

Federated Learning with Matched Averaging

Hongyi Wang

(University of Wisconsin-Madison)

Joint work with:



Mikhail Yurochkin
(IBM Research)



Yuekai Sun
(U Michigan)



Dimitris Papailiopoulos
(U Wisconsin-Madison)

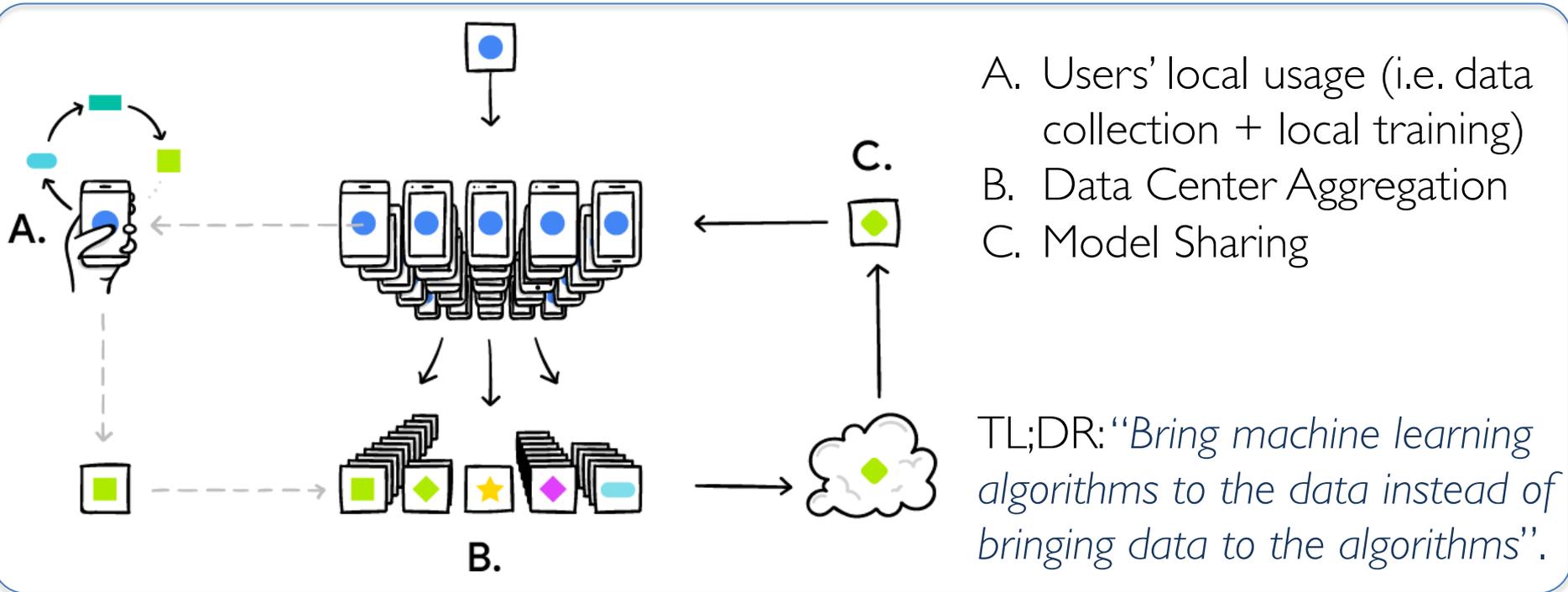


Yasaman Khazaeni
(IBM Research)



ICLR



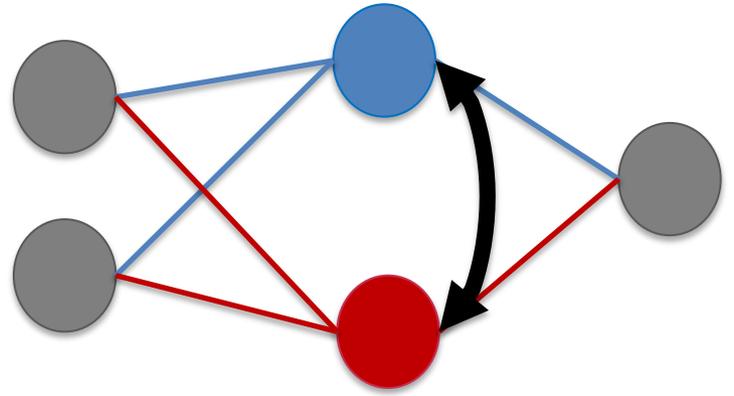
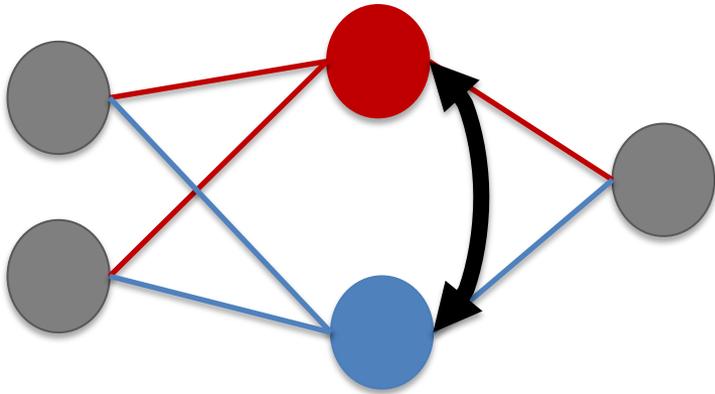


Challenges:

- Local data generated by users **heterogeneously**
- State-of-the-art federated learning approach (e.g. FedAvg) can lead to **poor convergence** performance under some cases
- Federated learning hasn't been well-explored on **large-scale applications**



Permutation Invariance of Neural Network



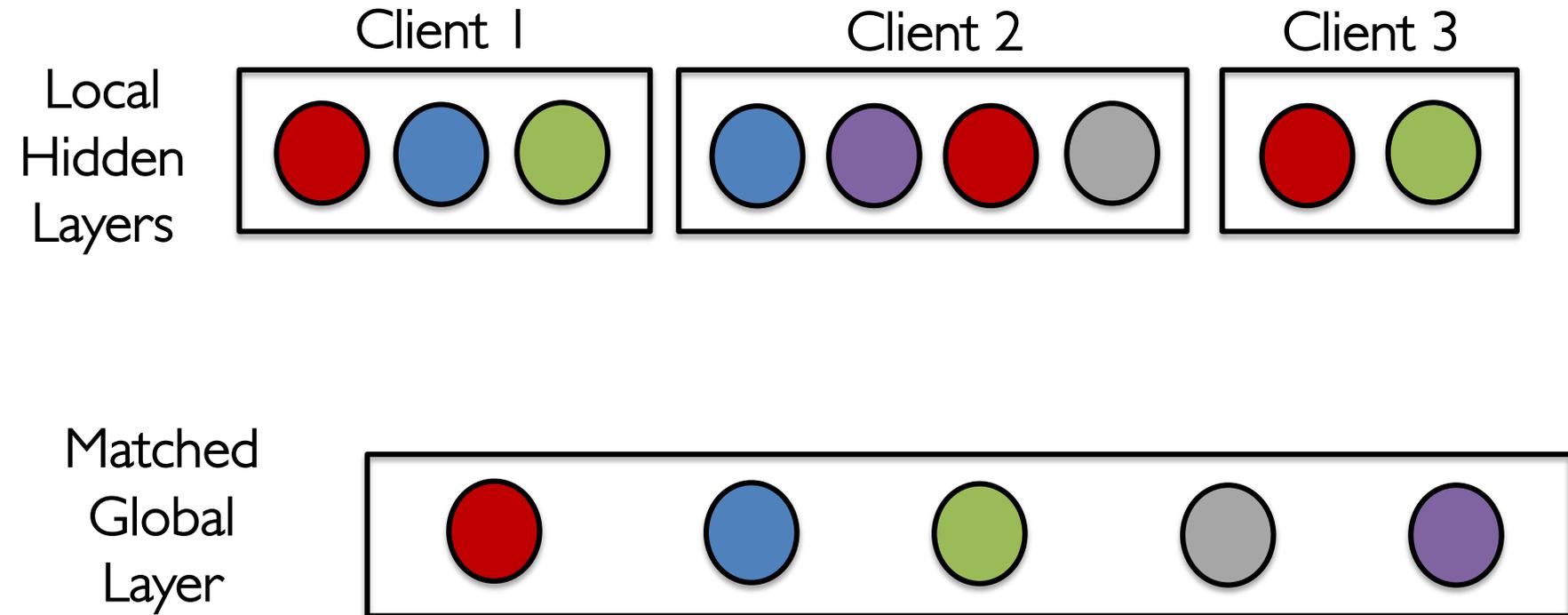
$$\begin{bmatrix} 1 & 2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 & 2 \end{bmatrix} \cdot \begin{bmatrix} 2 & 1 \\ 4 & 3 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ 3 \end{bmatrix}$$

TL;DR: “Due to **permutation invariance** of hidden neurons, direct element-wise averaging may not be a good idea. We should instead find a way to **align** hidden neurons before averaging them.”

