

# Personalized Federated Learning through Local Memorization

Othmane Marfoq<sup>1, 2</sup>

Giovanni Neglia<sup>1</sup>

Laetitia Kamani<sup>2</sup>

Richard Vidal<sup>2</sup>

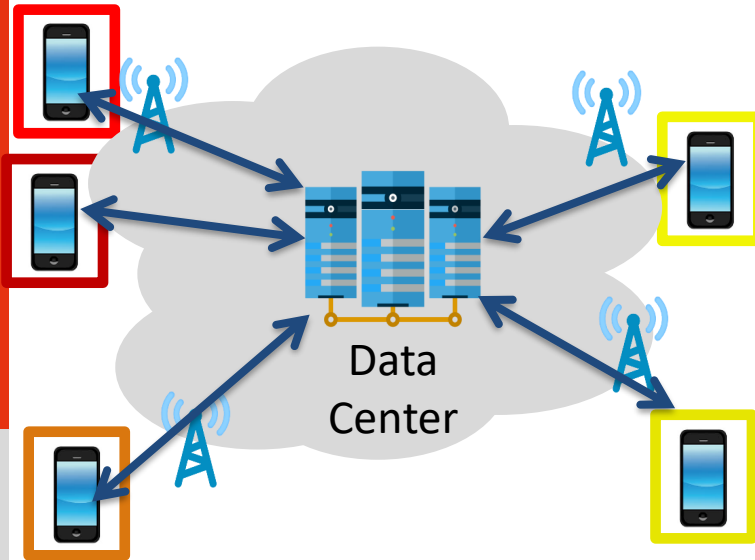
<sup>1</sup>Inria, Université Côte d'Azur and <sup>2</sup>Accenture Labs

*Inria*



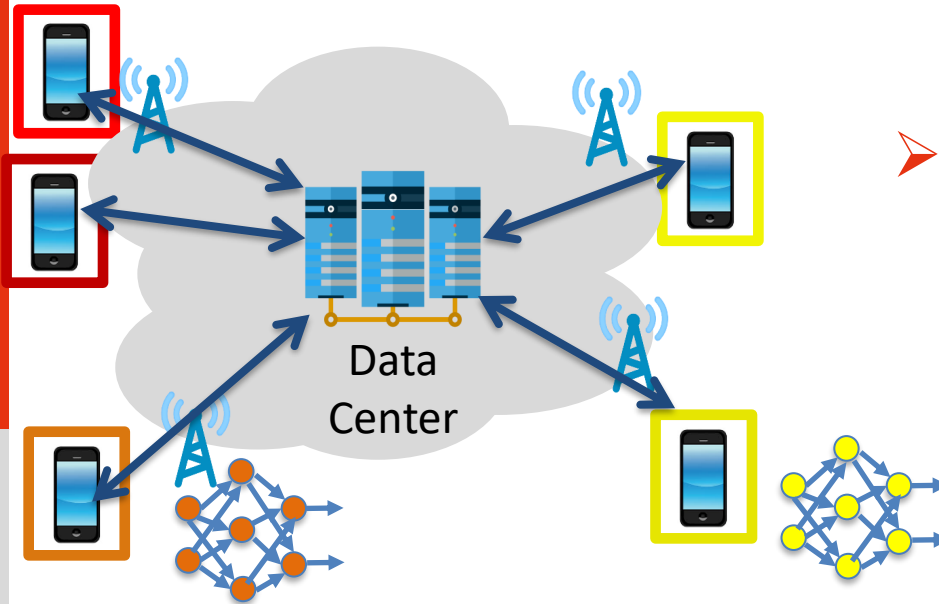
**Accenture Labs**

# Personalization



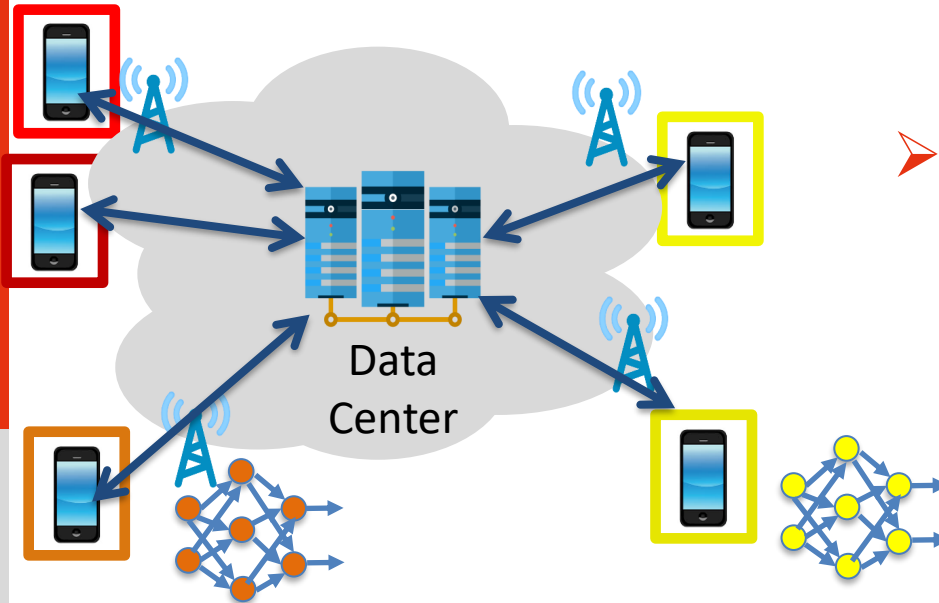
- Why a single model if local dataset come from a different distribution? Statistical heterogeneity

# Personalization

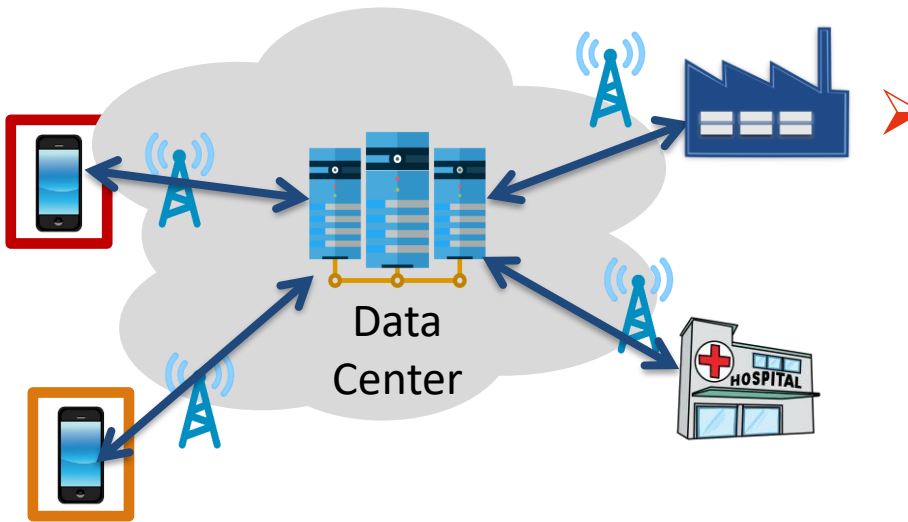


- Why a single model if local dataset come from a different distribution? Statistical heterogeneity

# Personalization

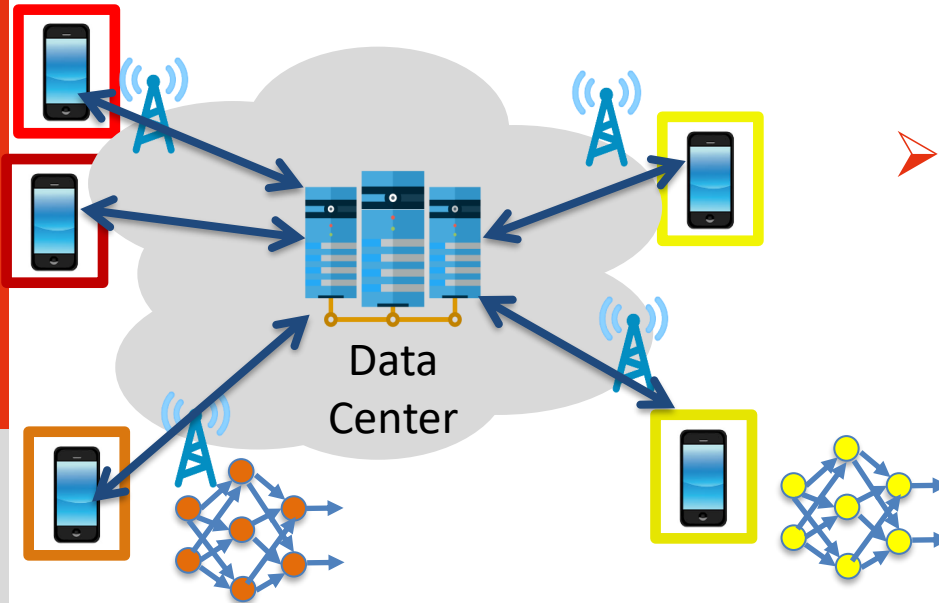


➤ Why a single model if local dataset come from a different distribution? Statistical heterogeneity

