

# Tackling Data Heterogeneity in Federated Learning with Class Prototypes

Yutong Dai<sup>1</sup>   Zeyuan Chen<sup>2</sup>   Junnan Li<sup>2</sup>  
Shelby Heinecke<sup>2</sup>   Lichao Sun<sup>1</sup>   Ran Xu<sup>2</sup>

Lehigh University<sup>1</sup> Salesforce Research<sup>2</sup>



February, 2023

Washington, DC, USA



# Outline

## ① Problem

## ② Methodology

- A Motivating Example
- Proposed Method

## ③ Numerical Results

- Class Semantics
- Accuracy

# Outline

## 1 Problem

## 2 Methodology

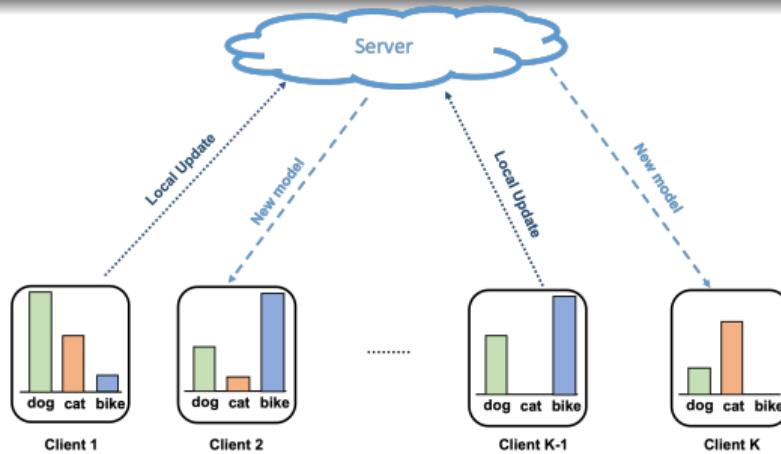
- A Motivating Example
- Proposed Method

## 3 Numerical Results

- Class Semantics
- Accuracy

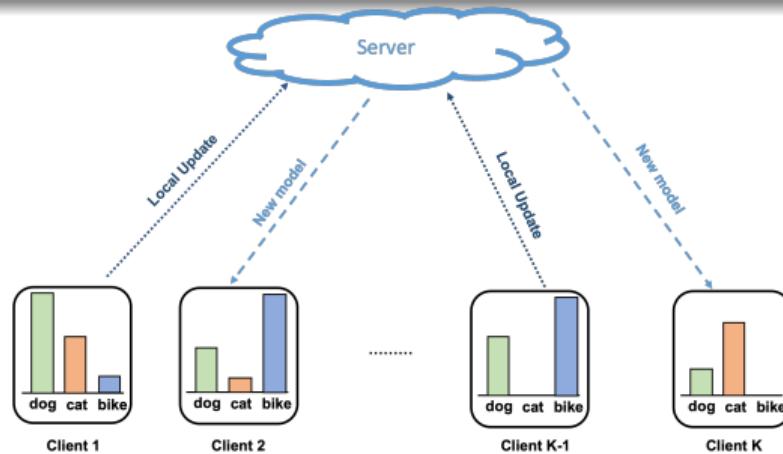
# Introduction

Federated learning and analytics are a distributed approach for collaboratively learning models (or statistics) from **decentralized data**, motivated by and designed for **privacy** protection.



# Introduction

Federated learning and analytics are a distributed approach for collaboratively learning models (or statistics) from **decentralized data**, motivated by and designed for **privacy** protection.



Address data heterogeneity with a particular focus on class imbalance.

# Introduction (Cont'd)

## Classic Problem Formulation

$$F^{\text{ERM}}(\mathbf{x}) = \sum_{i=1}^M p_i F_i^{\text{ERM}}(\mathbf{x}), \quad \text{where } F_i^{\text{ERM}}(\mathbf{x}) = \frac{1}{|D_i|} \sum_{\xi \in D_i} f_i(\mathbf{x}, \xi) \text{ and } \sum_{i=1}^M p_i = 1$$

- **Across clients:** Heterogeneous data distribution leads to inconsistent local objective functions, which imposes challenges into the optimization process.
- **Within a client:** Imbalanced data makes the local model likely to overfit **dominant classes**.

