

To Store or Not?

Online Data Selection for Federated Learning with Limited Storage

Chen Gong¹, Zhenzhe Zheng¹, Yunfeng Shao², Bingshuai Li², Fan Wu¹, Guihai Chen¹

Shanghai Jiao Tong University¹

Huawei Noah's Ark Lab²



NOAH'S ARK LAB

From Big Data to Deep Knowledge

Mobile Network needs ML models



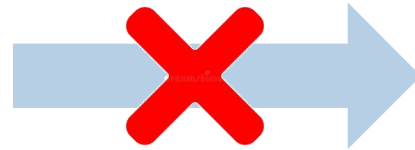
Large
Scale



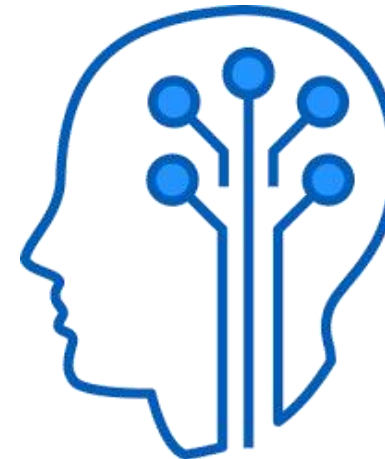
Complex
Correlation



Low
Latency

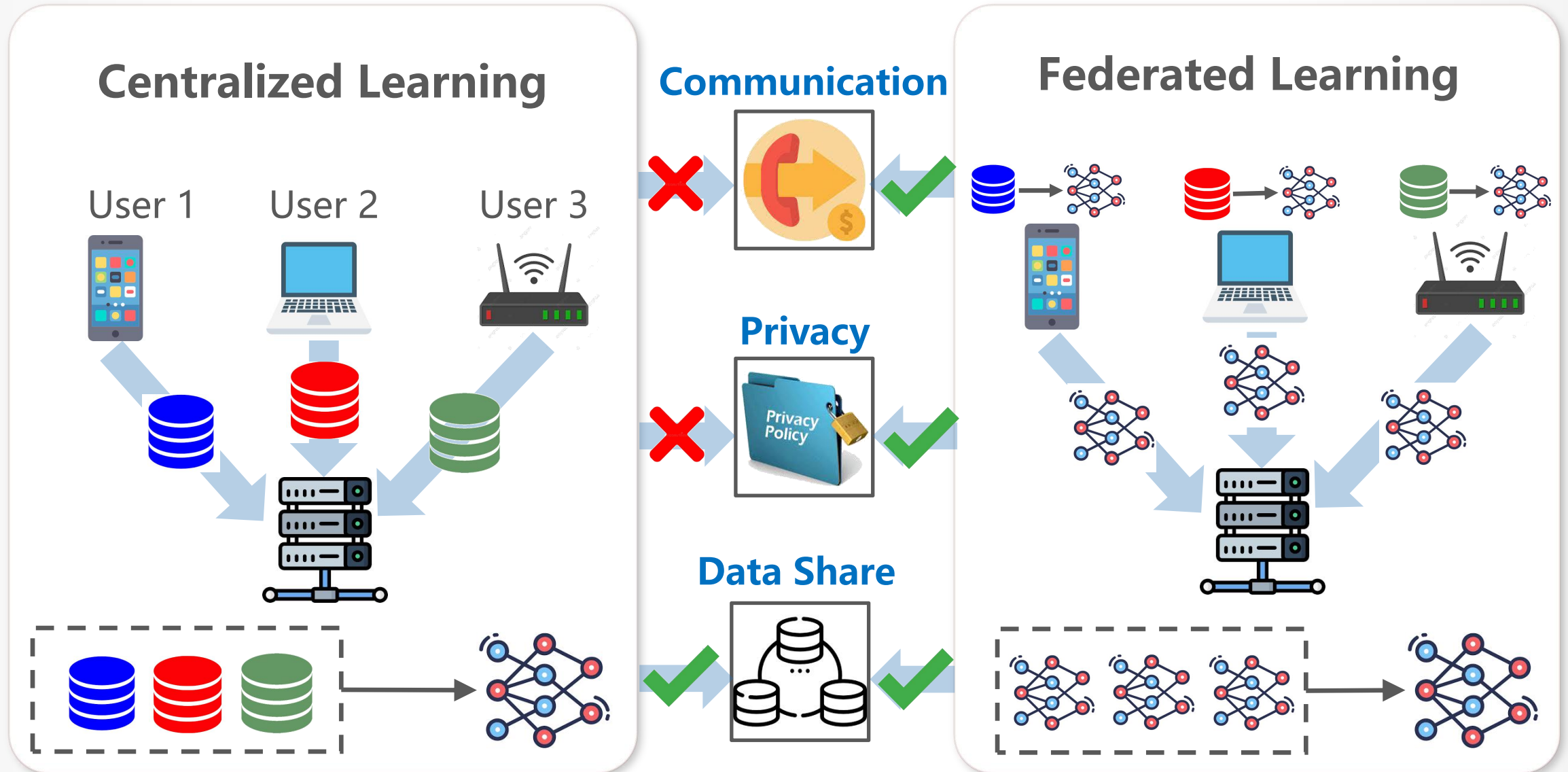


Human-in-
the-loop
Approach

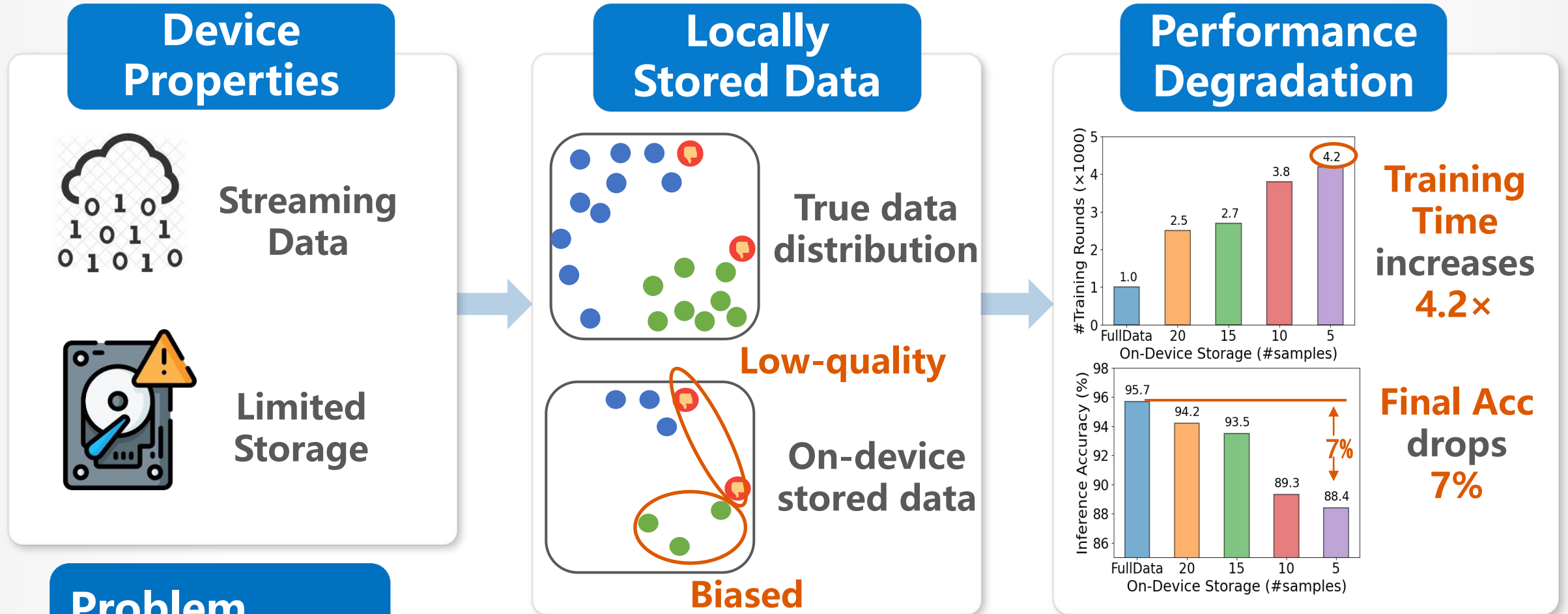