Neural Matching Formulation



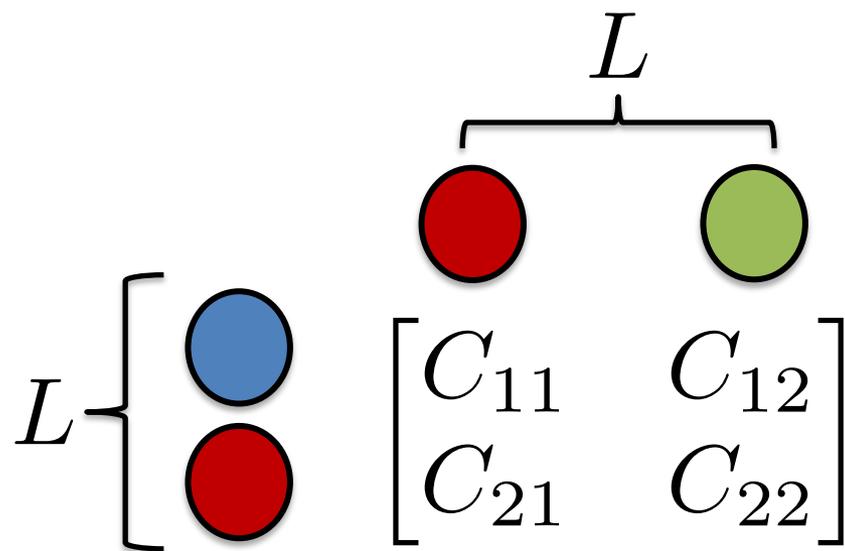
ICLR



Neural Matching Formulation



ICLR

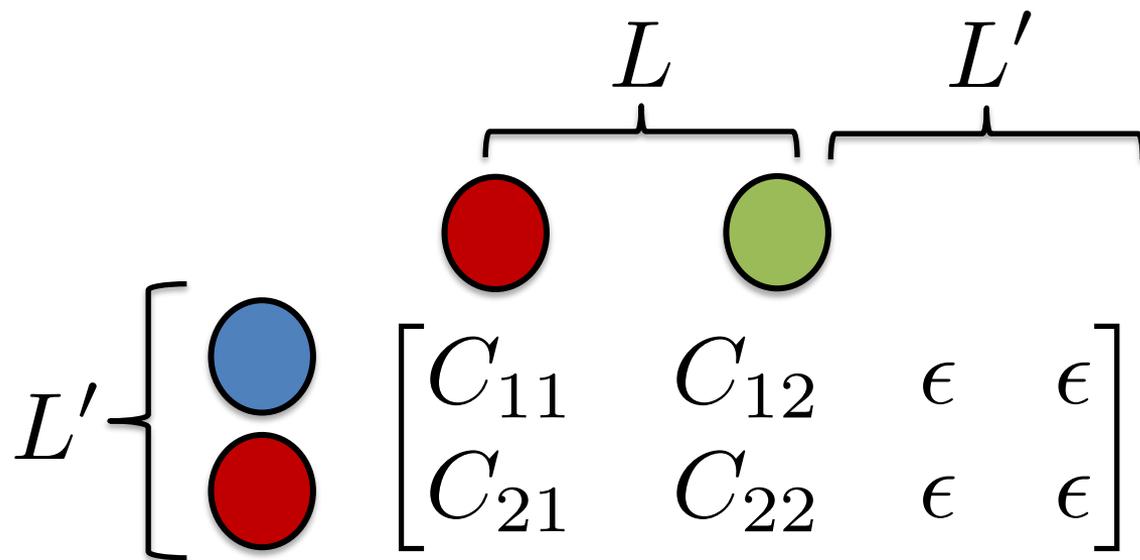


$$\begin{aligned} \min_B \quad & \sum_{l=1}^L \sum_{i=1}^L B_{li} C_{li} \\ \text{s.t.} \quad & \sum_i B_{li} = 1, \forall l, \\ & \sum_l B_{li} = 1, \forall i. \end{aligned}$$

Neural Matching Formulation



ICLR



$$\min_B \sum_{l=1}^{L'} \sum_{i=1}^{L+L'} B_{li} C_{li}$$

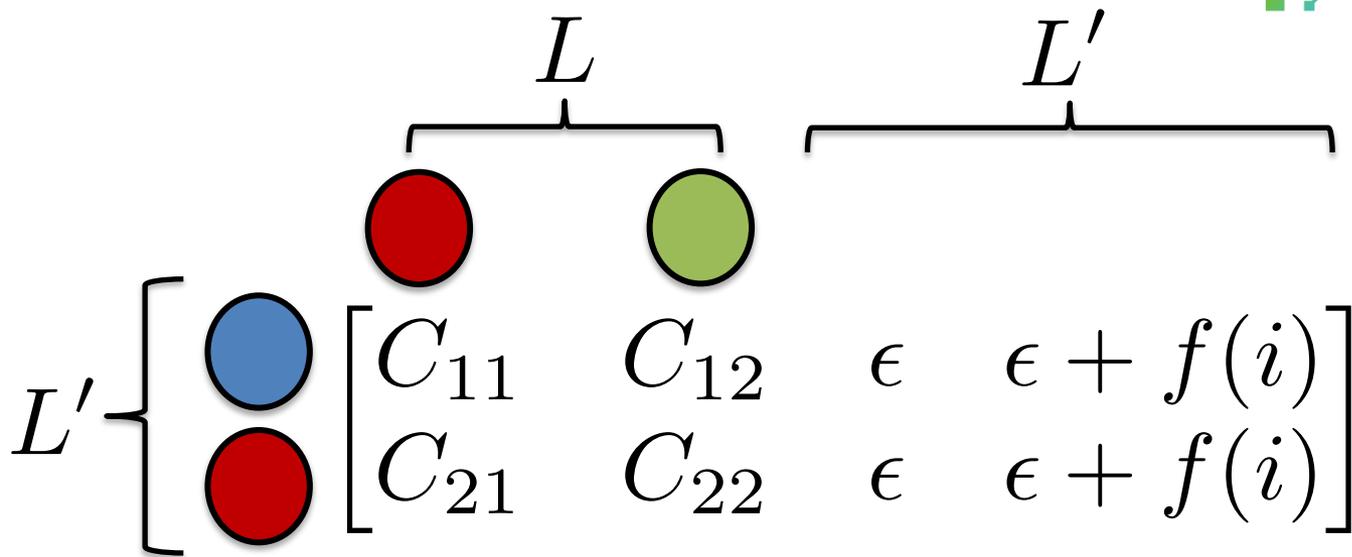
$$\text{s.t.} \sum_i B_{li} = 1, \forall l,$$

$$\sum_l B_{li} = 1, \forall i; B_{li} \in \{0, 1\}, \forall i, l.$$

Neural Matching Formulation



ICLR



$$\min_B \sum_{l=1}^{L'} \sum_{i=1}^{L+L'} B_{li} C_{li}$$

$$C_{li} = \begin{cases} c(\theta_l, \theta_i), & i \leq L \\ \epsilon, & L < i \leq L + L' \end{cases}$$

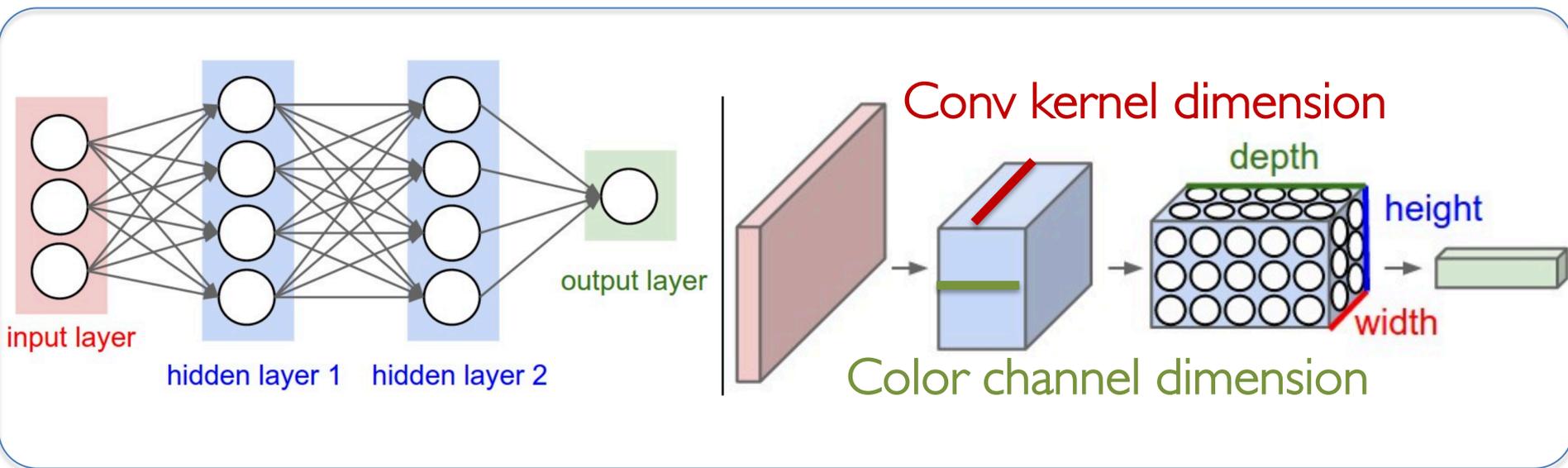
s.t. $\sum_i B_{li} = 1, \forall l,$ $f(i)$ increasing in i

