potential for various mobile applications with heterogeneous data.



- Despite extensive personalized FL proposals two drawbacks impede their practical usage for mobile applications.



- In response, we explore efficient shift-robust personalization for FL, which aims to train personalized models that are robust to test-time distribution shifts with low training overhead.

DM-PFL Overview

Review on Generic and Personalized FL

Generic FL:
- Train a global model with better generalization ability that performs well under a wider range of data distributions.

$$\min_{\theta_g} l(\theta_g) = \sum_{c=1}^C \frac{|D_c|}{|D|} \mathbb{E}[\mathcal{L}(\theta_g; D_c)]$$

θ_g : the global model

Optimizing performance on the aggregation of large amounts of data across clients.

- Good local performance
- Vulnerable to distribution shifts.

Personalized FL:
- Train client-specific personalized models for better performance on client's local data distributions.

$$\min_{\theta_1, \theta_2, \dots, \theta_C} l(\theta_1, \dots, \theta_C) = \sum_{c=1}^C \frac{|D_c|}{|D|} \mathbb{E}[\mathcal{L}(\theta_c; D_c)]$$

θ_c : each client's local model

Optimize performance on the clients' local data.

- Sub-optimal local performance
- Robust to distribution shifts.

Design Rationales:

- Hitchhike the global model (θ_g) to make our personalized models (θ_c) more shift-robust.
- Use mask-based training to enforce model sparsity and enable weight-level parameter sharing.

Our objectives and constraints:

$$\begin{aligned} \text{Objective: } & \min_{\theta_1, \theta_2, \dots, \theta_C} l(\theta_1, \dots, \theta_C) = \sum_{c=1}^C \frac{|D_c|}{|D|} \mathbb{E}[\mathcal{L}(\theta_c; D_c)] \\ \text{s.t.: } & \|\theta_c\|_0 \leq |\theta_c| * (1 - \$), \forall c \in \{1, 2, \dots, C\} \end{aligned}$$

Constraint: Take test-time distribution shift \bar{D}_c into consideration

Dual Masking Mechanism

Components

We use global mask m_g and a set of personalized masks $\{m_c\}_{c=1}^C$ with weights w_g and w_c to construct dual models.
 $\theta_g = w_g \odot m_g$
 $\theta_c = w_g \odot (m_g \cap m_c) + w_c \odot (m_c - m_g)$

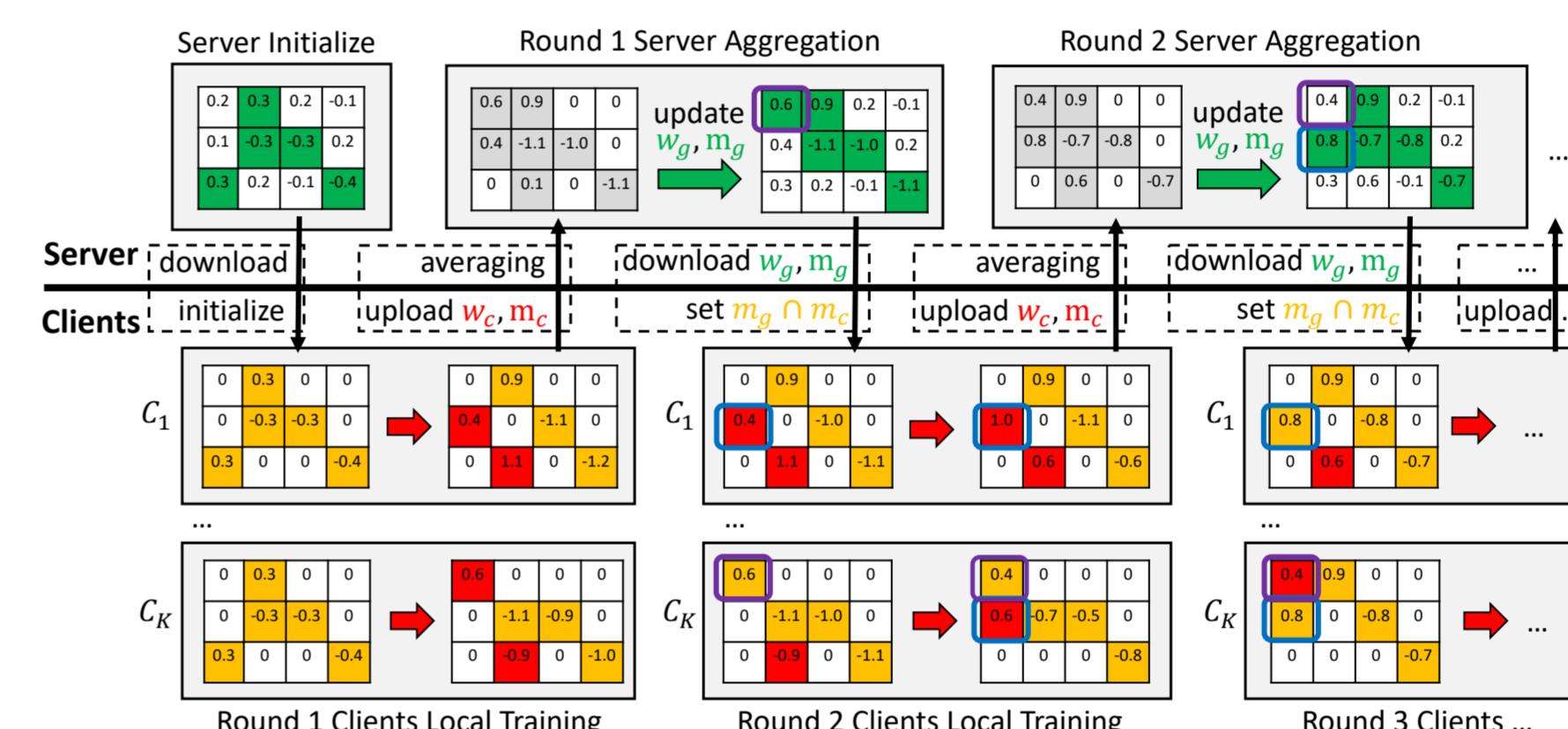
- Each client inherits a different portion of the global weights w_g according to its overlap in mask positions with the global mask, i.e., $m_g \cap m_c$.
- For the remaining positions $m_c - m_g$ in its personalized mask, each client inherits weights from its own w_c .