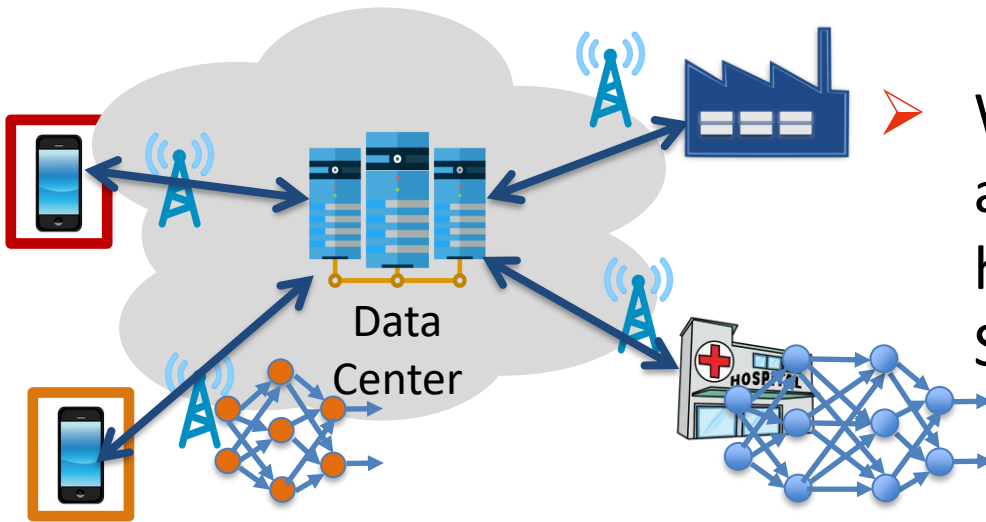


➤ Why the same model architecture when clients have different capabilities? System heterogeneity

# Personalization

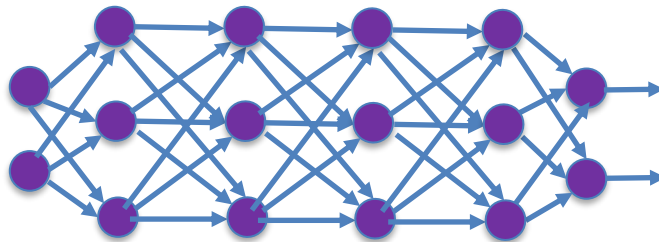


➤ Why a single model if local dataset come from a different distribution? Statistical heterogeneity



➤ Why the same model architecture when clients have different capabilities? System heterogeneity

# Personalized FL through Local Memorization: background

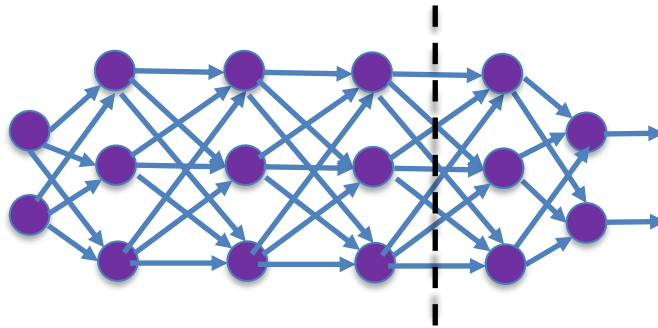


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

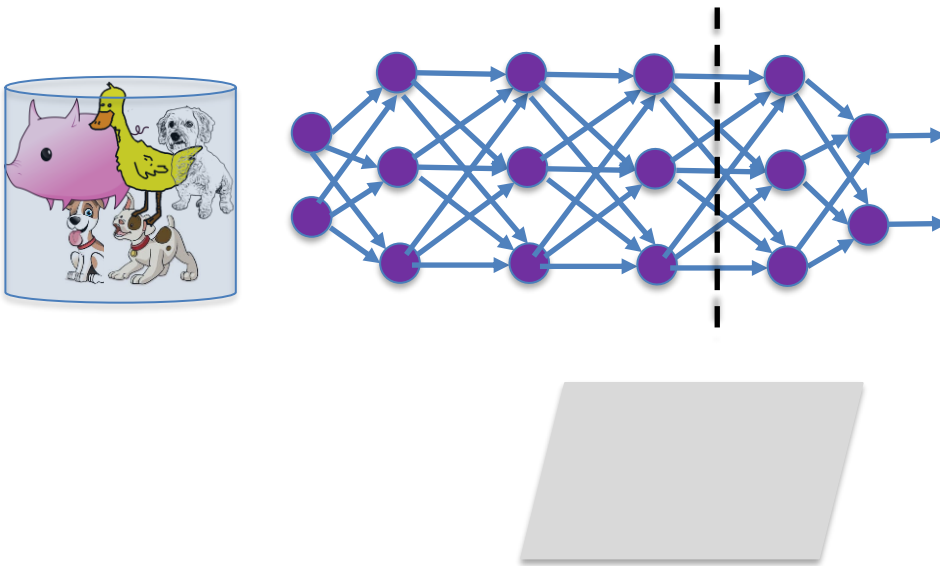


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background



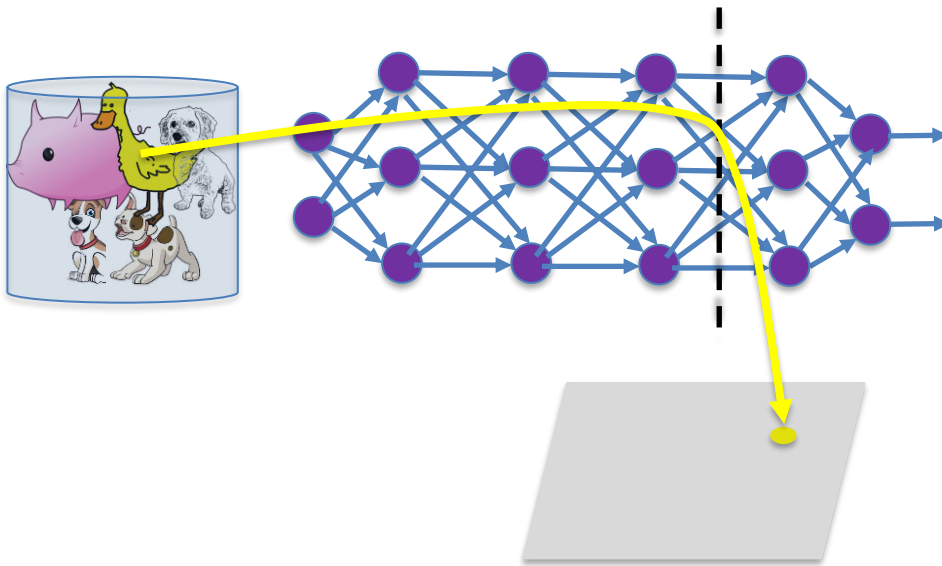
Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.