$$\sum_l B_{li} = 1, \forall i; B_{li} \in \{0, 1\}, \forall i, l.$$

Permutation Invariance Class in CNN



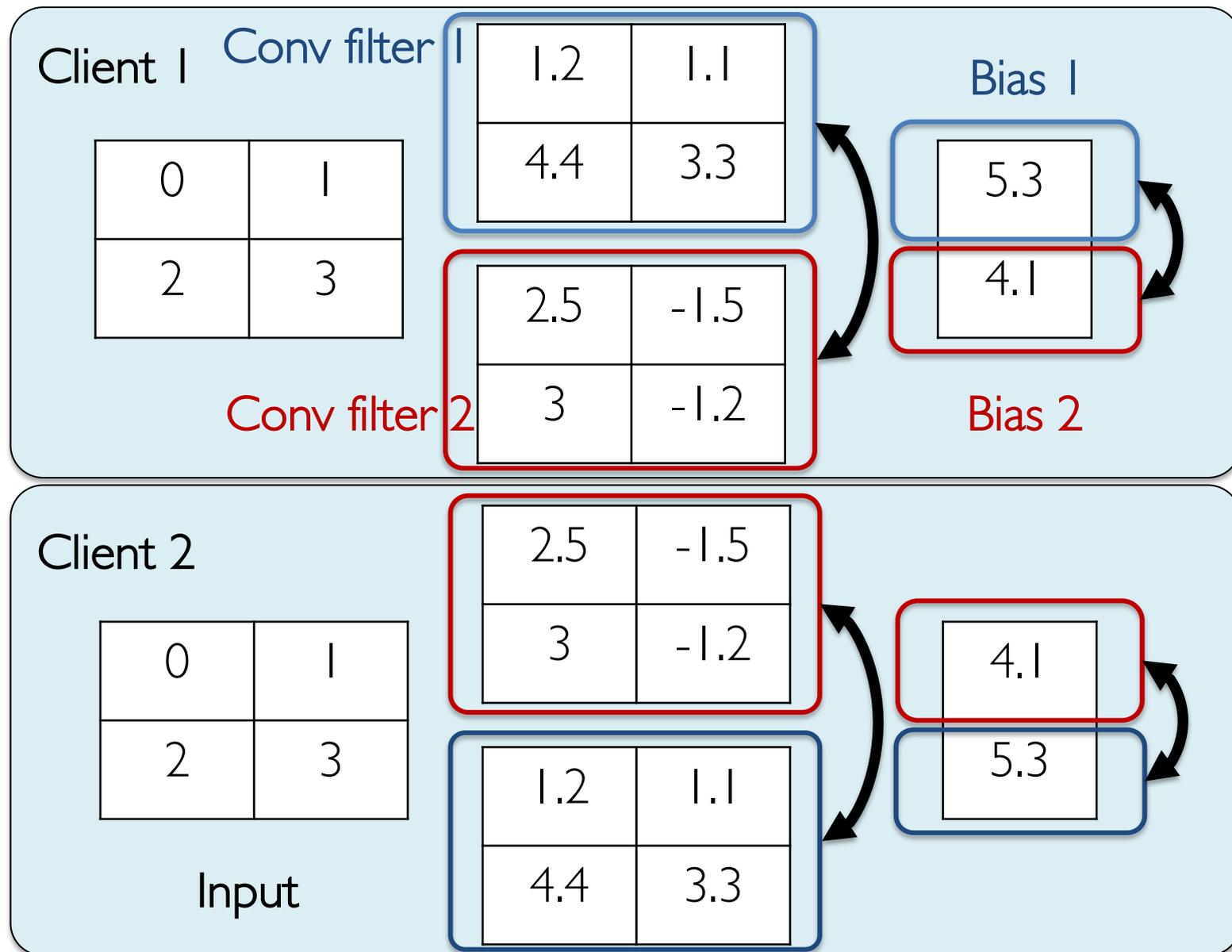
ICLR



- The **convolutional kernels** in a convolution layer are permutation invariant.
- To match convolutional layers, we conduct matching over the **kernel dimension**.

Permutation Invariance

Class in CNN



Permutation Invariance

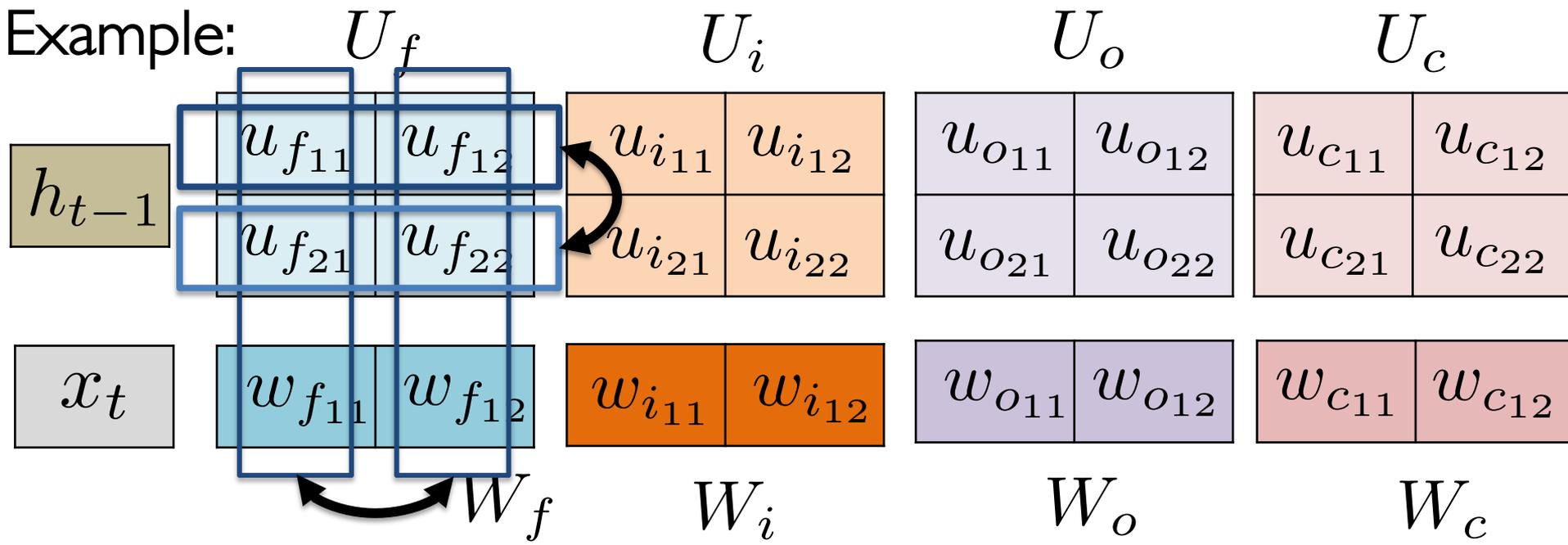
Class in LSTM



$$h_t = \sigma(h_{t-1}\Pi^\top U\Pi + x_t W\Pi)$$

Matching U requires $\|\Pi^\top U_j \Pi - U_{j'}\|_2^2$ which is NP-hard!

Fortunately, we can directly use the same permutation appear in $W\Pi$ to permute U .