DM-PFL Algorithm

Two-staged sparse training algorithm

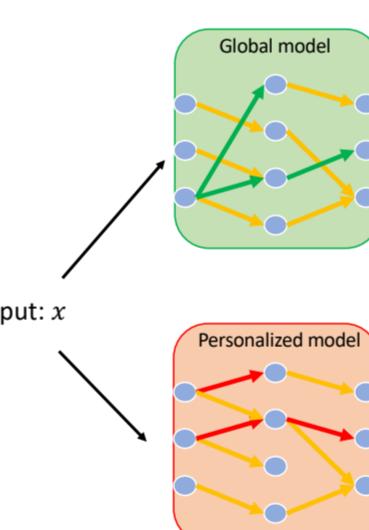
- In the first train dual masks stage, we train weights and masks of the dual models interactively.
- In the refine dual weights stage, we first refine w_g then refine w_c under fixed masks.

Example of train dual masks stage



- Dual masks and weights are trained interactively. Shared weights can transform into personalized weights (purple box), and vice versa (blue box).

Adaptive inference



- We may further utilize θ_g to improve shift-robustness, which adaptively ensembles θ_g and θ_c based on the estimated degree of test-time shift.

Experimental Evaluation

Performance with and without shifts

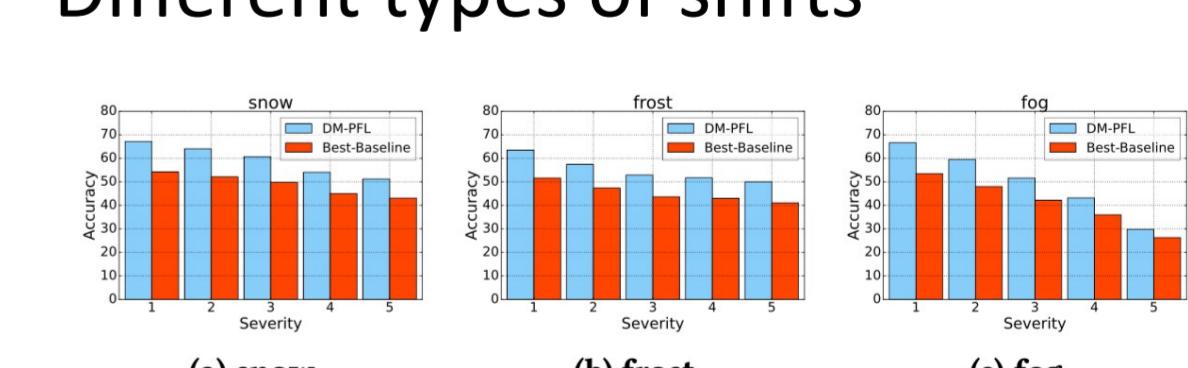
Different Non-IID settings

Non-IID	Dirichlet distribution-based CIFAR10		Pathological CIFAR10		Realistic HAR		Realistic FEMNIST			
	Dataset	HAR	FEMNIST	Dataset	HAR	FEMNIST	Dataset	HAR	FEMNIST	
Methods	w/o	w/ 10%	w/o	w/ 10%	w/o	w/ 10%	w/o	w/ 10%	w/o	w/ 10%
FedAvg	97.85	97.53	97.84	97.45	97.45	98.05	98.55	98.55	98.54	98.76
L-G-FedAvg	95.52	94.13	95.52	94.13	94.13	95.12	95.12	95.12	95.12	95.20