Machine
Learning
Models

ML Model Training needs FL



Device Properties Degrade FL Performance

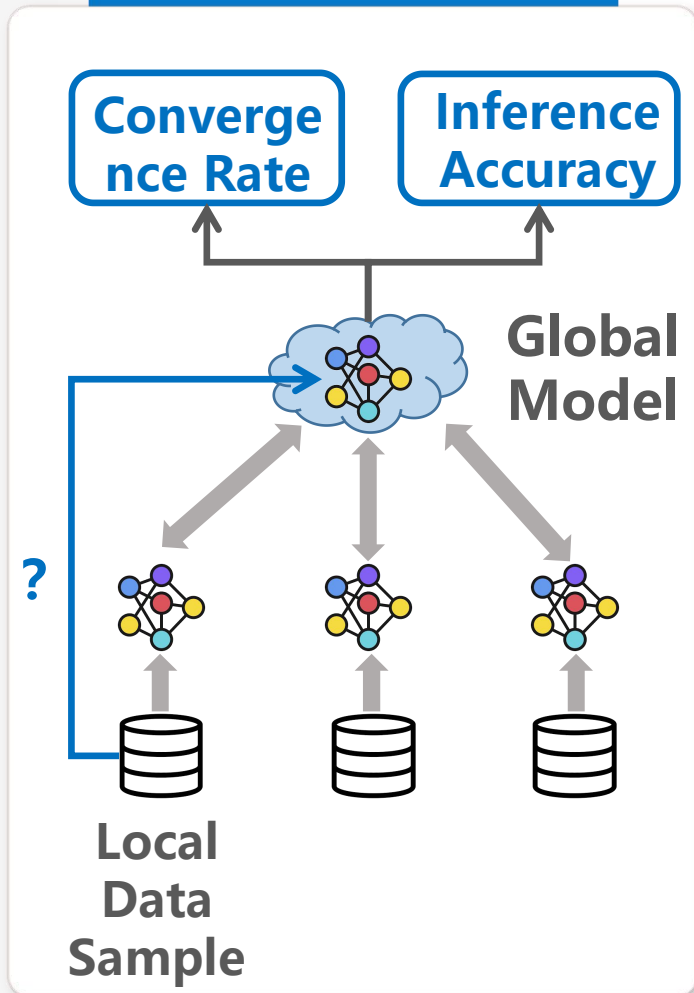


Problem

Filter valuable on-device data to simultaneously improve training convergence and inference accuracy of FL model