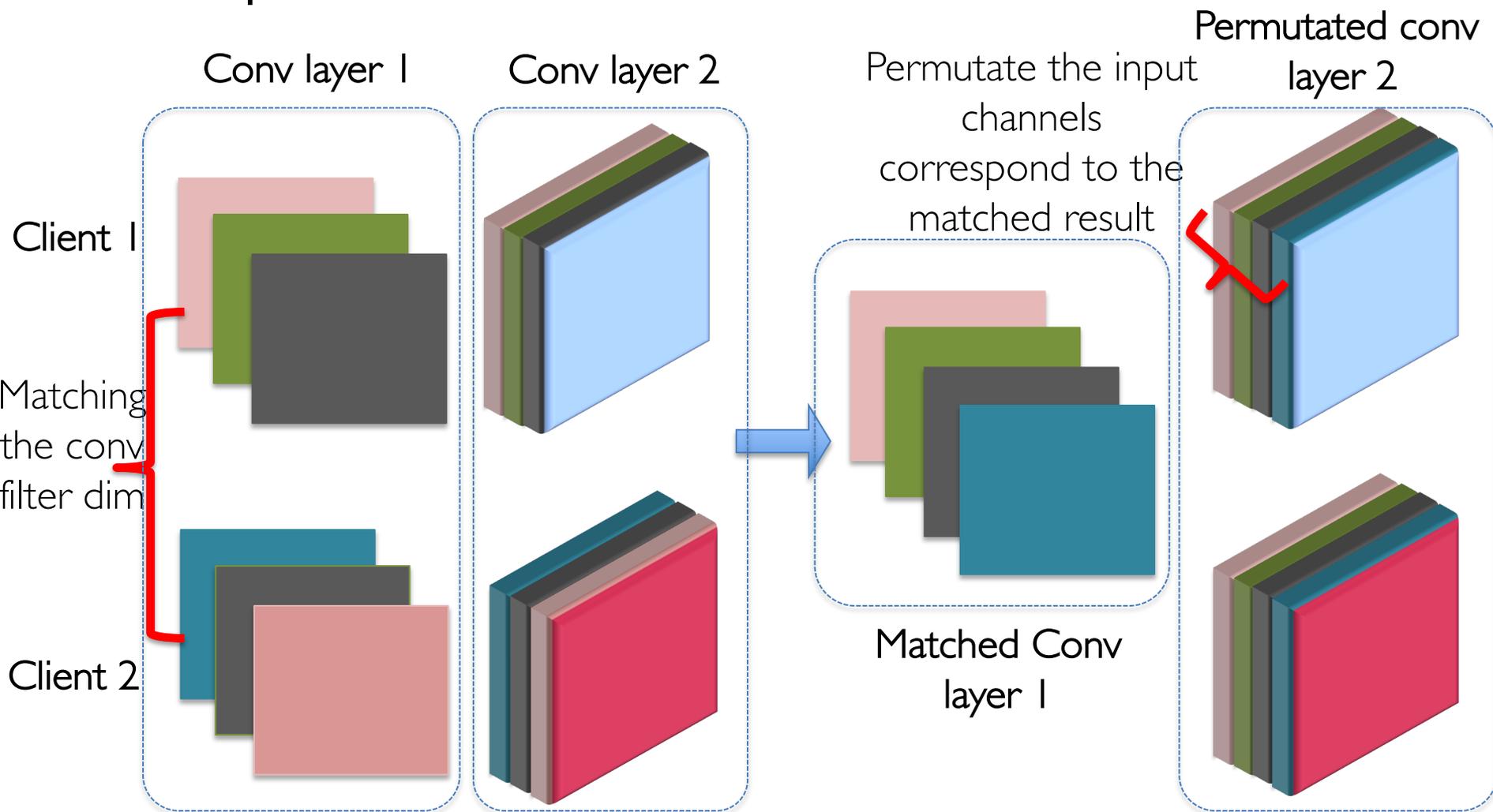


One-shot Matching



ICLR

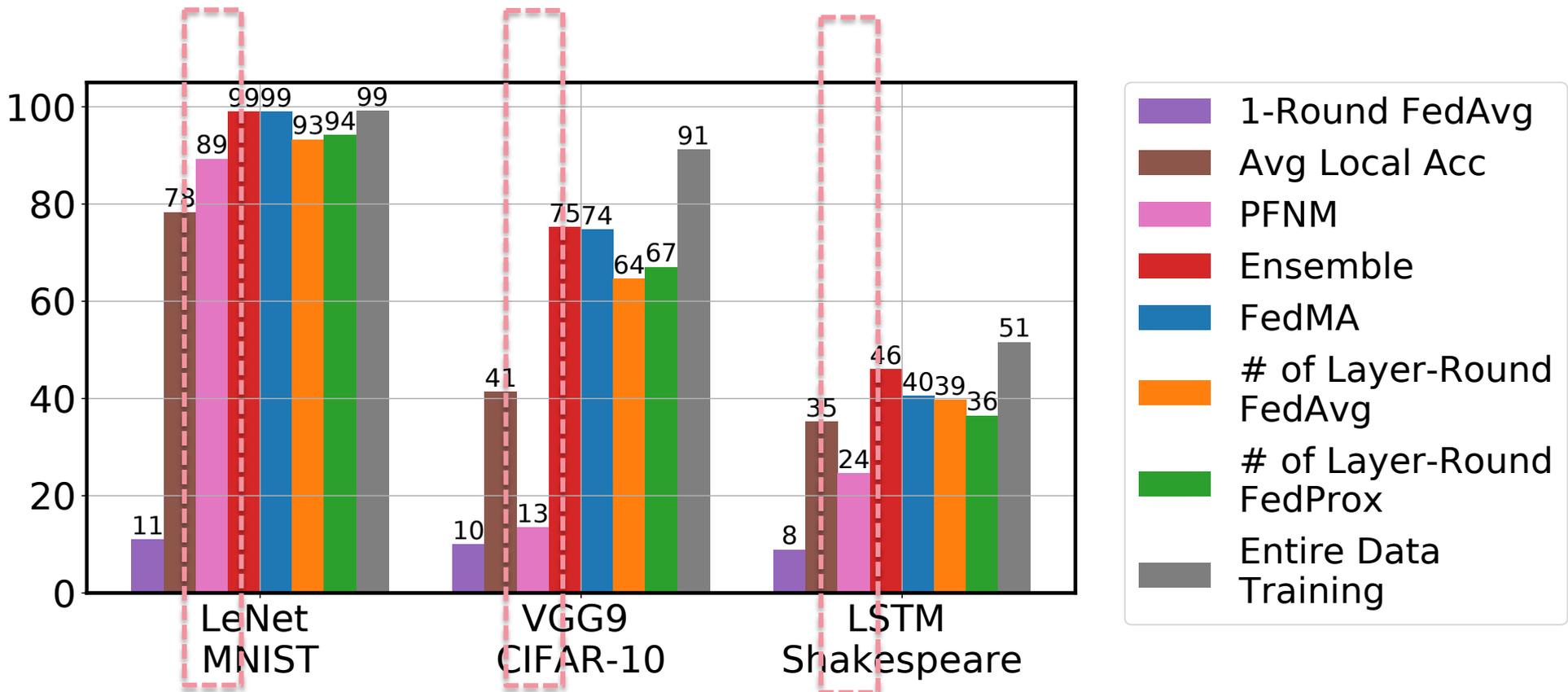
An Example on CNN:



First Result on One-shot Matching



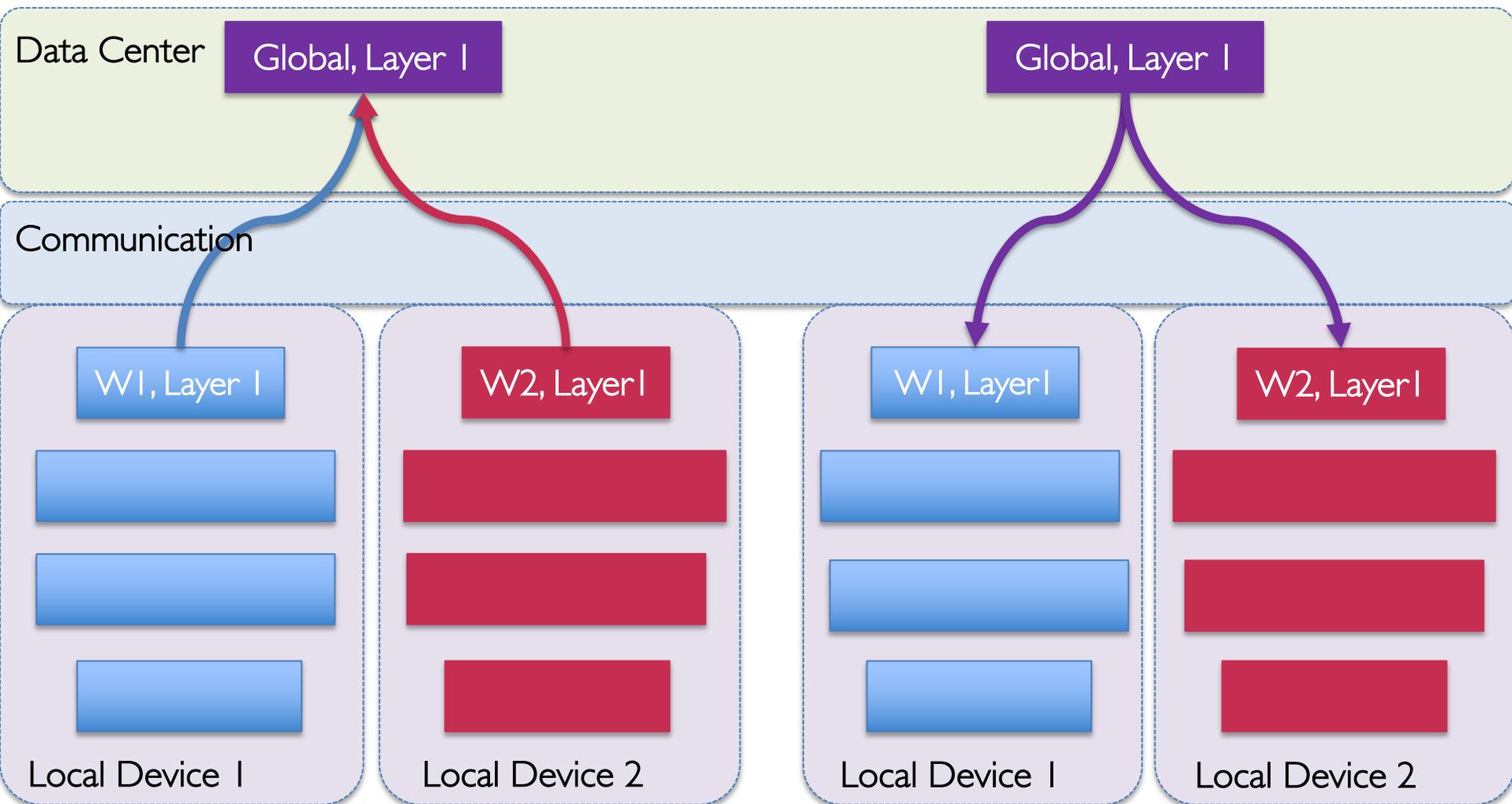
ICLR

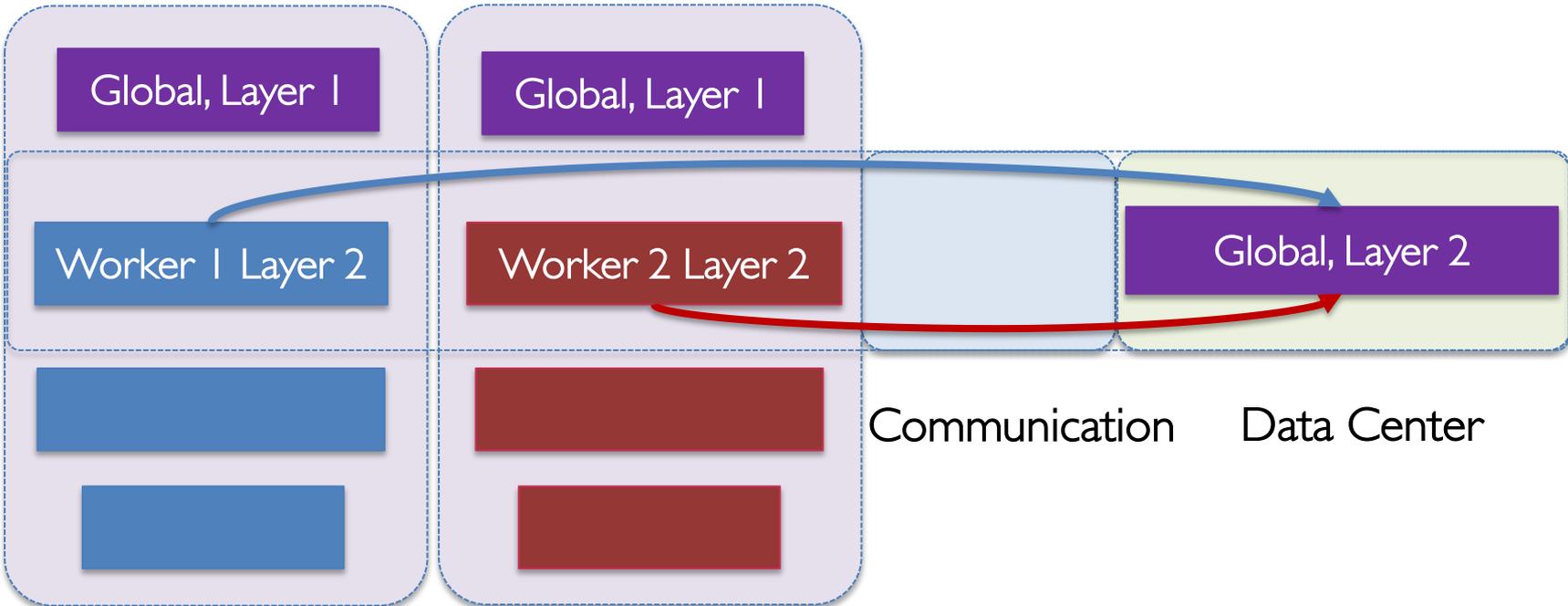
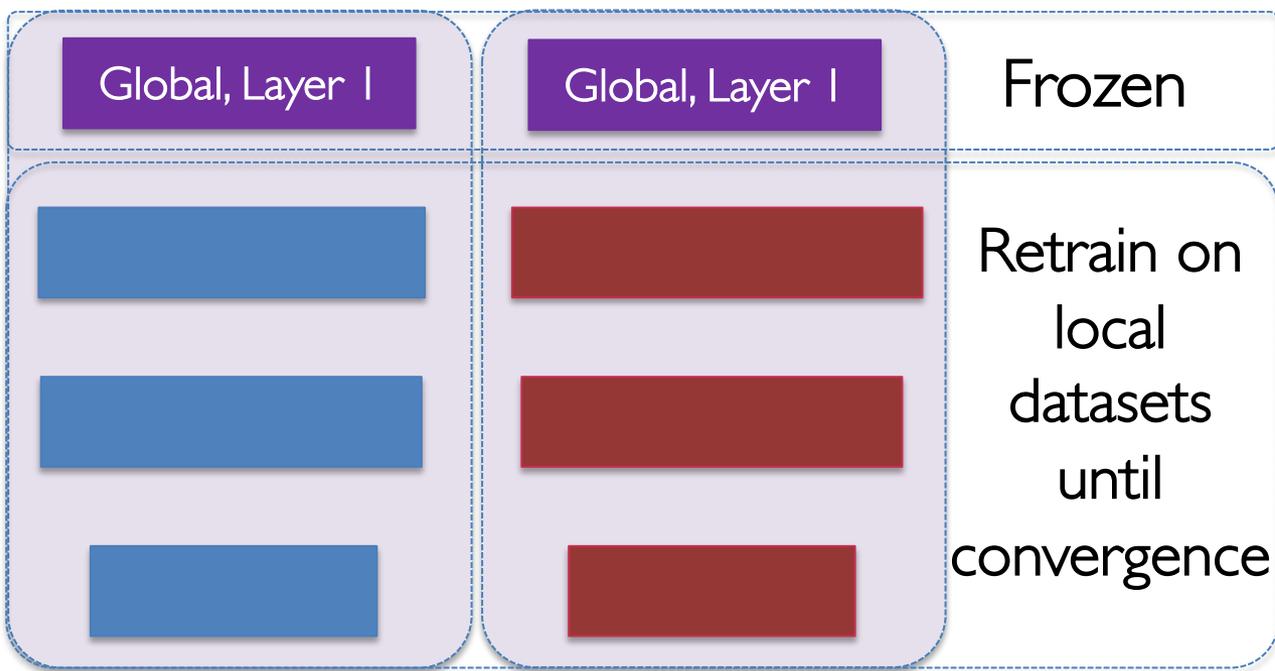