# Personalized FL through Local Memorization: background

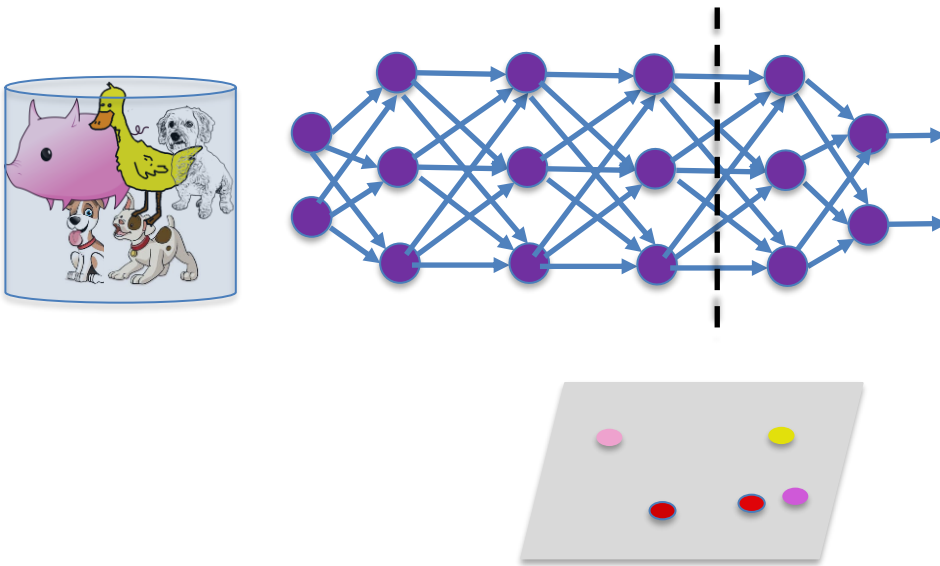


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

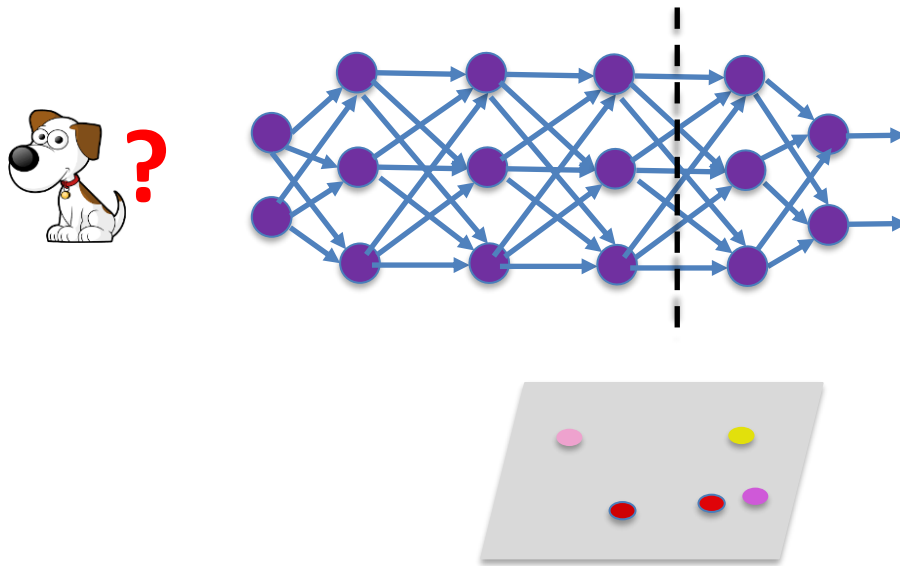


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

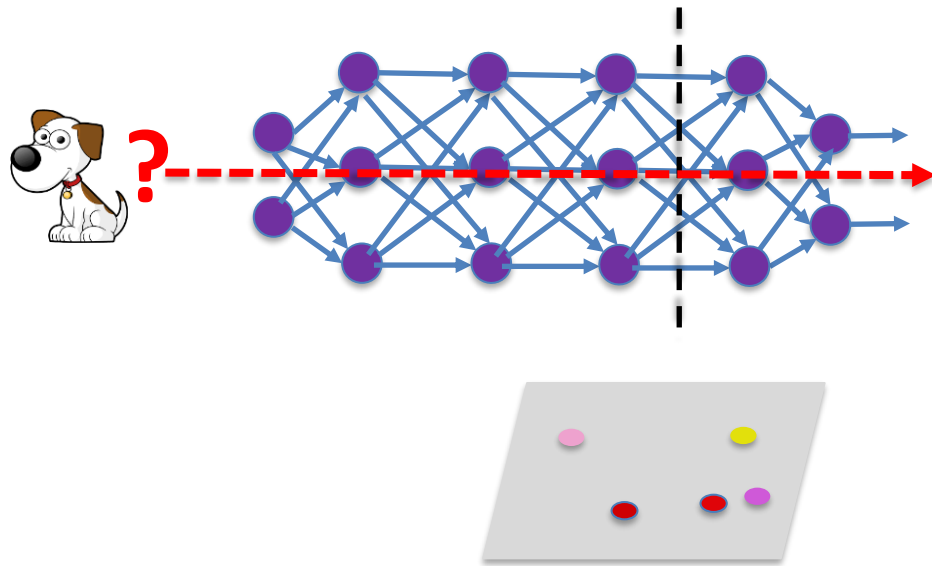


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

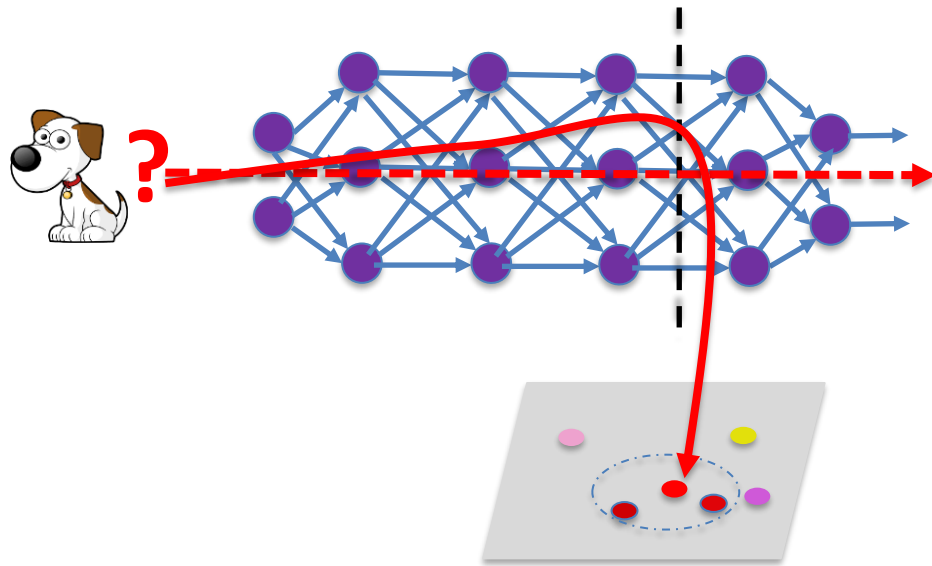


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

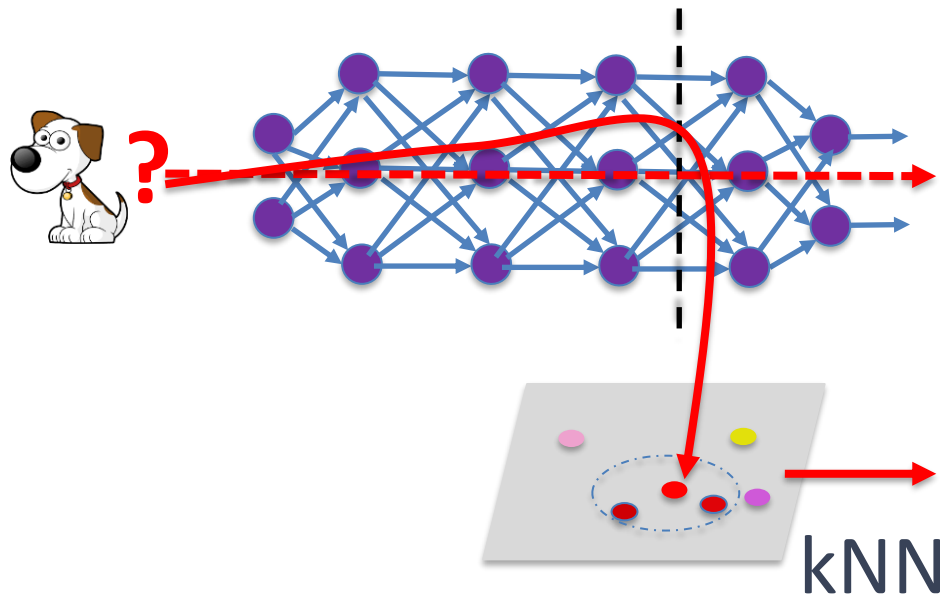