# Introduction (Cont'd)

## Classic Problem Formulation

$$F^{\text{ERM}}(\mathbf{x}) = \sum_{i=1}^M p_i F_i^{\text{ERM}}(\mathbf{x}), \quad \text{where } F_i^{\text{ERM}}(\mathbf{x}) = \frac{1}{|D_i|} \sum_{\xi \in D_i} f_i(\mathbf{x}, \xi) \text{ and } \sum_{i=1}^M p_i = 1$$

- **Across clients:** Heterogeneous data distribution leads to inconsistent local objective functions, which imposes challenges into the optimization process.
- **Within a client:** Imbalanced data makes the local model likely to overfit **dominant classes**.

Personalized federated learning comes to the rescue.

# Brief Literature Review

- Personalized federated learning
  - **goal:** tailor personalized models to client-specific tasks;
  - **methods:** parameter decoupling, regularization, model interpolation, and more<sup>1</sup>;
    - parameter decoupling: **body** (representation learning) + **head** (classification task).
- Class-imbalance learning
  - **data-level:** Over-sampling minority classes or under-sampling majority classes
  - **algorithm-level:**
    - sample-wise<sup>2</sup> or class-wise<sup>3</sup> class-balanced losses;
    - decoupling the training procedure into the representation learning and classification phases<sup>4</sup>.

---

<sup>1</sup>Alysa Ziying Tan et al. "Towards personalized federated learning". In: *IEEE Transactions on Neural Networks and Learning Systems* (2022).

<sup>2</sup>Tsung-Yi Lin et al. "Focal loss for dense object detection". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2980–2988.

<sup>3</sup>Yin Cui et al. "Class-balanced loss based on effective number of samples". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019, pp. 9268–9277.

<sup>4</sup>Bingyi Kang et al. "Decoupling Representation and Classifier for Long-Tailed Recognition". In: *International Conference on Learning Representations*. 2019.

# Brief Literature Review

- Class-imbalance learning + FL
  - **CReFF**<sup>1</sup> adapts the idea of<sup>2</sup> into FL setting while observing the privacy.
  - **FedROD**<sup>3</sup> designs a two-head-one-body architecture, where one head is trained with class-balanced loss while the other head is trained with empirical loss.

---

<sup>1</sup>Xinyi Shang et al. “Federated Learning on Heterogeneous and Long-Tailed Data via Classifier Re-Training with Federated Features”. In: *arXiv preprint arXiv:2204.13399* (2022).

<sup>2</sup>Bingyi Kang et al. “Decoupling Representation and Classifier for Long-Tailed Recognition”. In: *International Conference on Learning Representations*. 2019.

<sup>3</sup>Hong-You Chen and Wei-Lun Chao. “On Bridging Generic and Personalized Federated Learning for Image Classification”. In: *International Conference on Learning Representations*. 2021.