FedPer	88.82	82.12	88.82	82.12	82.12	89.73	89.73	89.73	89.73	89.73
FedRep	89.19	10.48	89.19	10.48	10.48	89.69	91.34	91.34	91.34	91.51
Ditto	88.20	9.12	88.20	9.12	9.12	89.25	89.25	89.25	89.25	89.25
APFL	89.73	10.48	89.73	10.48	10.48	89.80	90.55	90.55	90.55	90.55
DM-PFL	89.29	11.65	89.29	11.65	11.65	90.06	91.39	91.39	91.39	91.39
DM-PFL+	88.16	12.46	87.33	12.46	12.46	88.59	87.65	87.65	87.65	87.65

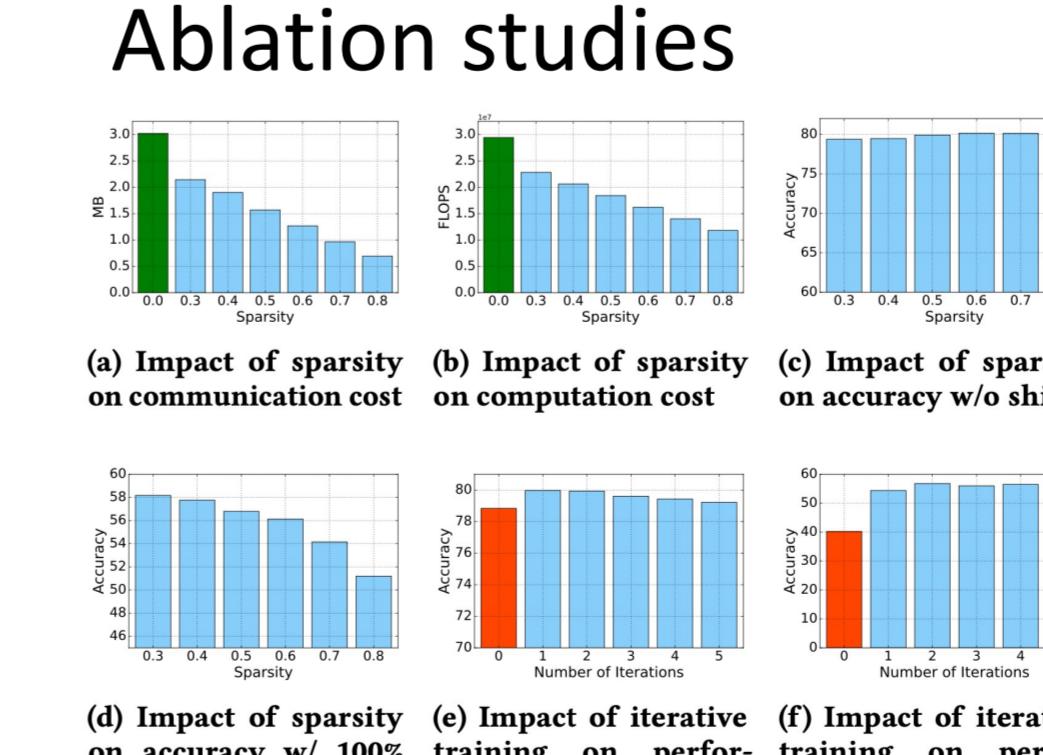
Different Non-IID settings

Type	Local	Dirichlet distribution-based CIFAR10	Pathological CIFAR10	Realistic HAR	Realistic FEMNIST	
Dataset	w/o shift	w/ 20% shift	w/ 40% shift	w/ 40% shift	w/ 100% shift	
Method	30.70	17.07	22.93	19.81	14.81	
FedAvg	54.11	54.93	54.11	54.11	54.11	
GFL	60.61	68.21	59.84	84.84	59.53	19.78
FedFox	58.64	10.60	58.52	10.50	58.52	10.46
PL	69.73	19.29	60.53	14.52	51.09	15.76
FedPer	76.73	7.27	68.08	4.45	51.02	3.67
Ditto	79.70	8.88	73.17	3.24	61.50	1.59
APFL	78.54	10.48	78.54	10.48	78.54	10.48
Ours	79.41	9.61	75.36	9.55	66.99	11.46
DM-PFL	77.22	10.51	73.81	10.80	64.64	17.27
DM-PFL+	77.22	10.51	73.81	10.80	64.64	17.27