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

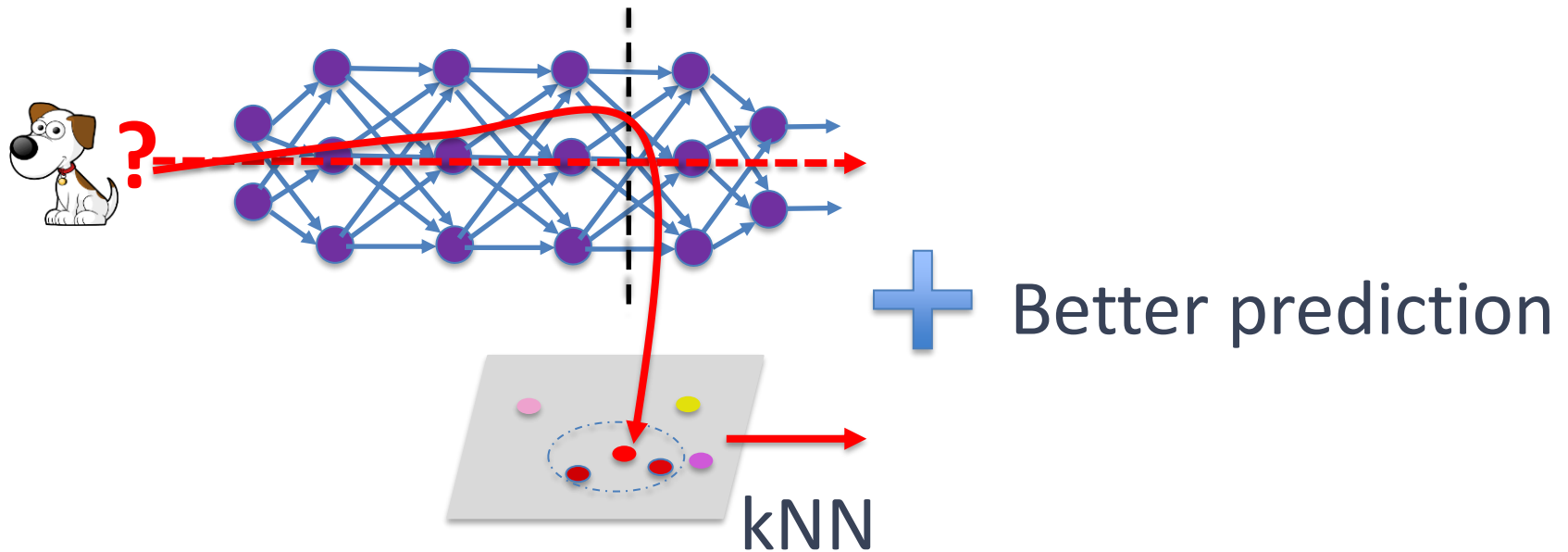


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Personalized FL through Local Memorization: background

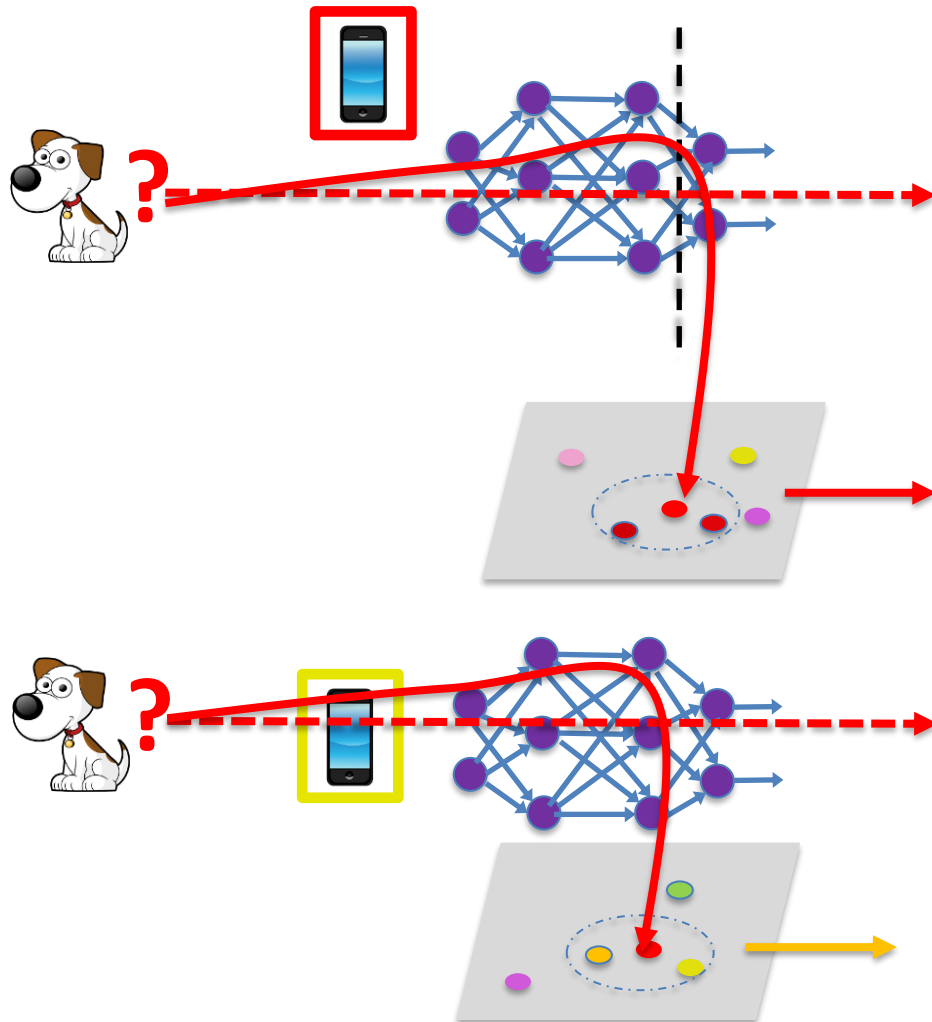


Orhan. A simple cache model for image recognition. NeurIPS, 2018.

Khandelwal, Levy, Jurafsky, Zettlemoyer, Lewis.

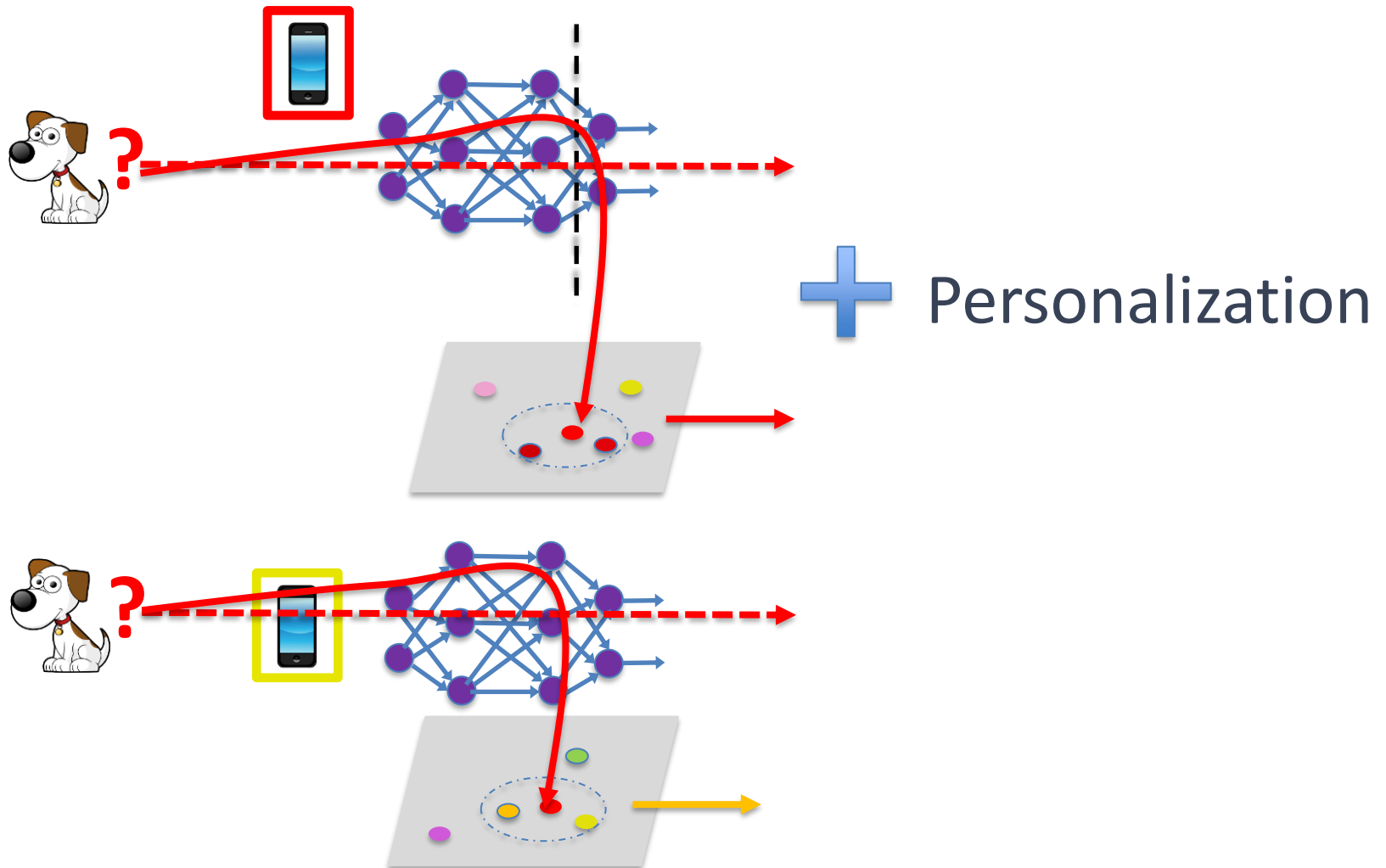
Generalization through Memorization: Nearest Neighbor Language Models. ICLR, 2020.