# Outline

1 Problem

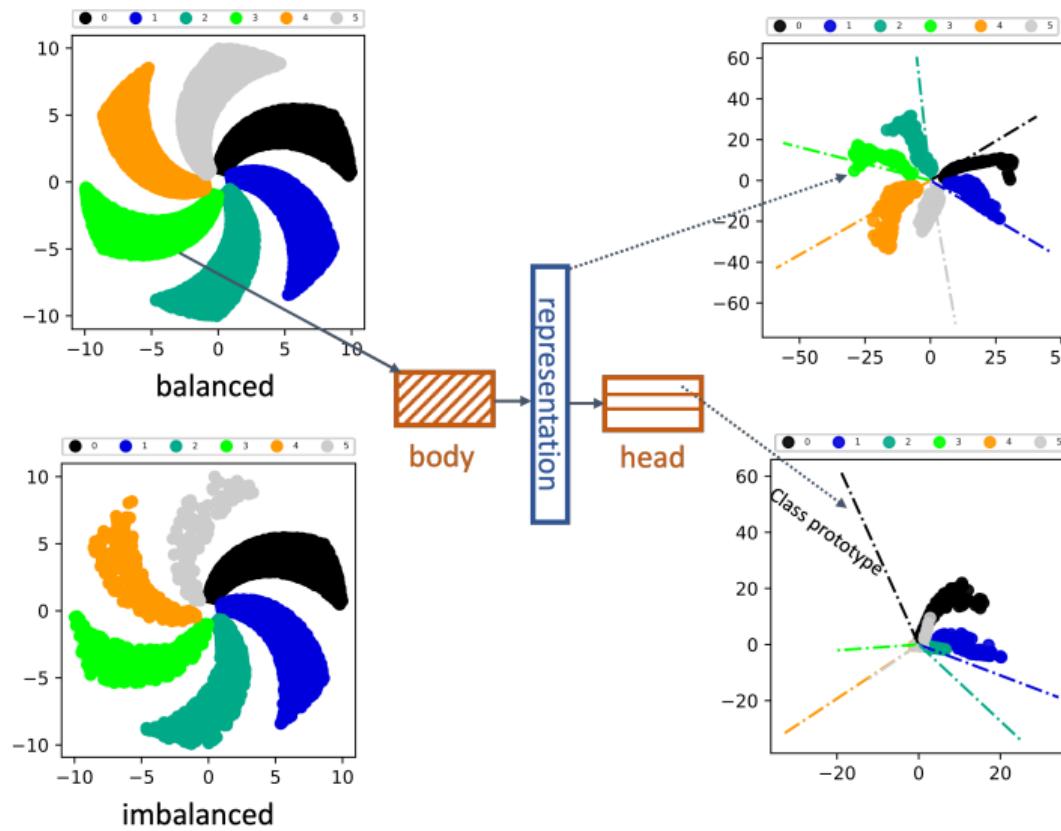
2 Methodology

- A Motivating Example
- Proposed Method

3 Numerical Results

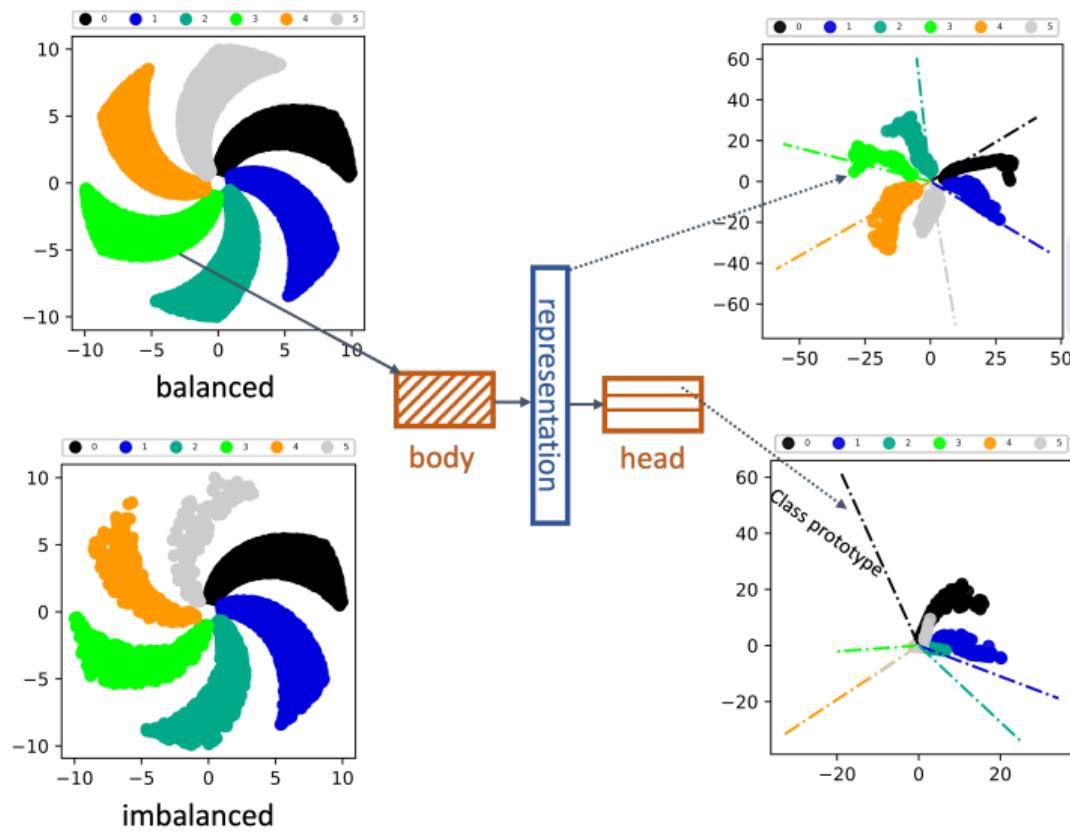
- Class Semantics
- Accuracy

## Centralized Training on A Toy Dataset



Visualization is over the same **balanced** testing dataset for a fully trained MLP.

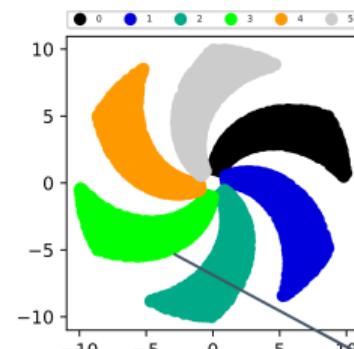
# Centralized Training on A Toy Dataset



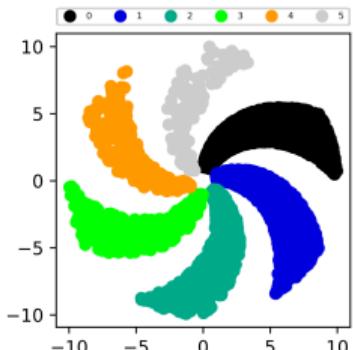