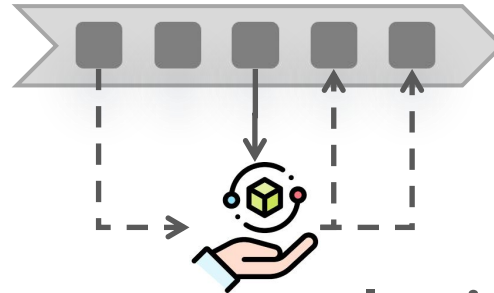
Three Key Challenges

Theoretical Guarantee



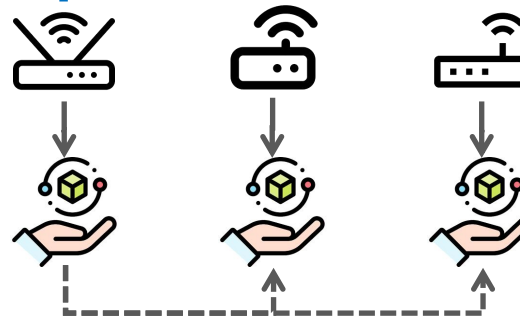
Lack of Multi-Dim Information

Temporal Information



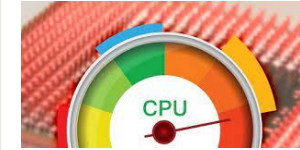
Inaccurate Evaluation

Spatial Information



Redundant Stored Data

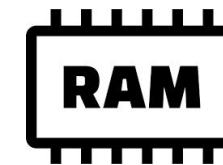
Constrained Resource



Computation Cost



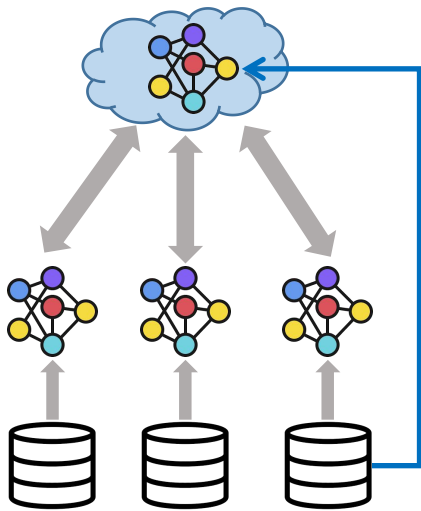
Evaluation Delay



Memory Footprint

Theoretical Analysis

Local Data Impact



Convergence Rate

$$\underbrace{\left| F(w_{\text{fed}}^{t-1}) - F(w_{\text{fed}}^t) \right|}_{\text{global loss reduction}} \geq \sum_{c \in \mathcal{C}_t} \sum_{i=0}^{m-1} \sum_{(x,y) \in B_c} \left[-\alpha_c \underbrace{\| \nabla_w l(w_c^{t,i}, x, y) \|_2^2}_{\text{term 1}} \right. \\
 \left. + \beta_c \underbrace{\langle \nabla_w l(w_c^{t,i}, x, y), \nabla_w F(w_{\text{fed}}^{t-1}) \rangle}_{\text{term 2}} \right],$$

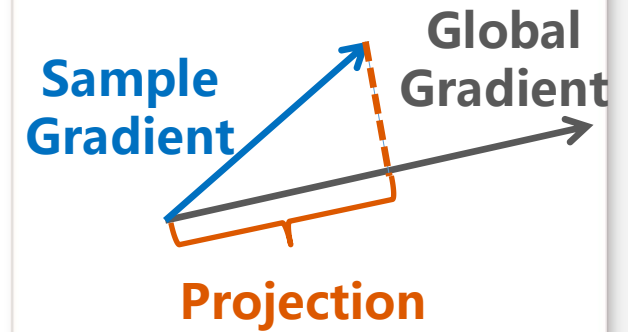
Dominant Term

Inference Accuracy

$$\underbrace{\| w_{\text{fed}}^t - w_{\text{cen}}^{mt} \|_2}_{\Delta \text{model of CL and FL}} \leq (1 + \eta L)^m \| w_{\text{fed}}^{t-1} - w_{\text{cen}}^{m(t-1)} \|_2 \\
 + \sum_{c \in \mathcal{C}_t} \zeta_c^t \left[\eta \sum_{i=0}^{m-1} (1 + \eta L)^{m-1-i} \underbrace{G_c(w_c^{t,i})}_{\text{Dominant Term}} \right],$$