# Our idea: Datastore for Personalization

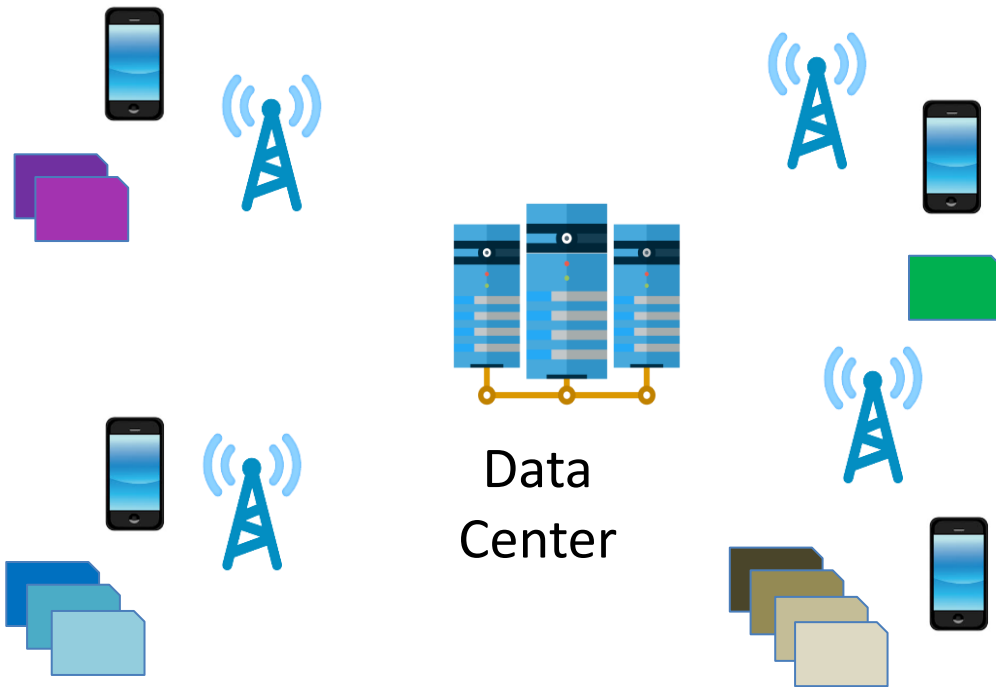




# Our idea: Datastore for Personalization



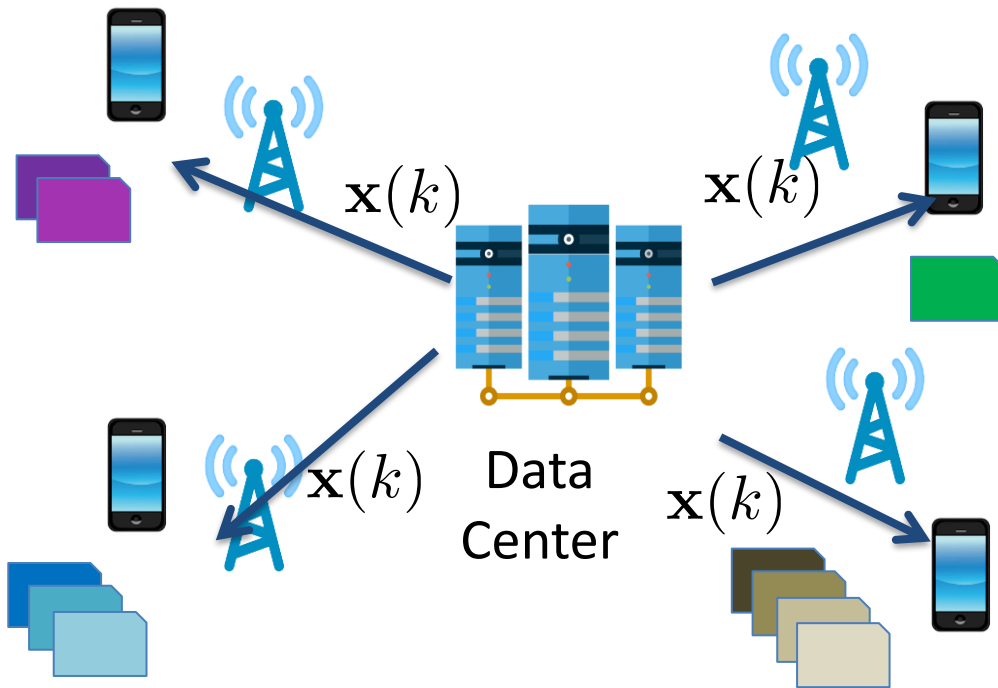
# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)

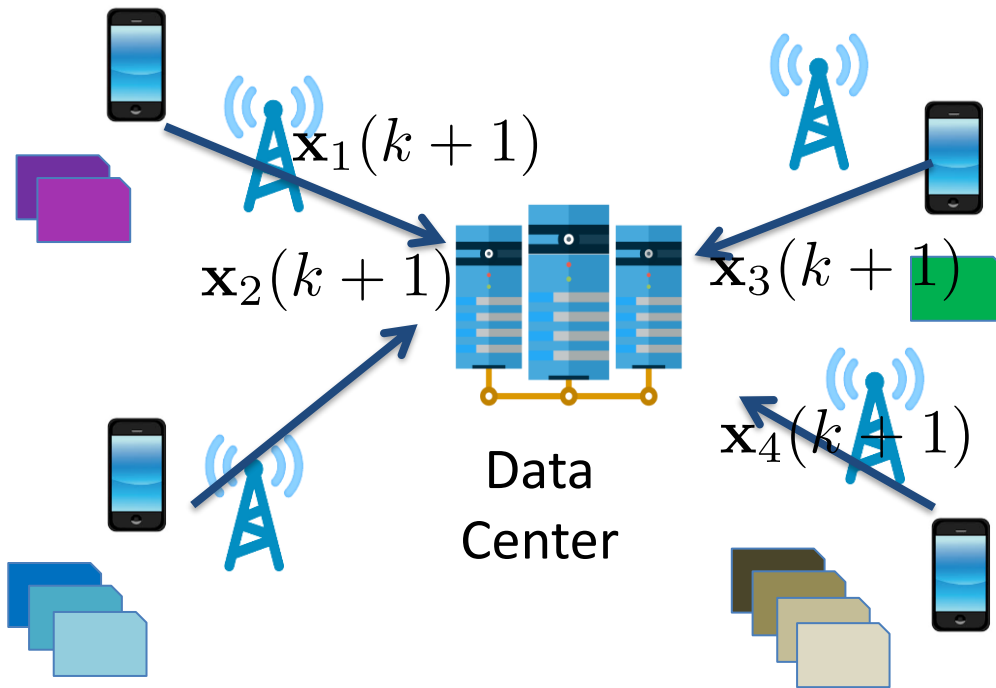
# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)

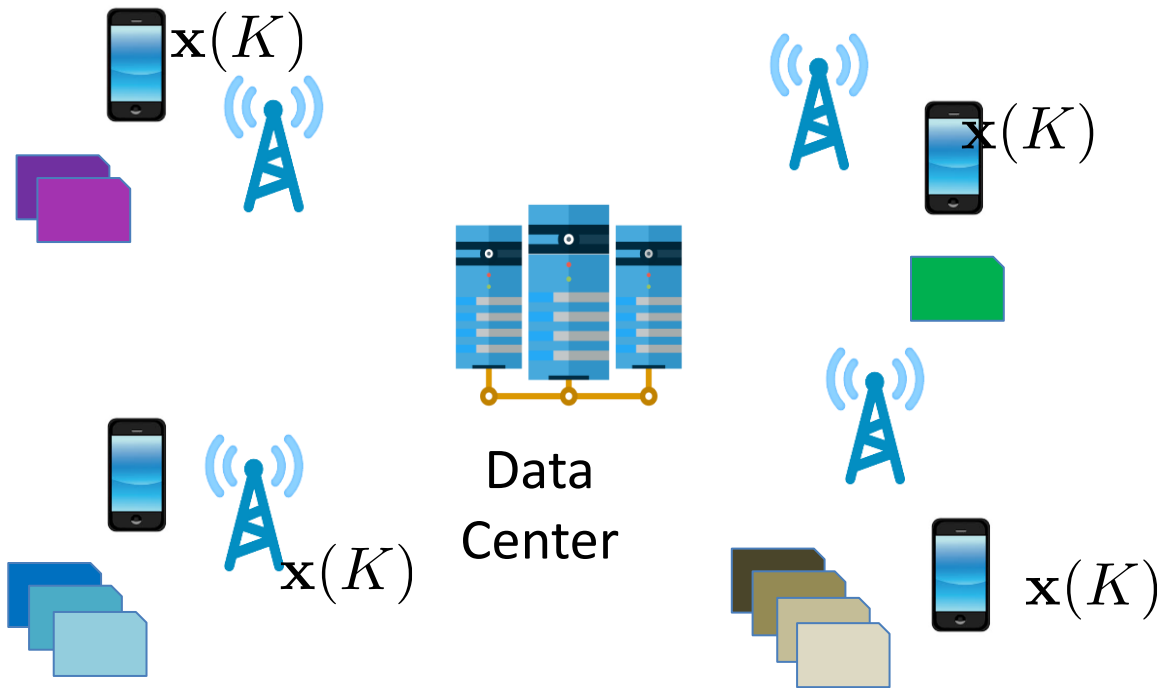
# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)