Visualization is over the same **balanced** testing dataset for a fully trained MLP.

- Balanced Dataset
  - uniformly distributed class prototypes
  - separated representations
- Imbalanced Dataset
  - crowded class prototypes
  - overlapped representations

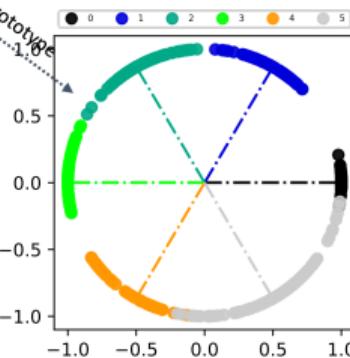
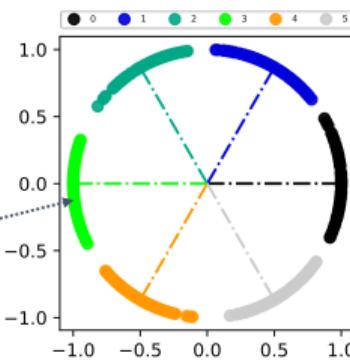
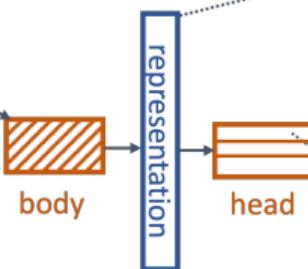
## Centralized Training on A Toy Dataset (Cont'd)