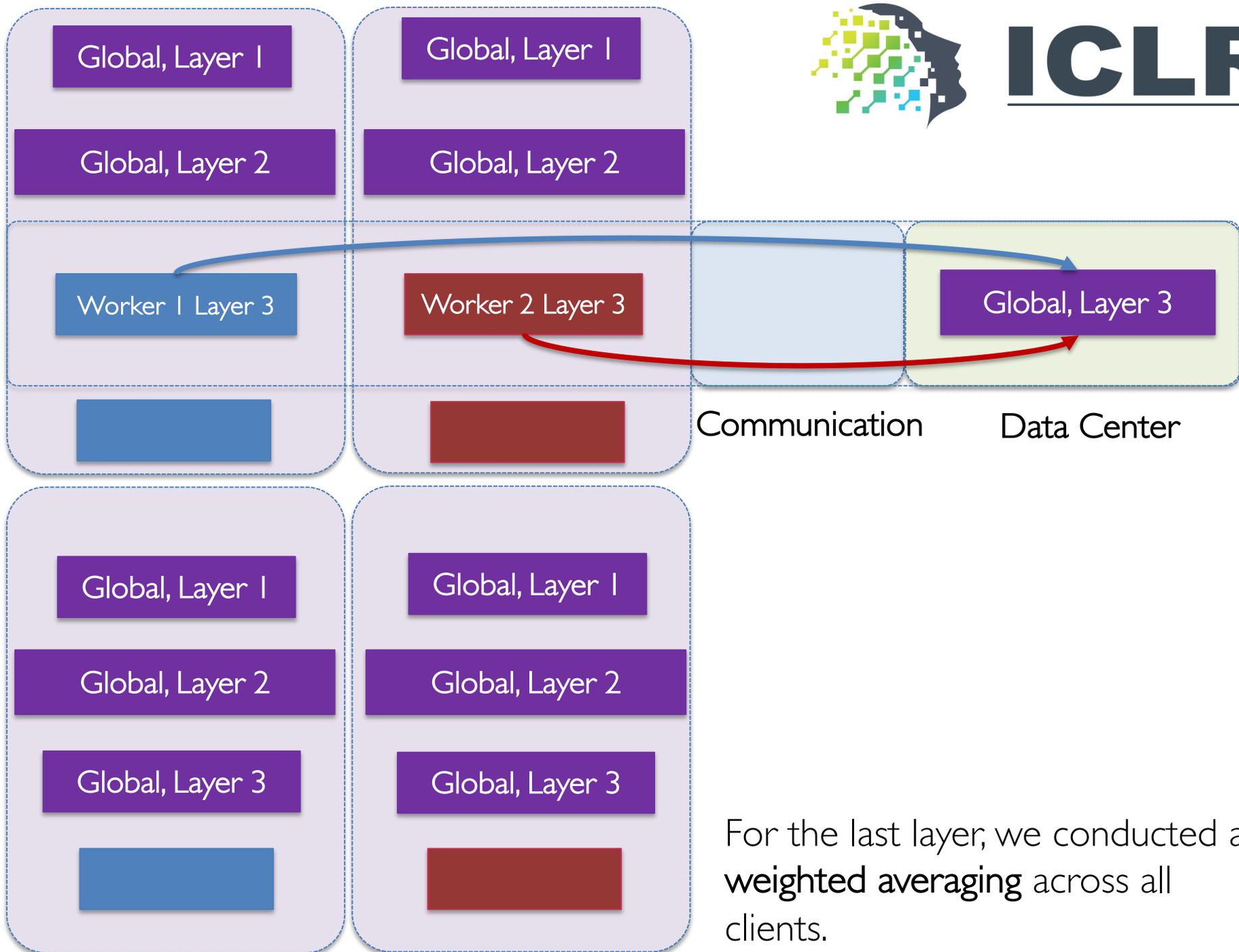
FedMA: Layer-wise matched averaging



ICLR





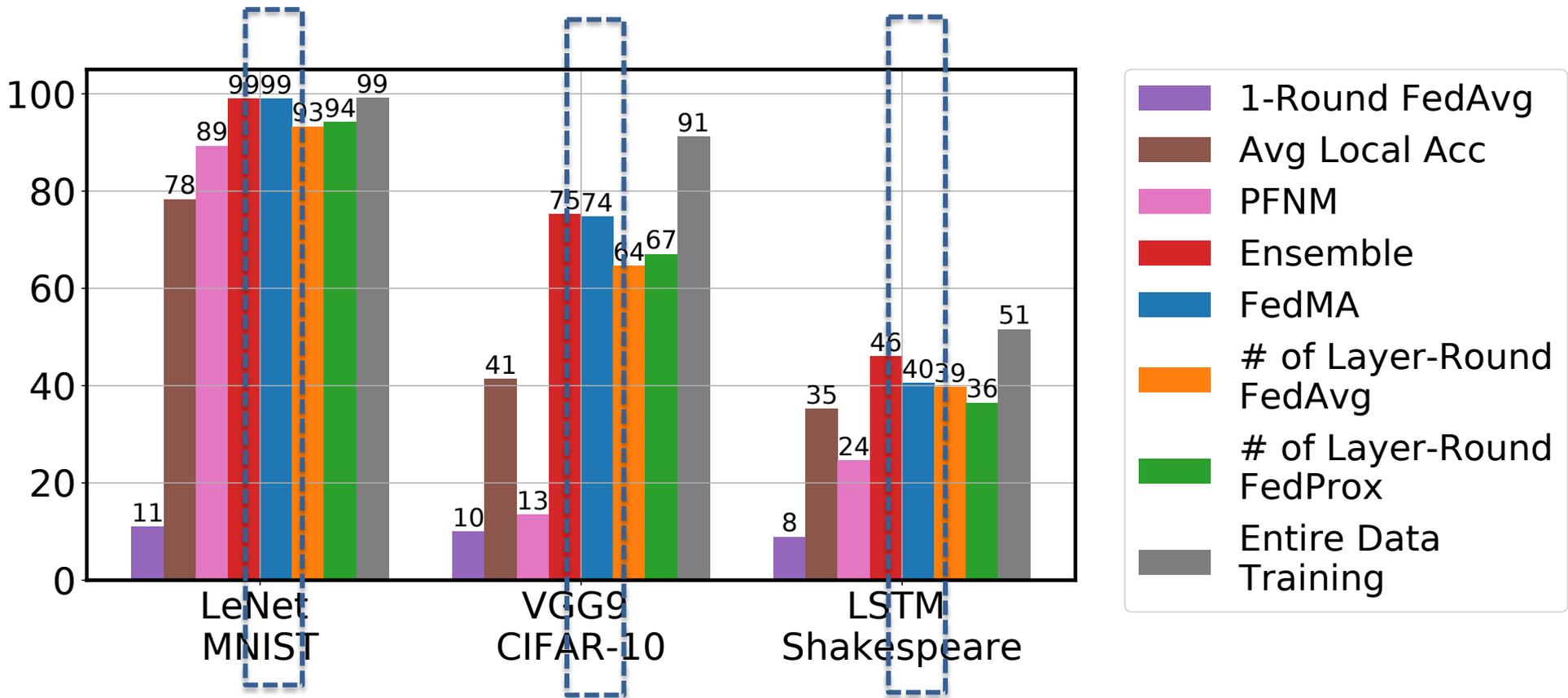


For the last layer, we conducted a **weighted averaging** across all clients.

First Result on One-shot Matching



ICLR

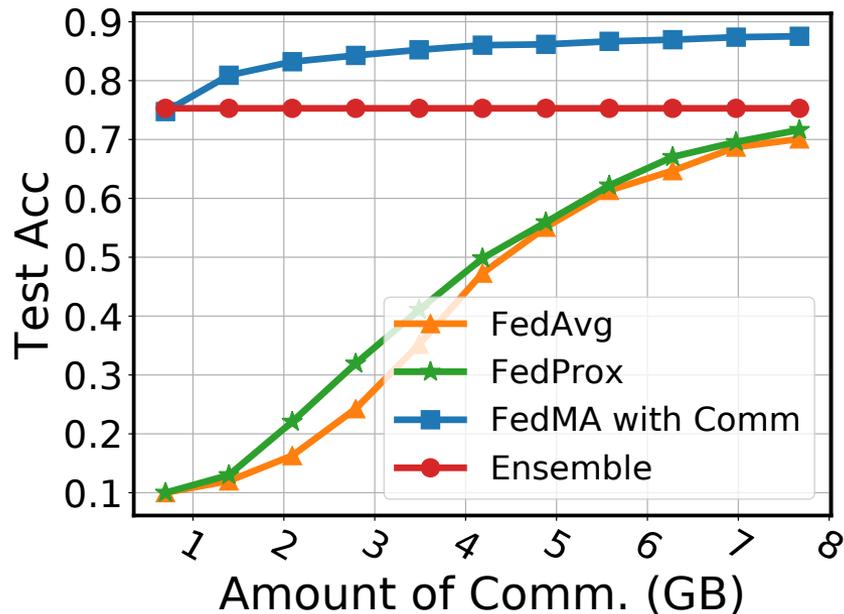


FedMA: matching with comm.

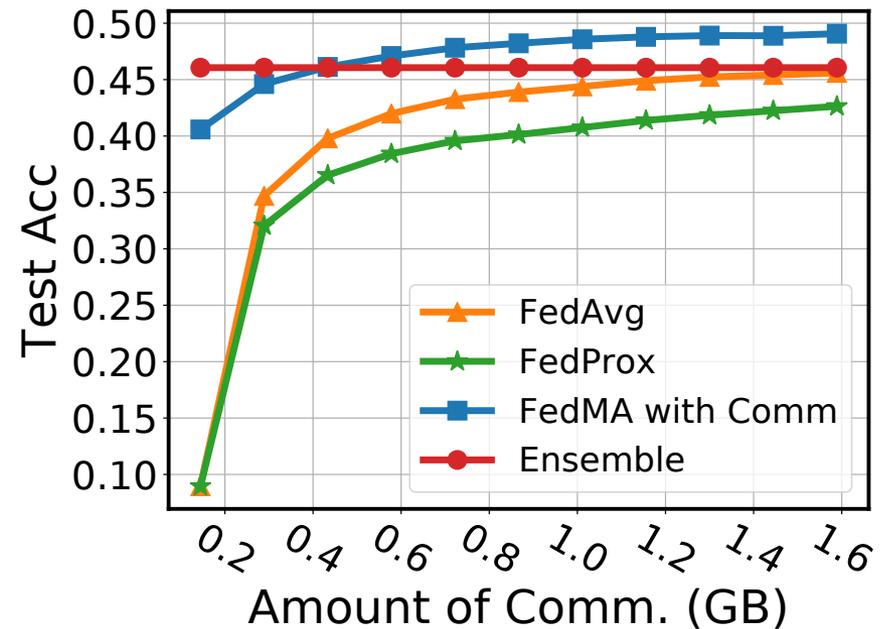


ICLR

- After the first round of layer-wise matching, we broadcast the matched model back to local workers.
- We slice the matched global model to recover the original local model size.
- We then repeat the layer-wise matching process.



VGG-9 trained on CIFAR-10



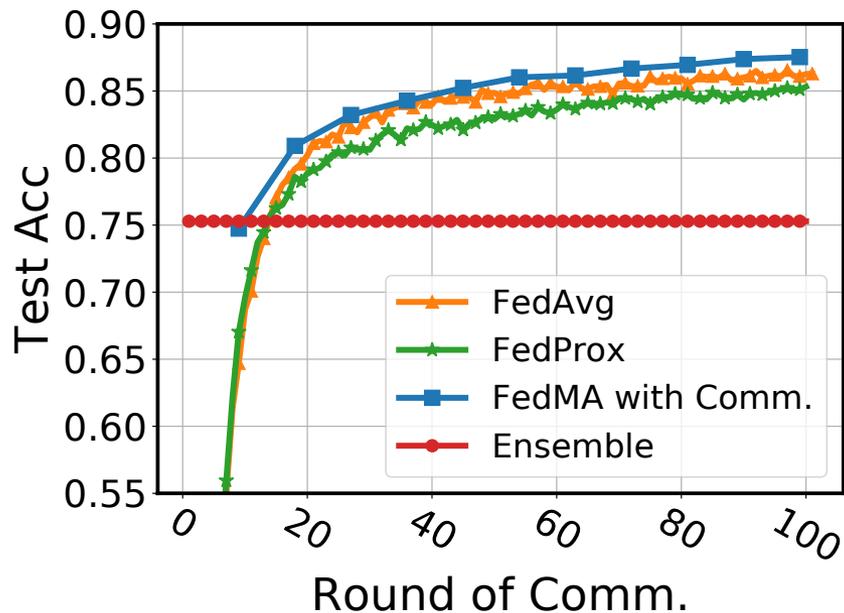
LSTM trained on Shakespeare

16 local workers; heterogeneous setup; convergence w.r.t. message size

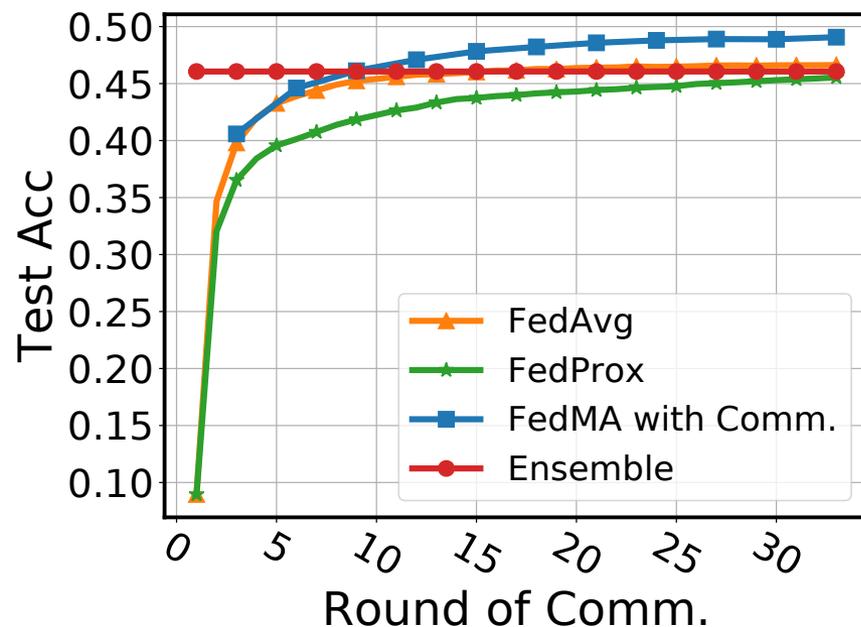
FedMA matching with comm.