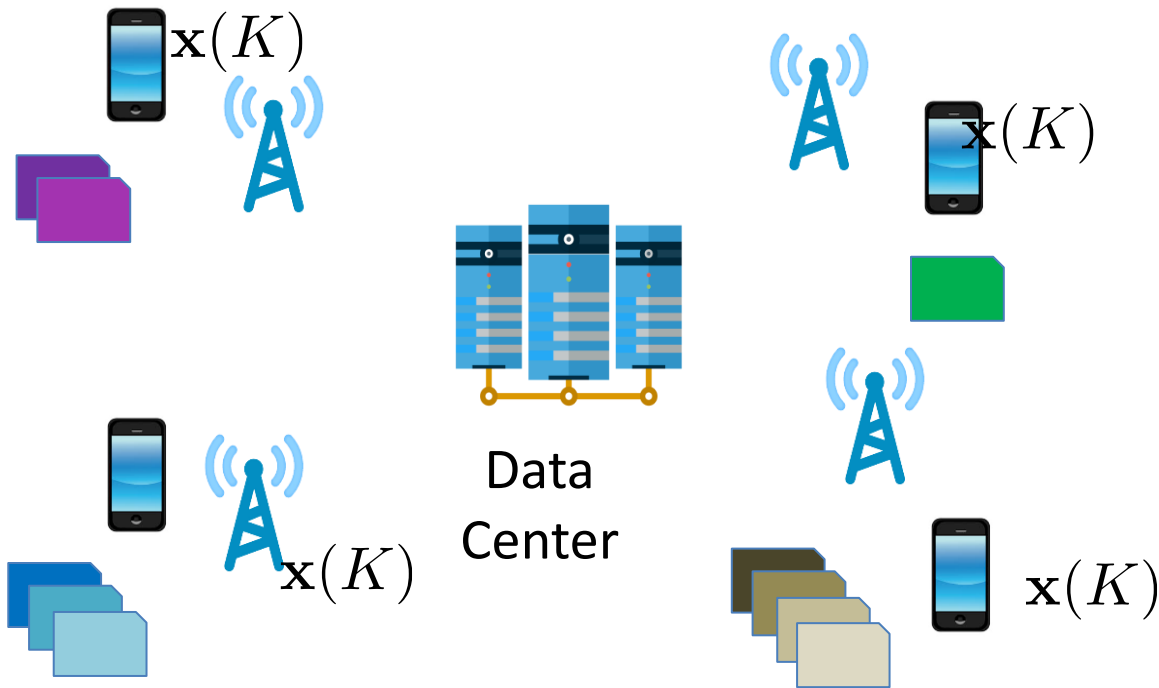
# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)

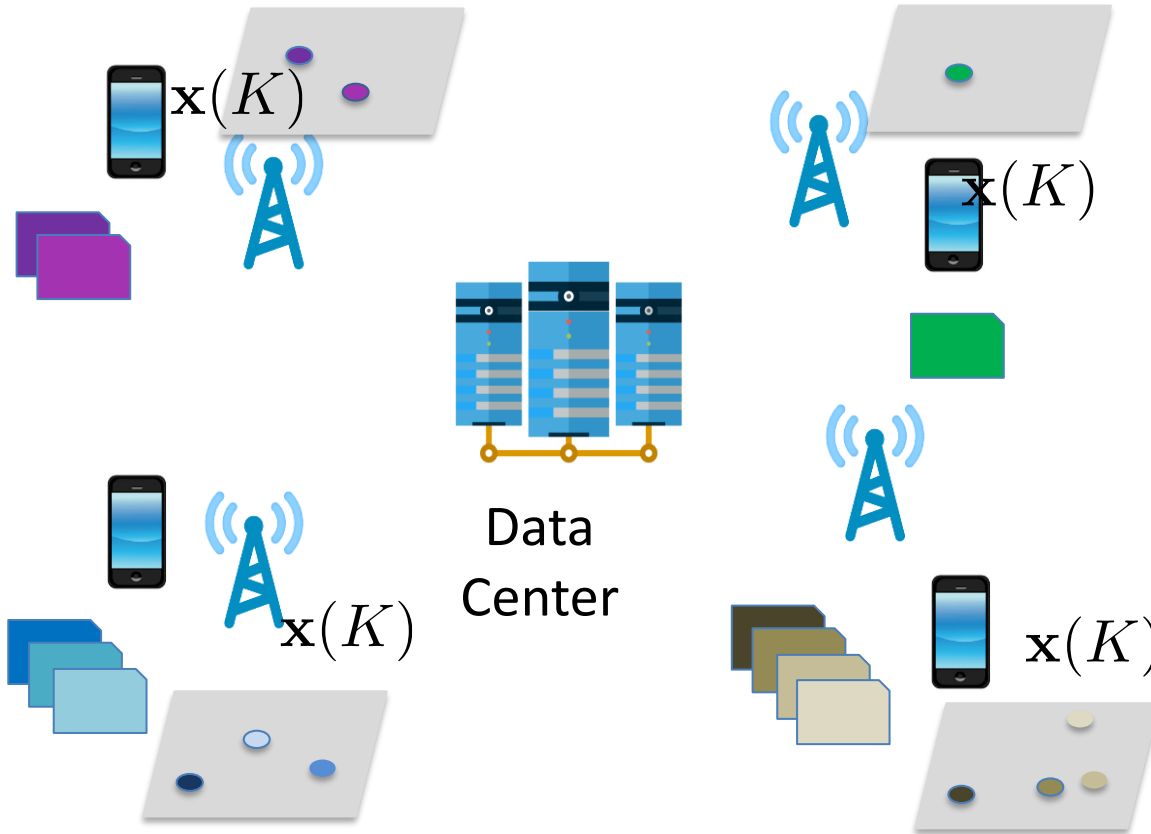
# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)
2. Each client creates its local datastore

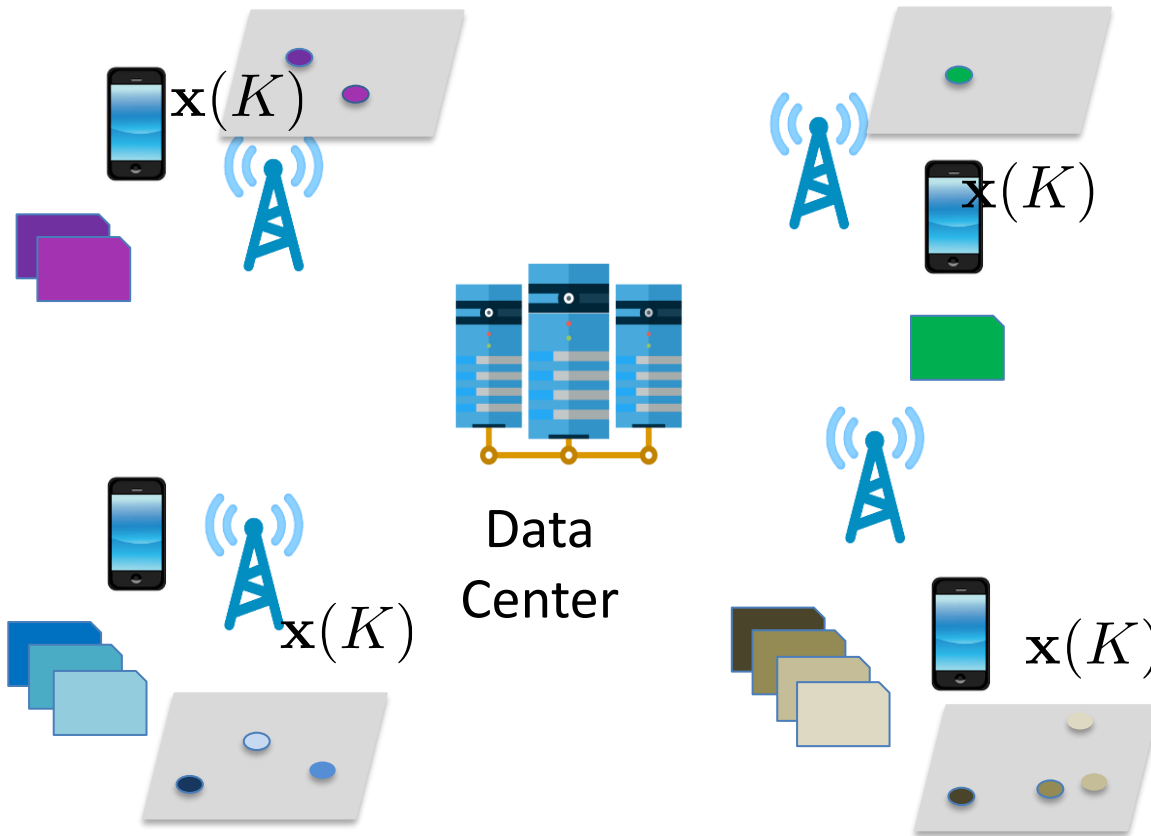
# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)
2. Each client creates its local datastore

# kNN-Per



## kNN-Per

1. Clients train a global model using a federated learning algorithm (e.g. FedAvg)
2. Each client creates its local datastore
3. A linear interpolation is used at inference

$$(1 - \lambda)h_{\text{glob}}(\mathbf{x}(K), \chi) + \lambda h_{i,k\text{NN}}(\mathbf{x}(K), \chi)$$



# kNN-Per

- Enjoys global model's convergence properties

# kNN-Per

- Enjoys global model's convergence properties
- What about generalization properties?