balanced



imbalanced



- no harm on the balanced dataset
- rescue the minority classes

## FL with A Toy Dataset

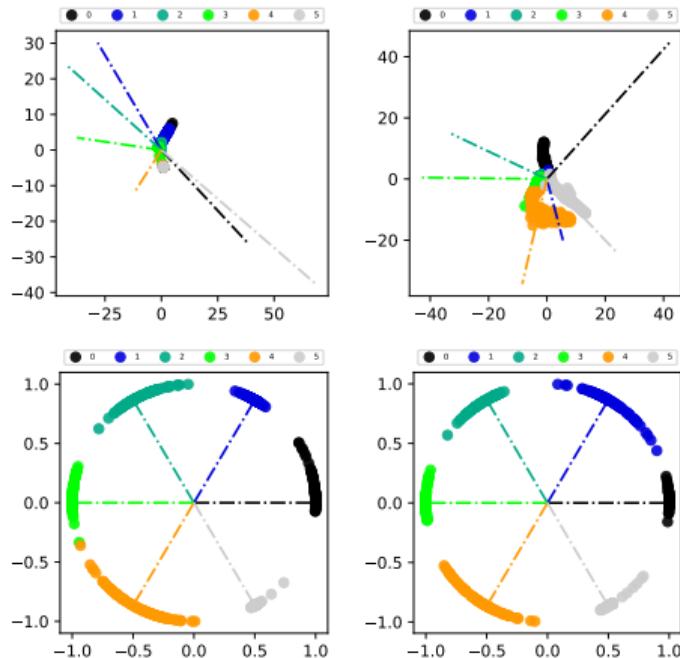


Figure: Top: Visualization of representations and prototypes on two of FedAvg clients. Bottom: Visualization of representations and prototypes (fixed at prior and never updated) on the same two of FedAvg clients.

# Uniformity in Class Prototypes

## Initialization Strategy

$$\begin{aligned} & \max_{\{w_1, \dots, w_{|\mathcal{C}|}, M\}} M^2 \\ \text{s.t. } & \|w_i - w_j\|^2 \geq M^2, \|w_i\|^2 = 1 \text{ for all } i \in [|\mathcal{C}|], i \neq j. \end{aligned} \tag{1}$$

- Can be solved with any constrained optimization solver, e.g., interior point method.
- Need only to be solved once.
- The solution can be approximated with an orthonormal base, which is similar to FedBABU<sup>4</sup>.

<sup>4</sup> Jaehoon Oh, SangMook Kim, and Se-Young Yun. “FedBABU: Toward Enhanced Representation for Federated Image Classification”. In: International Conference on Learning Representations. 2021.

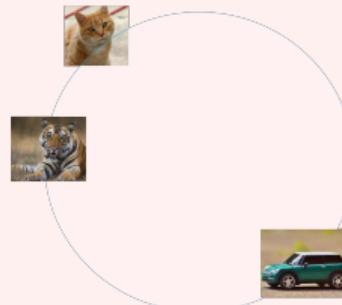
# Uniformity in Class Prototypes

## Initialization Strategy

$$\begin{aligned} & \max_{\{w_1, \dots, w_{|\mathcal{C}|}, M\}} M^2 \\ \text{s.t. } & \|w_i - w_j\|^2 \geq M^2, \|w_i\|^2 = 1 \text{ for all } i \in [|\mathcal{C}|], i \neq j. \end{aligned} \tag{1}$$

- Can be solved with any constrained optimization solver, e.g., interior point method.
- Need only to be solved once.
- The solution can be approximated with an orthonormal base, which is similar to FedBABU<sup>4</sup>.

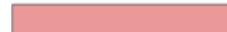
However, this is not enough...



<sup>4</sup> Jaehoon Oh, SangMook Kim, and Se-Young Yun. “FedBABU: Toward Enhanced Representation for Federated Image Classification”. In: International Conference on Learning Representations. 2021.