- After the first round of layer-wise matching, we broadcast the matched model back to local workers.
- We slice the matched global model to recover the original local model size.
- We then repeat the layer-wise matching process.



VGG-9 trained on CIFAR-10



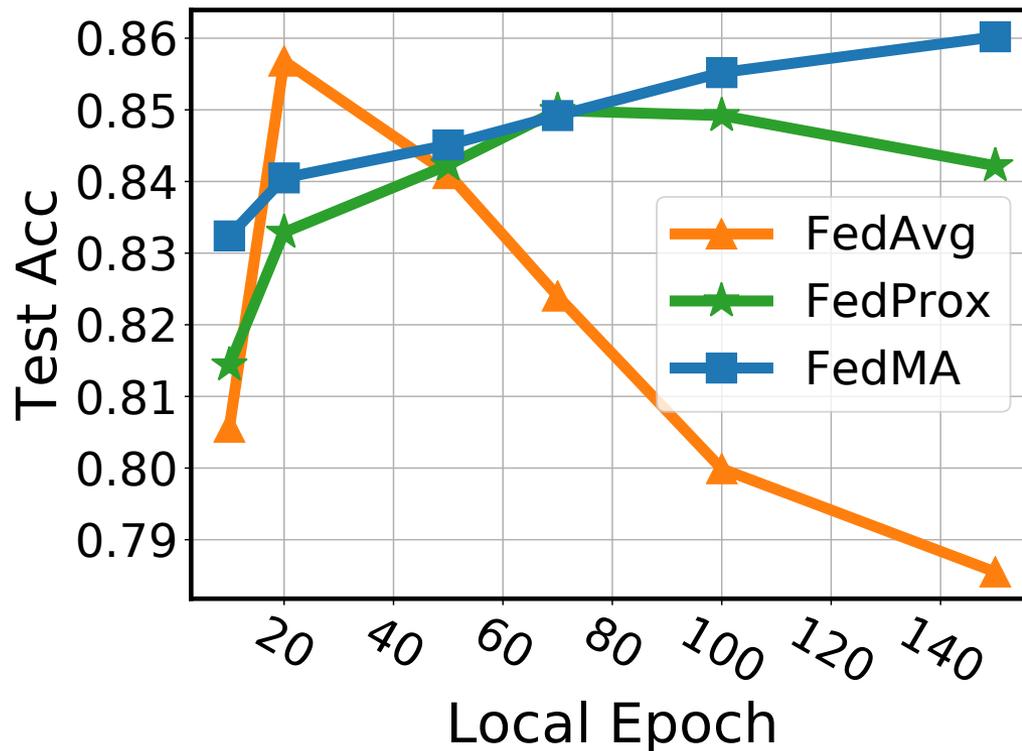
LSTM trained on Shakespeare

16 local workers; heterogeneous setup; convergence w.r.t. communication round

Effect of local training period



- Local training epochs (denoted by E) is a hyper-parameter shared among FedMA, FedAvg, and FedProx.
- It has been studied that a “too-long” local training period leads to bad global model in FedAvg.
- We study the effect of E on FedMA.



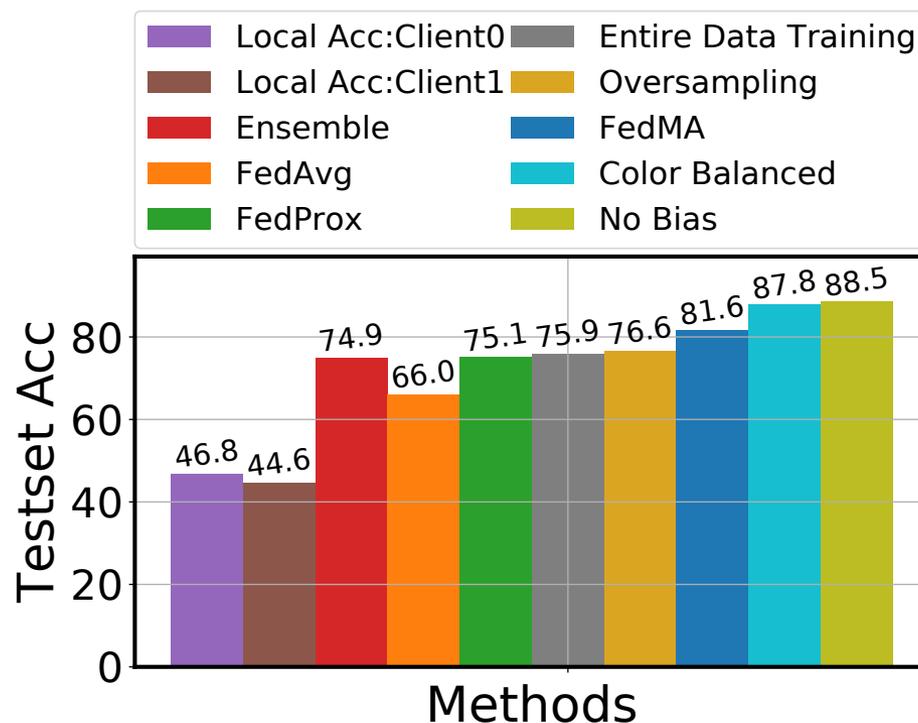
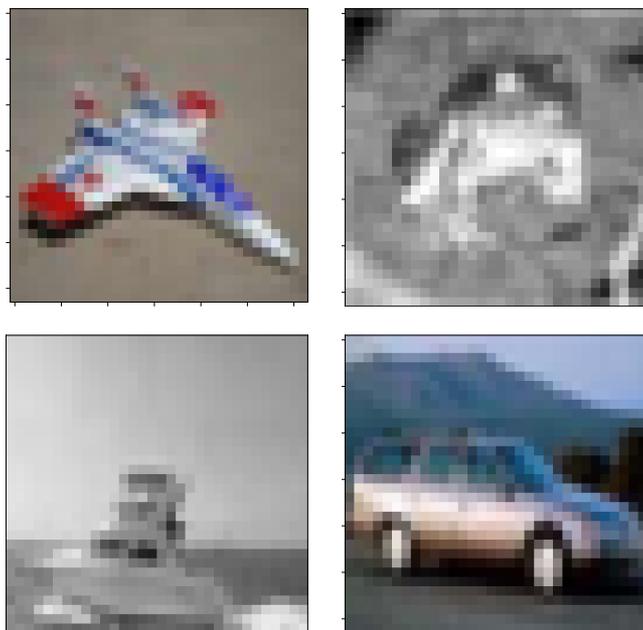
16 local clients;
heterogeneous data partition

Handling Data Bias



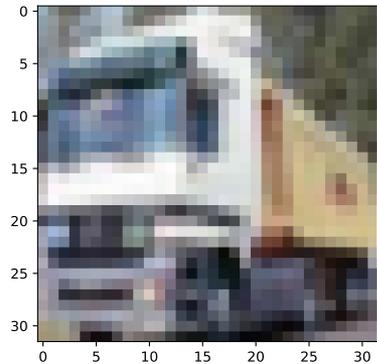
ICLR

- Randomly sample 5 classes where 95% images are grayscale and the rest 5% are color in the training set.
- In the test set, the fraction of grayscale and color images are balanced.
- We split all grayscale images to one client and all color image to the other client.