Dominant Term

Common Term

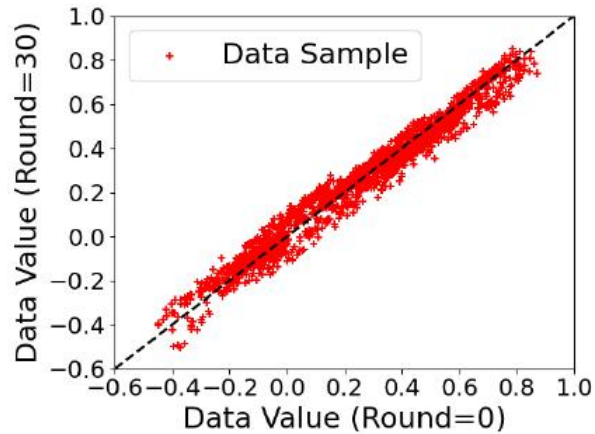


Intuition

- Preserve **personal information**
- Reduce global data **heterogeneity**

7 On-Device Data Selection

Lack of Latest Global Model



(b) Round 0 vs Round 30

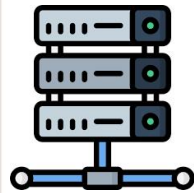
Data value remains stable across few training rounds

Lack of Accurate Global Gradient



Local Gradient Estimator

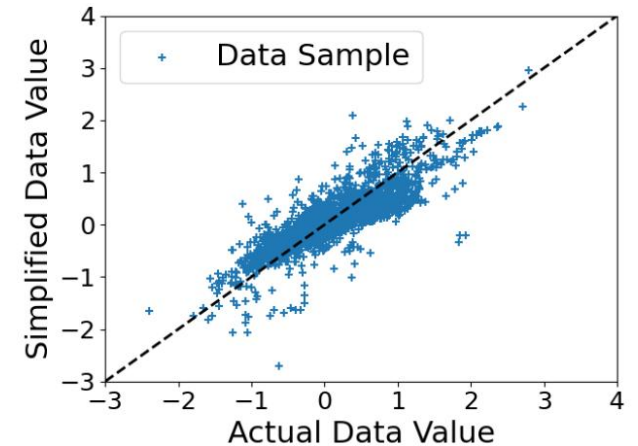
$$\hat{g}_c \leftarrow \frac{n-1}{n} \hat{g}_c + \frac{1}{n} \nabla_w l(w_{\text{fed}}^{t_c, \text{last}-1}, x, y).$$



Global Gradient Estimator

$$\hat{g}^t \leftarrow \hat{g}^{t-1} + \sum_{c \in \mathcal{C}_t} \zeta_c (\hat{g}_c - \hat{g}_c^{t_{\text{last}}}),$$