# kNN-Per

- Enjoys global model's convergence properties
- What about generalization properties?

$$\begin{aligned} \mathbb{E}_{\mathcal{S} \sim \otimes_{m=1}^M \mathcal{D}_m^{n_m}} [\mathcal{L}_{\mathcal{D}_m}(h_{m,\lambda})] &\leq (1 + \lambda) \cdot \mathcal{L}_{\mathcal{D}_m}(h_m^*) \\ &+ c_1 (1 - \lambda) \cdot \text{disc}_{\mathcal{H}}(\bar{\mathcal{D}}, \mathcal{D}_m) + c_3 (1 - \lambda) \cdot \sqrt{\frac{d}{n}} \cdot \sqrt{c_4 + \log\left(\frac{n}{d}\right)} \\ &+ c_2 \lambda \cdot \frac{\sqrt{p}}{p+1\sqrt{n_m}} \cdot \text{disc}_{\mathcal{H}}(\bar{\mathcal{D}}, \mathcal{D}_m) + c_5 \lambda \cdot \sqrt{\frac{d}{n}} \cdot \sqrt{c_4 + \log\left(\frac{n}{d}\right)} \cdot \frac{\sqrt{p}}{p+1\sqrt{n_m}} \end{aligned}$$

# kNN-Per

**Assumption.** Let  $h_m^* \in \arg \min_{h \in \mathcal{H}} \mathcal{L}_{D_m}(h)$ . There exist constants  $\gamma_1, \gamma_2 > 0$ , such that for any dataset  $\mathcal{S}$  drawn from  $\mathcal{X} \times \mathcal{Y}$  and any data points  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , we have

$$|\eta_m(\mathbf{x}) - \eta_m(\mathbf{x}')| \leq d(\phi_{h_{\mathcal{S}}}(\mathbf{x}), \phi_{h_{\mathcal{S}}}(\mathbf{x}')) \times (\gamma_1 + \gamma_2(\mathcal{L}_{D_m}(h_{\mathcal{S}}) - \mathcal{L}_{D_m}(h_m^*))).$$

# kNN-Per

**Assumption.** Let  $h_m^* \in \arg \min_{h \in \mathcal{H}} \mathcal{L}_{D_m}(h)$ . There exist constants  $\gamma_1, \gamma_2 > 0$ , such that for any dataset  $\mathcal{S}$  drawn from  $\mathcal{X} \times \mathcal{Y}$  and any data points  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , we have

$$|\eta_m(\mathbf{x}) - \eta_m(\mathbf{x}')| \leq d(\phi_{h_{\mathcal{S}}}(\mathbf{x}), \phi_{h_{\mathcal{S}}}(\mathbf{x}')) \times (\gamma_1 + \gamma_2(\mathcal{L}_{D_m}(h_{\mathcal{S}}) - \mathcal{L}_{D_m}(h_m^*))).$$

$\mathbf{x}$  &  $\mathbf{x}'$   
in same class?

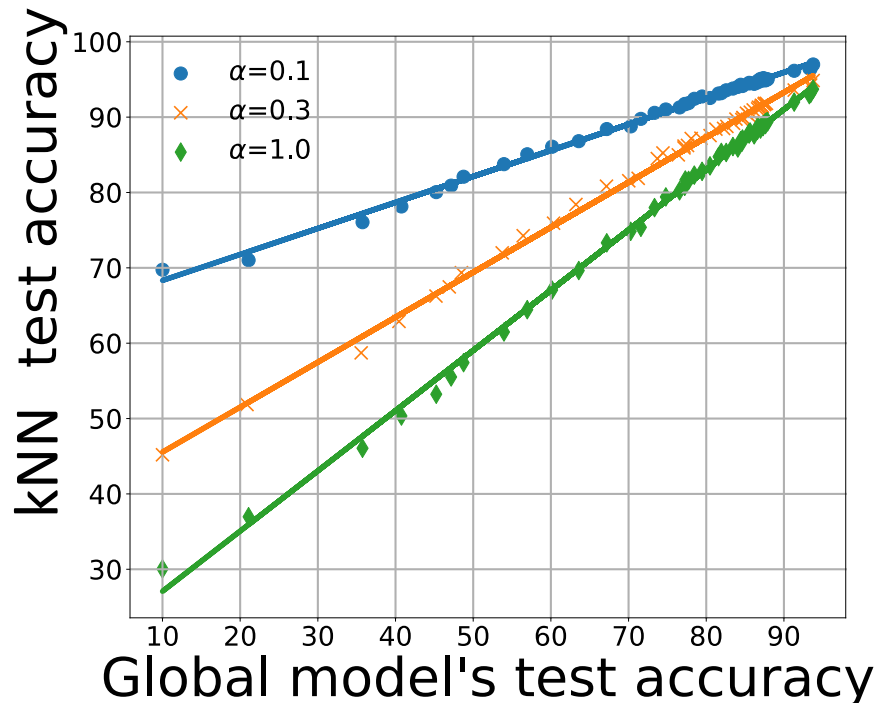
representations'  
distance

global model's  
quality for client  $m$

# kNN-Per

**Assumption.** Let  $h_m^* \in \arg \min_{h \in \mathcal{H}} \mathcal{L}_{D_m}(h)$ . There exist constants  $\gamma_1, \gamma_2 > 0$ , such that for any dataset  $\mathcal{S}$  drawn from  $\mathcal{X} \times \mathcal{Y}$  and any data points  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ , we have

$$|\eta_m(\mathbf{x}) - \eta_m(\mathbf{x}')| \leq d(\phi_{h_{\mathcal{S}}}(\mathbf{x}), \phi_{h_{\mathcal{S}}}(\mathbf{x}')) \times (\gamma_1 + \gamma_2(\mathcal{L}_{D_m}(h_{\mathcal{S}}) - \mathcal{L}_{D_m}(h_m^*))).$$



CIFAR-10

# kNN-Per

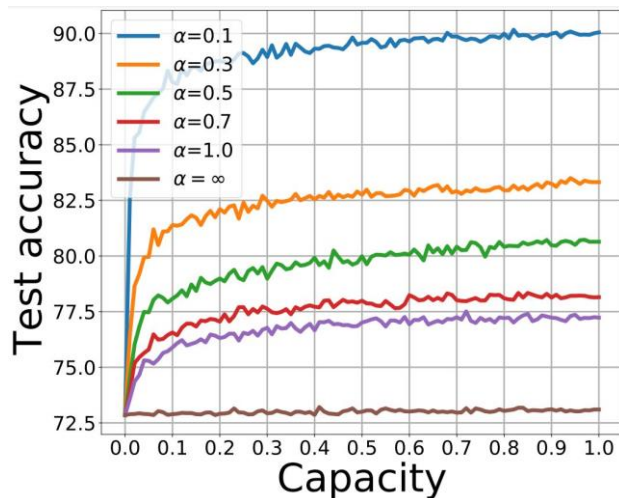
Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>