# Infuse Class Semantics

At the client side:



representations  
of the dog class  
(sample level)

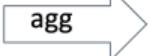


an aggregated local  
representation of the dog class

At the server side:

broadcasting

Step 1



aggregated representations  
of dogs from different clients  
(clients level)

Step 2

$$p \text{ } \textcolor{brown}{\boxed{\text{---}}} + (1-p) \text{ } \textcolor{violet}{\boxed{\text{---}}}$$

the latest knowledge  
about the dog class

past knowledge about  
the dog class

the latest knowledge  
about the dog class

# Algorithm

---

## Algorithm FedNH - Skeleton

---

```

1: Initialization: the body  $\theta$ ; the head  $W \in \mathbb{R}^{|\mathcal{C}| \times d}$ ;
2: for  $t = 0, \dots, R - 1$  communication rounds do
3:   Select a subset of clients  $S^t$ .
4:   ...
5:   for each selected client  $k \in S^t$  in parallel do
6:      $(\theta_k^{t+1}, \mu_k^{t+1}) \leftarrow \text{ClientUpdate}(\theta^t, W^t)$ .    [ $\mu_k^{t+1}$  is aggregated representations of classes of the local training dataset.]
7:   end for
8:   Head Update:  $W_c^{t+1} \leftarrow \rho W_c^t + (1 - \rho) \sum_{k \in S^t} \alpha_k^{t+1} \mu_{k,c}^{t+1}$  for all  $c \in \mathcal{C}$ .
9:   ...
10:  Body Update:  $\theta^{t+1} = \frac{1}{|S^t|} \sum_{k \in S^t} \theta_k^t$ .
11:  ...
12: end for

```

---

# Convergence Result

## Theorem 1 (Informal)

Let the  $k$ th client uniformly at random returns an element from  $\{\theta_k^{t,j}\}$  as the solution, denoted as  $\theta_k^*$ . Further, let  $W^*$  share the same round index as  $\theta_k^*$ . Then for any  $\epsilon > 0$ , set  $\rho \in (\nu_1(\epsilon, M_G, M_f), 1)$  and  $\eta \in (0, \nu_2(\epsilon, L_g, \sigma^2, \rho, M_G, M_f))$ , if  $R > \mathcal{O}(\epsilon^{-1})$ , one gets

$$\mathbb{E} \left[ \|\nabla_{\theta} F_k(\theta_k^*; W^*)\|^2 \right] \leq \epsilon,$$

where  $\nu_1(\epsilon, M_G, M_f)$ ,  $\nu_2(\epsilon, L_g, \sigma^2, \rho, M_G, M_f)$ ,  $M_G$ , and  $M_f$  are some positive constants.

# Outline

1 Problem

2 Methodology

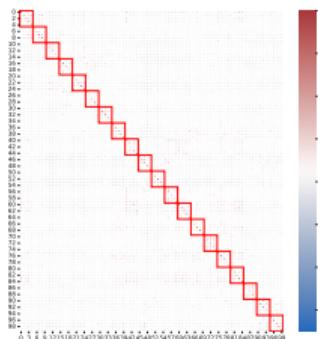
- A Motivating Example
- Proposed Method

3 Numerical Results

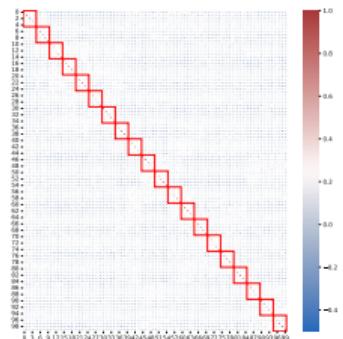
- Class Semantics
- Accuracy

# Learned Class Semantics

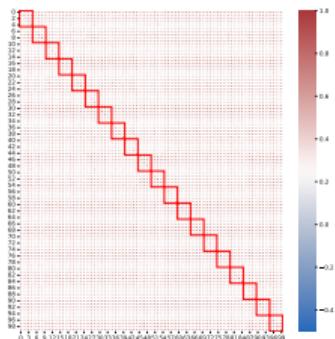
Visualize the pair-wise cosine similarity of class prototypes on Cifar100.