Simplification Version

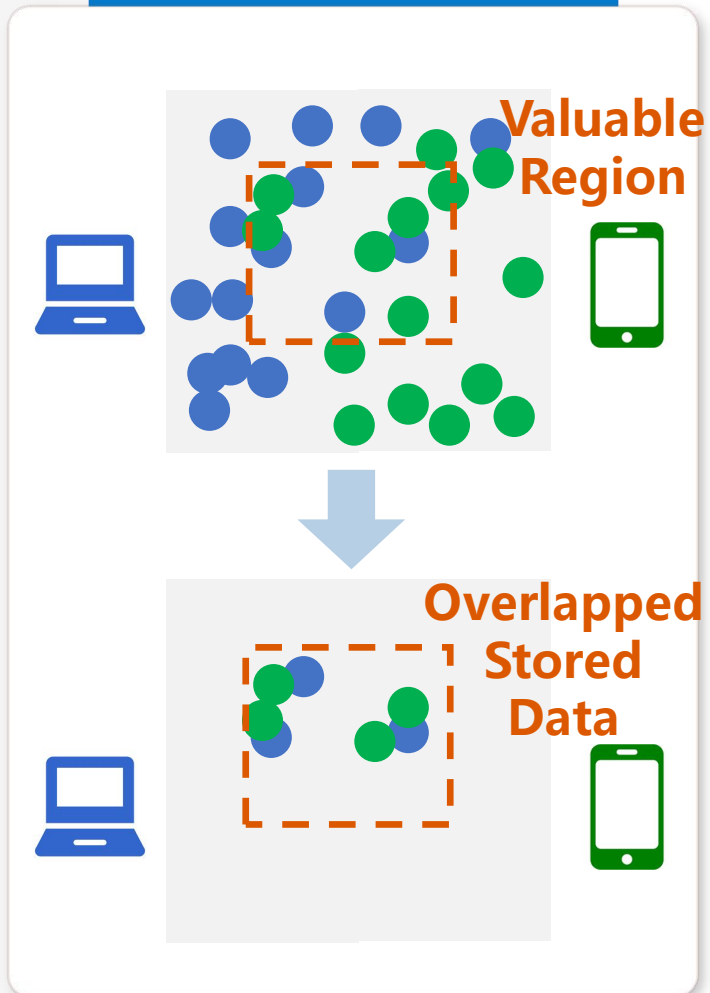


(c) Full model vs Last layer

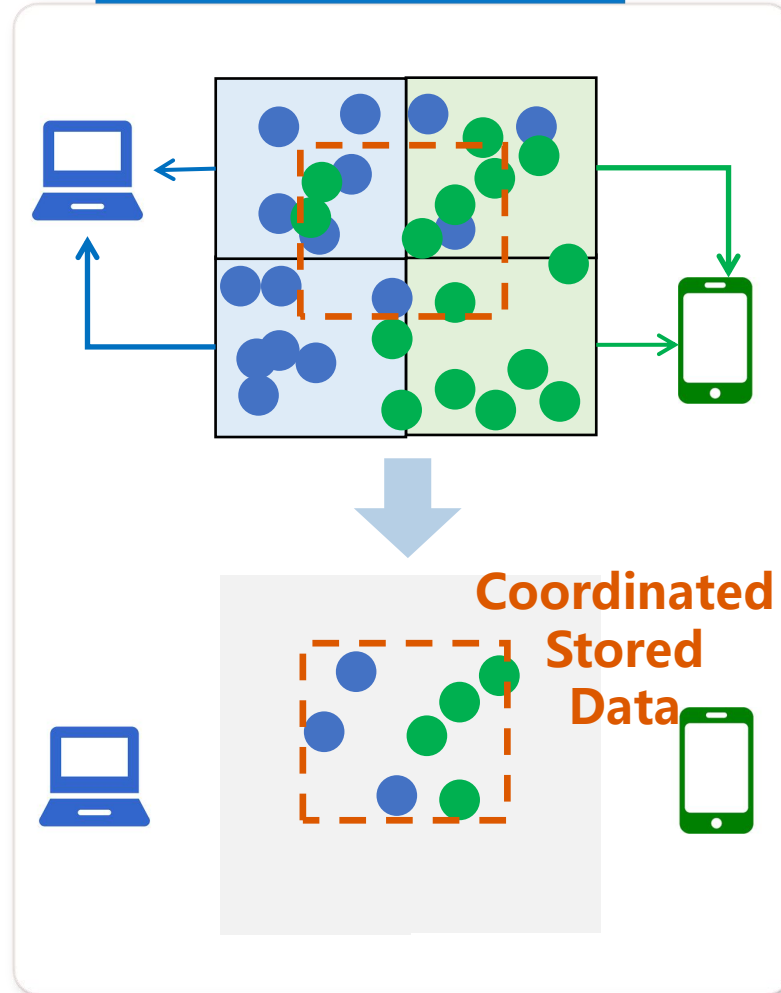
Partial model gradient is able to reflect the **overall tendency**