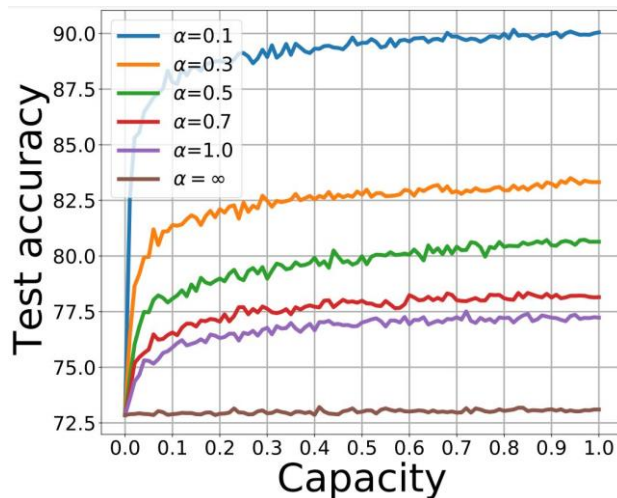
CIFAR-10



# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>



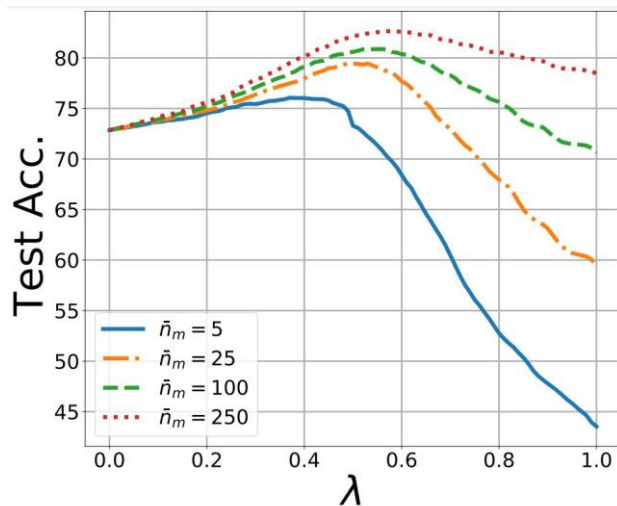
The benefit of kNN-Per is larger when data distributions are more heterogenous

CIFAR-10

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>

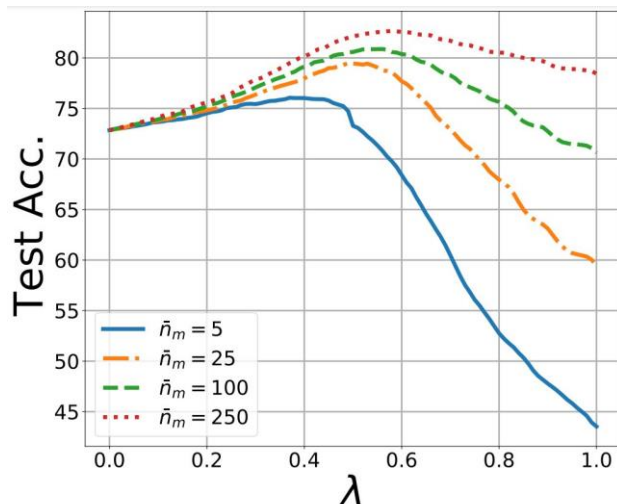


CIFAR-10

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>



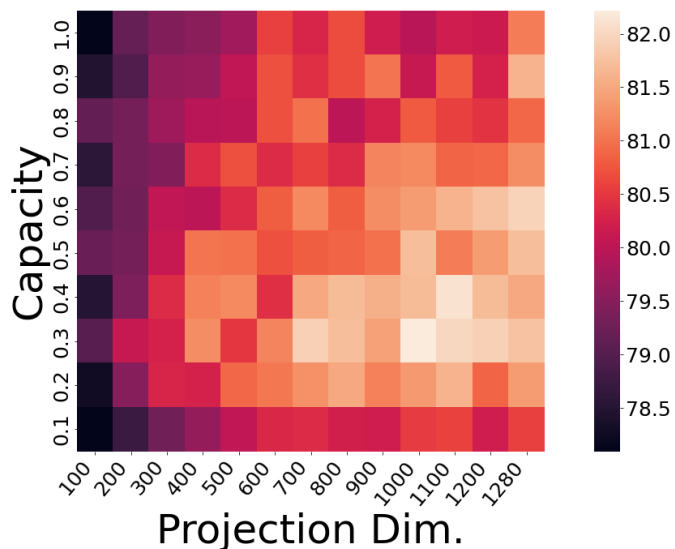
kNN-Per relies mostly on kNN  
for datasets with more than 100  
samples

CIFAR-10

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>



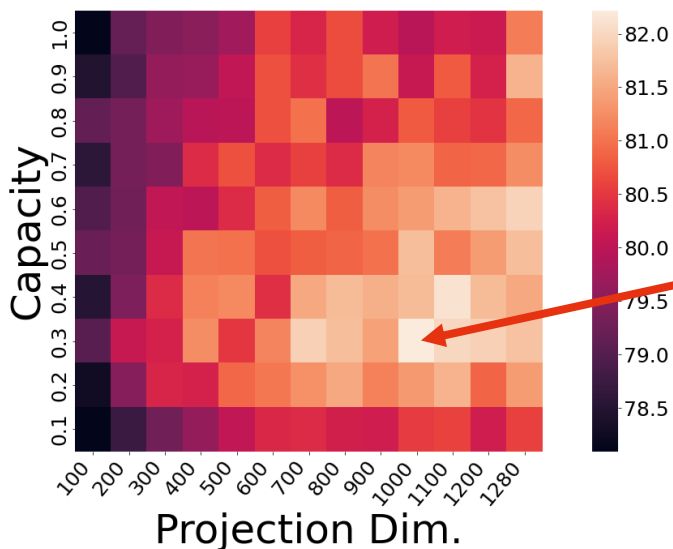
ProtoNN-like datastore compression

CIFAR-10

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>



ProtoNN-like datastore compression

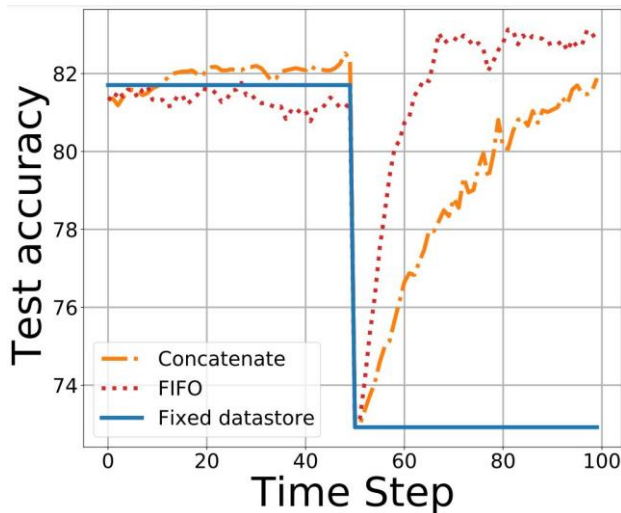
4x memory savings with  
limited accuracy loss (0.7pp)

CIFAR-10

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>

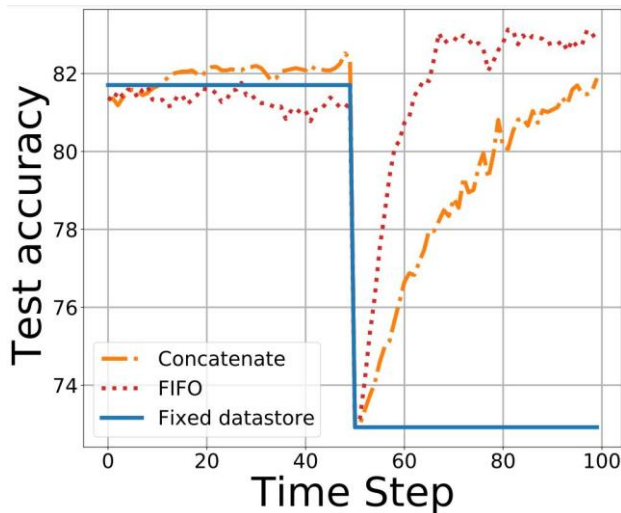


CIFAR-10

# kNN-Per

Table 2: Test accuracy: average across clients / bottom decile.

Dataset	Local	FedAvg	FedAvg+	ClusteredFL	Ditto	FedRep	APFL	kNN-Per (Ours)
FEMNIST	71.0 / 57.5	83.4 / 68.9	84.3 / 69.4	83.7 / 69.4	84.3 / 71.3	85.3 / 72.7	84.1 / 69.4	<b>88.2 / 78.8</b>
CIFAR-10	57.6 / 41.1	72.8 / 59.6	75.2 / 62.3	73.3 / 61.5	80.0 / 66.5	77.7 / 65.2	78.9 / 68.1	<b>83.0 / 71.4</b>
CIFAR-100	31.5 / 19.8	47.4 / 36.0	51.4 / 41.1	47.2 / 36.2	52.0 / 41.4	53.2 / 41.7	51.7 / 41.1	<b>55.0 / 43.6</b>
Shakespeare	32.0 / 16.0	48.1 / 43.1	47.0 / 42.2	46.7 / 41.4	47.9 / 42.6	47.2 / 42.3	45.9 / 42.4	<b>51.4 / 45.4</b>



kNN-Per is robust to distribution shift

CIFAR-10

# Questions?



Paper



Code