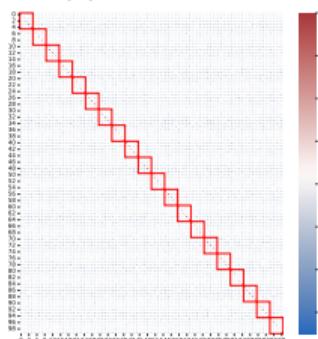
(a) Our method



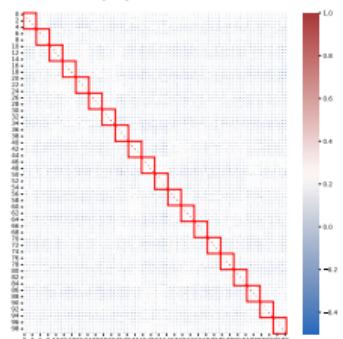
(b) FedAvg



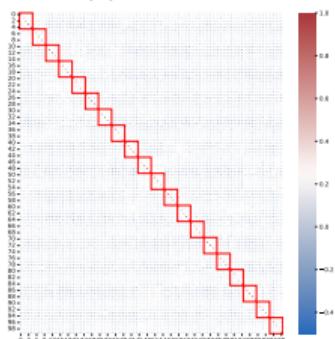
(c) FedProto



(d) FedBABU



(e) FedRod



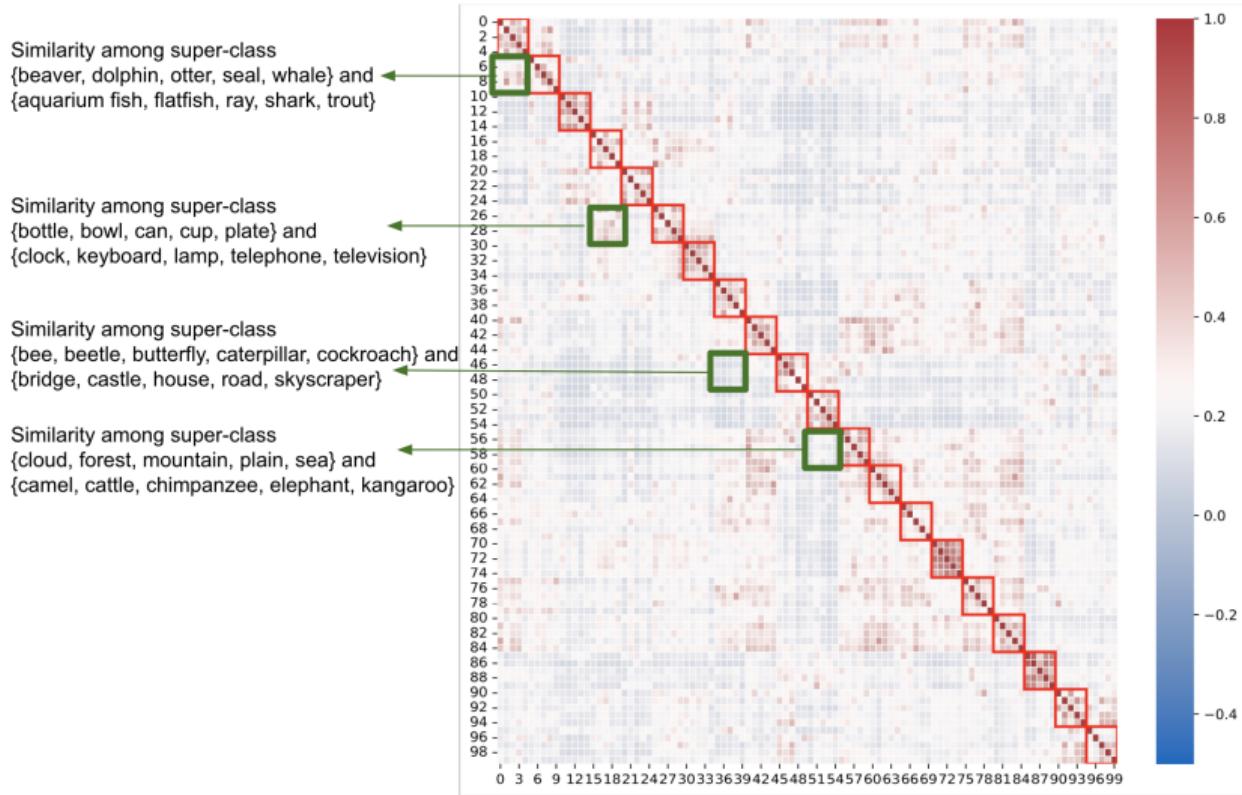
(f) CReFF

(AAAI 2023)