Cross-Device Data Storage

Overlapped Stored Data



Coordination Policy



Problem Formulation

- **Efficient Data Selection**
- **Redundant Label Assignment**
- **Limited Storage**
- **Unbiased global data distribution**

$$\begin{aligned}
 \max_D \quad & \sum_{c \in C} \sum_{y \in Y} D_{c,y} \cdot V_{c,y} = \| DV \|_1, \\
 \text{s.t.} \quad & \| D_y^T \|_0 \geq n_y^{\text{label}}, & \forall y \in Y, \\
 & \| D_c \|_0 \leq n_c^{\text{client}}, & \forall c \in C, \\
 & \| D_c \|_1 = |B_c|, & \forall c \in C, \\
 & \frac{\sum_{c \in C} D_{c,y}}{\| B \|_1} = \frac{\sum_{c \in C} V_{c,y}}{\| V \|_1}, & \forall y \in Y.
 \end{aligned}$$

Experiments

▶ Learning Tasks and Models

Synthetic
Image Classification
Activity Recognition
Traffic Classification

Datasets	Models	#Samples	#Labels	#Devices	n_y^{label}	$\frac{ C_t }{ C }$	η	m	$ B_c $
Synthetic Dataset [31]	LogReg	1,016,442	10	200	5	5%	$1e^{-4}$	5	10
Fashion-MNIST [32]	LeNet [58]	70,000	10	50	5	10%	$1e^{-3}$	5	5
HARBOX [33]	Customized DNN	34,115	5	120	5	10%	$1e^{-3}$	5	5
Industrial Dataset	Customized CNN	37,853	20	30	5	20%	$5e^{-3}$	5	10

Videos



YouTube



Games



Downloading



App Store



Communication



Experiments

Overall Performance

Task	Model Training Speedup							
	RS	HL	GN	FB	SLD	ODE-Exact	ODE-Est	FD
ST	1.0×	—	4.87×	—	4.08×	9.52×	5.88×	2.67×
IC	1.0×	—	—	—	—	1.35×	1.20×	1.01×
HAR	1.0×	—	—	—	—	2.22×	1.55×	4.76×
TC	1.0×	—	—	—	—	2.51×	2.50×	3.92×

Task	Inference Accuracy							
	RS	HL	GN	FB	SLD	ODE-Exact	ODE-Est	FD
ST	79.56%	78.44%	83.28%	78.56%	82.38%	87.12%	82.80%	88.14%
IC	71.31%	51.95%	41.45%	60.43%	69.15%	72.71%	72.70%	71.37%
HAR	67.25%	48.16%	51.02%	48.33%	56.24%	73.63%	70.39%	77.54%
TC	89.3%	69.00%	69.3%	72.19%	72.30%	95.3%	95.30%	96.00%

Time:
2.51× speedup

Accuracy:
6% increase

Memory
< 15MB

Delay
1ms

Task	Memory Footprint (MB)				
	RS	HL	GN	ODE-Est	ODE-Simplified
IC	1.70	11.91	16.89	18.27	16.92
HAR	1.92	7.27	12.23	13.46	12.38
TC	0.75	10.58	19.65	25.15	14.47

Task	Evaluation Time (ms)				
	RS	HL	GN	ODE-Est	ODE-Simplified
IC	0.05	11.1	21.1	22.8	11.4
HAR	0.05	0.36	1.04	1.93	0.53
TC	0.05	1.03	9.06	9.69	1.23

Robustness Analysis

Number of Local Training Epochs (m)

With the increasing of local training epoch number, ODE achieves higher final accuracy improvement

Device Participation Rate

With 10% participation rate, ODE achieves $2.57\times$ training time speed up and 6.6% increase on inference accuracy.

On-Device Storage Capacity

ODE has stable performance with different device capacity, and can reduce up to 50% storage compared with baseline .

Conclusion

- Identify **two practical properties** of mobile devices and demonstrate the enormity
- **Theoretically** analyze impact of an individual local data sample on global model
- Design a collaborative data selection **framework** for FL to **simultaneously** improve convergence rate and final inference accuracy
- Achieve **remarkable** training speedup and accuracy improvement on **industrial** traffic classification task



Thanks for Watching !

Q & A

Please refer to our paper for more details !