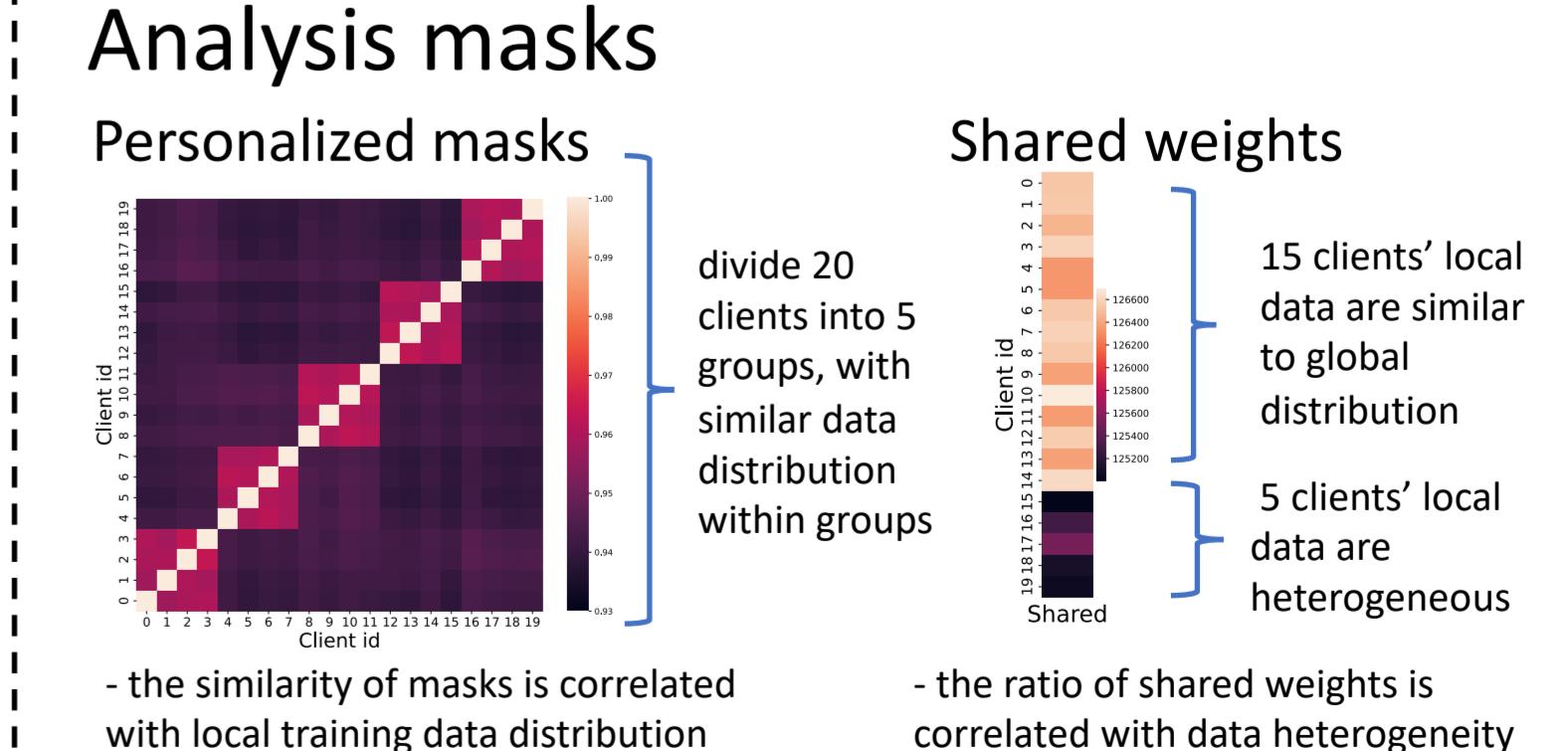
Different types of shifts



Ablation studies



Analysis masks



Acknowledgment

We are grateful to anonymous reviewers for their constructive comments. This work is partially supported by National Science Foundation of China (NSFC) under Grant No. U21A20516 and 62076017, the Beihang University Basic Research Funding No. YWF-22-L-531, the Funding No. 22-TQ23-14-ZD-01-001, the CityU APRC grant No. 9610633, and WeBank Scholars Program. Yongxin Tong is the corresponding author.