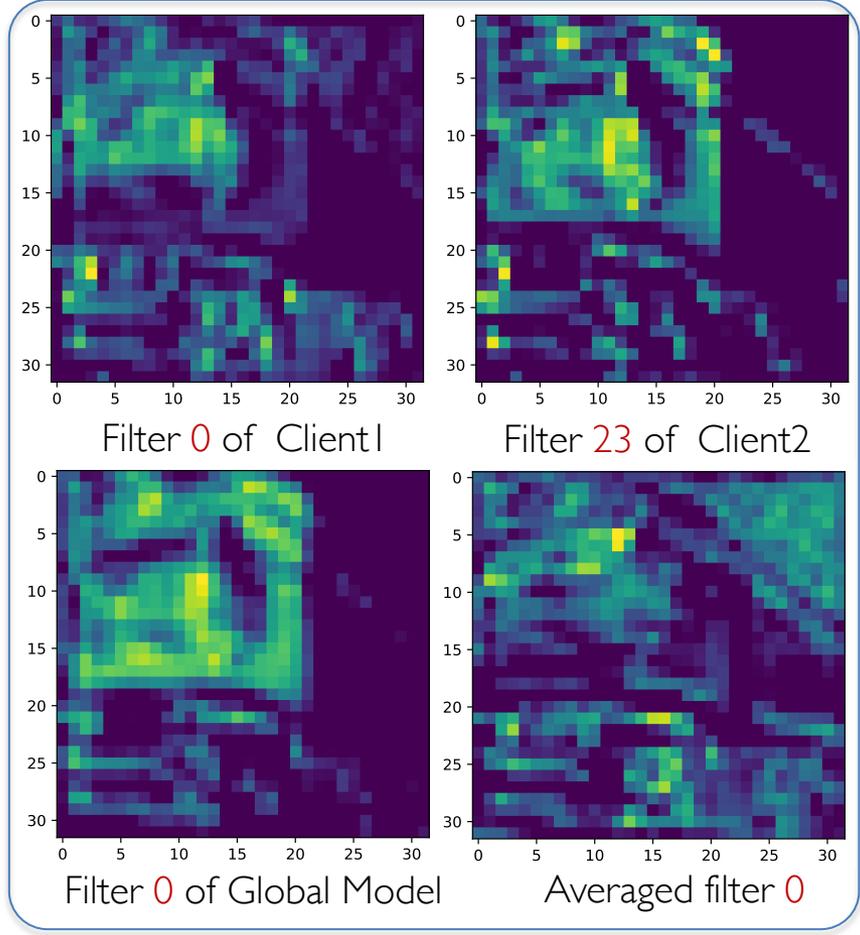
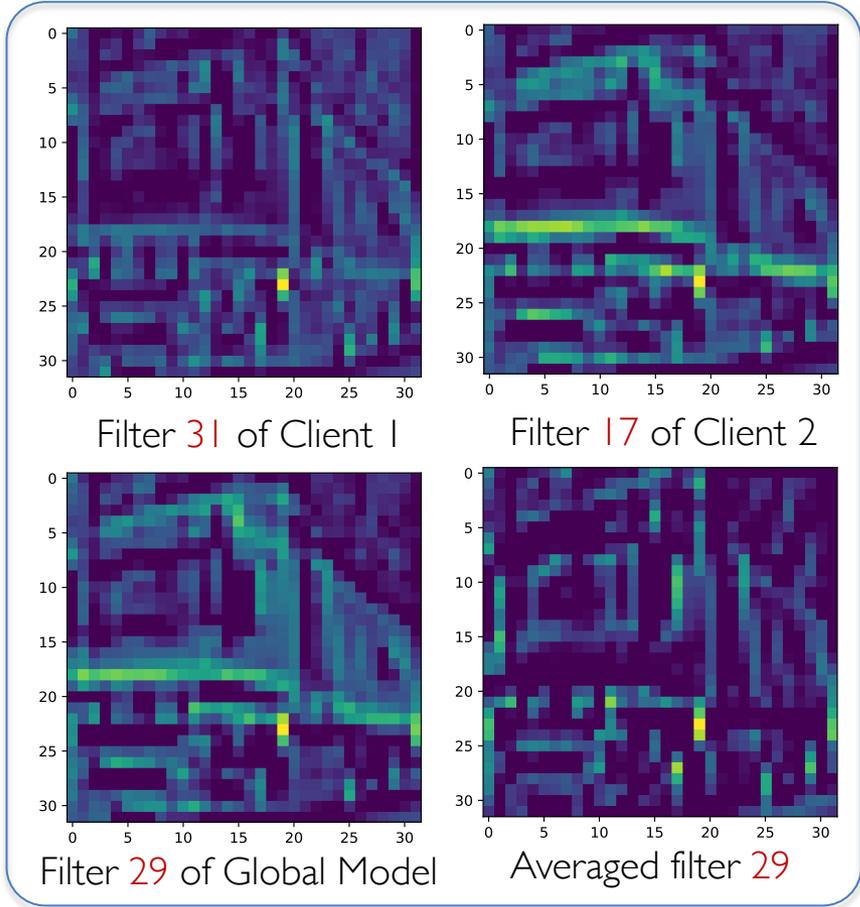


Inspired by "identifying and understanding deep learning phenomena" workshop on ICML19.

Explainability:



Example: Datapoint #2 of CIFAR-10

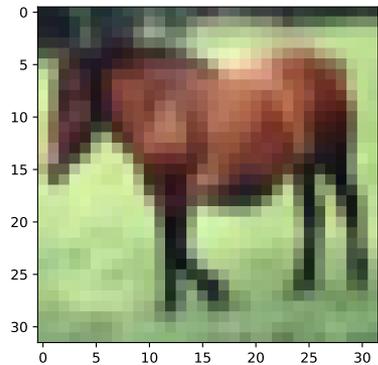


Since plotting the **conv filters** is not easy, we visualize the **representations** coming out of the matched filters

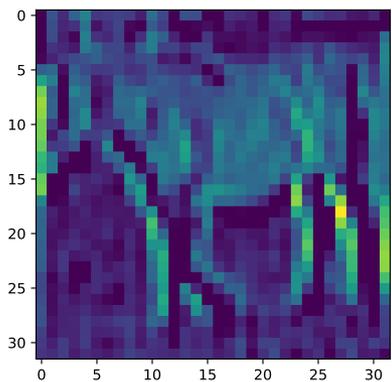


ICLR

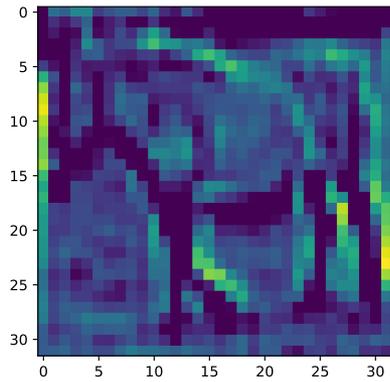
Matching outcome:



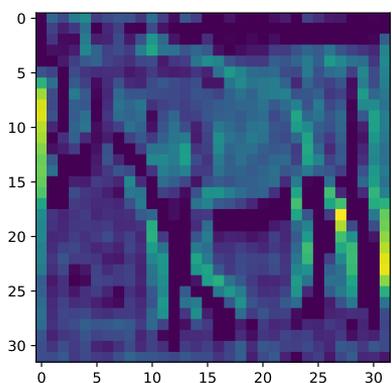
Example: Datapoint #8 of CIFAR-10



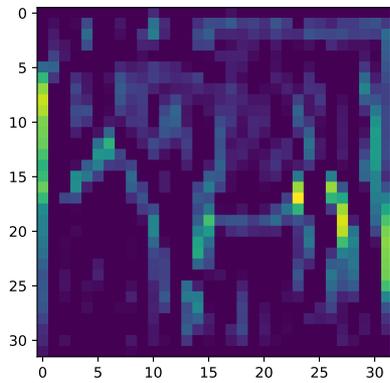
Filter 31 of Client 1



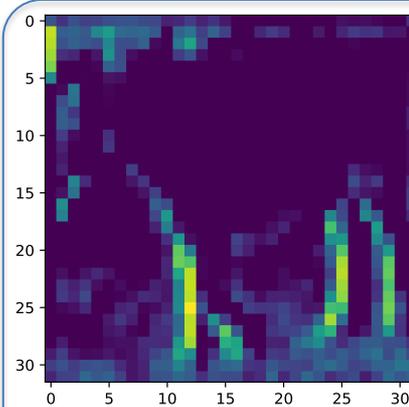
Filter 17 of Client 2



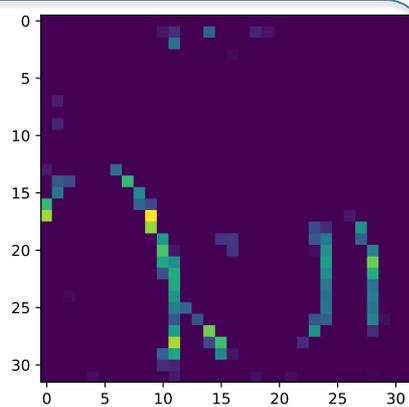
Filter 29 of Global Model



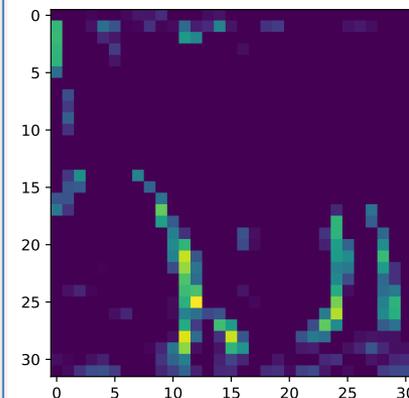
Averaged filter 29



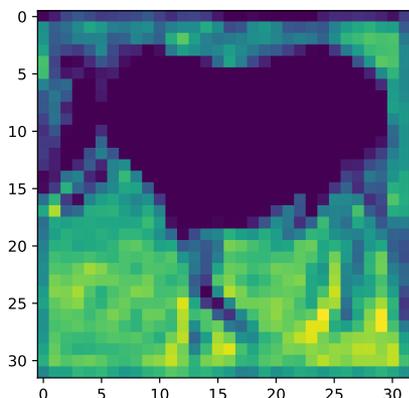
Filter 0 of Client 1



Filter 23 of Client 2



Filter 0 of Global Model



Averaged filter 0

Averaging the conv filters can break the “feature maps”.



ICLR



ICLR

Code available at:

<https://github.com/IBM/FedMA>

Thank you!