FL Heterogeneity with Class Prototypes

- Cifar100 has 20 super-classes.
- Each block along the diagonal contains 5 sub-classes within one super-class.

## Learned Class Semantics (Cont'd)



# Classification Accuracy

Dataset	Method	Dir(0.3)			Dir(1.0)		
		GM	PM(V)	PM(L)	GM	PM(V)	PM(L)
Cifar100	Local	—	13.63 ± 2.45	30.89 ± 1.82	—	9.44 ± 1.27	16.71 ± 1.03
	FedAvg	35.14 ± 0.48	31.85 ± 1.33	50.77 ± 2.31	36.07 ± 0.41	28.86 ± 1.23	38.35 ± 2.11
	FedPer	15.04 ± 0.06	16.15 ± 2.34	33.10 ± 1.50	14.69 ± 0.03	11.61 ± 2.17	19.08 ± 1.36
	Ditto	35.14 ± 0.48	26.19 ± 1.11	45.91 ± 2.17	36.07 ± 0.41	22.92 ± 1.77	32.81 ± 2.16
	FedRep	5.42 ± 0.03	13.59 ± 2.31	29.45 ± 2.45	6.37 ± 0.04	9.47 ± 2.27	16.07 ± 1.27
	FedProto	—	10.64 ± 1.02	19.11 ± 1.75	—	9.24 ± 1.33	12.61 ± 1.78
	CReFF	22.90 ± 0.30	31.85 ± 1.33	50.77 ± 2.31	22.21 ± 0.15	28.86 ± 1.23	38.35 ± 2.11
	FedBABU	32.41 ± 0.40	28.96 ± 2.16	47.86 ± 1.03	32.34 ± 0.49	25.84 ± 1.44	34.85 ± 1.80
	FedROD	33.83 ± 0.25	28.53 ± 1.27	42.93 ± 1.03	35.20 ± 0.19	27.58 ± 1.98	33.44 ± 1.76
	FedNH	41.34 ± 0.25	38.25 ± 1.23	55.21 ± 2.11	43.19 ± 0.24	36.88 ± 1.15	45.46 ± 2.14

**Metric:** the accuracy of the  $i$ th personalized model is computed as

$$\text{acc}_i = \frac{\sum_{(x_j, y_j) \sim \mathcal{D}^{\text{test}}} \alpha_i(y_j) \mathbf{1}(y_j = \hat{y}_j)}{\sum_{(x_j, y_j) \sim \mathcal{D}^{\text{test}}} \alpha_i(y_j)}.$$

- $\mathcal{D}^{\text{test}}$  is a unified and **balanced** dataset.
- **PM(L):**  $\alpha_i(y)$  to 1 if the class  $y$  appears in  $i$ th client's training dataset and 0 otherwise.
- **PM(V):**  $\alpha_i(y) = \mathbb{P}_i(y = c)$ , the probability of the sample  $y$  is from class  $c$  in the  $i$ th client.
- $\hat{y}$  is the predicted value and  $\mathbf{1}(\cdot)$  is the indicator function.

# Summary

- We proposed FedNH to address the data heterogeneity with class imbalance. FedNH combines the uniformity and semantics of class prototypes to learn high-quality representations for classification.
- Our idea currently only applies to the classification task, and the inductive bias from uniformity and semantics of class prototypes can only be imposed on the head of neural network architecture.

# Thank you and Questions?

Contact: [yud319@lehigh.edu](mailto:yud319@lehigh.edu)

