

ILLINOIS PerAda: Parameter-Efficient Federated Learning Personalization with Generalization Guarantees



Chulin Xie¹ De-An Huang² Wenda Chu³ Daguang Xu² Chaowei Xiao^{2,4} Bo Li^{1,5} Anima Anandkumar³
¹ University of Illinois Urbana-Champaign ² NVIDIA ³ Caltech ⁴ University of Wisconsin-Madison ⁵ University of Chicago

<https://github.com/NVlabs/PerAda>



Background & Motivation

Personalized Federated Learning (FL)

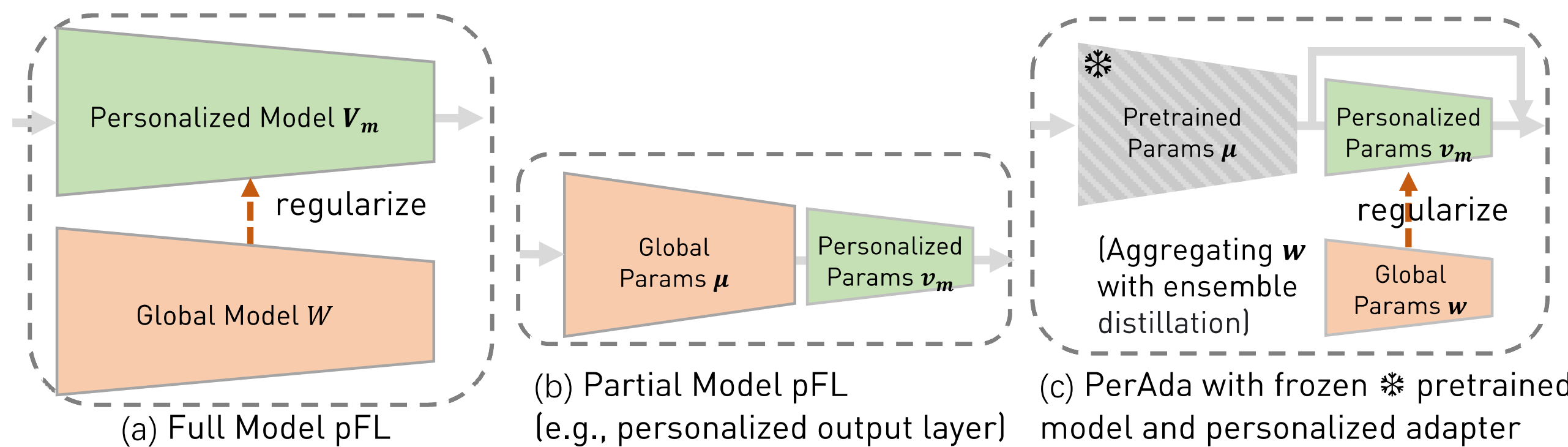
- **Personalization:** each client trains a personalized model on local data
- **Generalization:** clients leverage aggregated knowledge from other clients

Existing work

- **Full model** personalization: Each client trains a personalized model and a copy of global model aggregated by the server for regularization
- **Partial model** personalization:
 - Split model into personalized and shared parameters
 - Only shared parameters are aggregated

Challenge

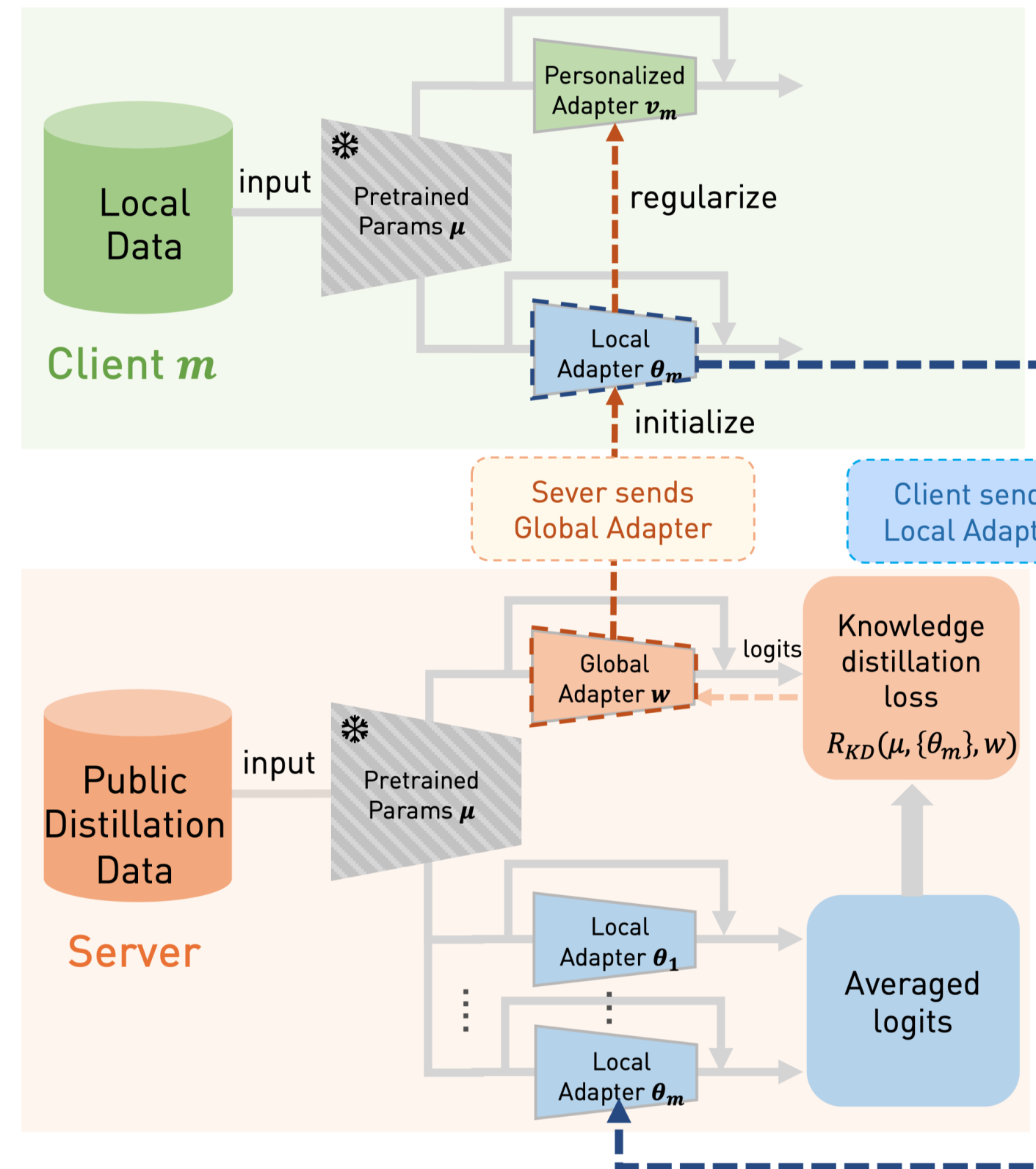
- **High communication and computational costs:** Full model personalization doubles the memory at each client
- **Limited generalization:** Partial model personalization overfits to local data and struggles with distribution shifts
 - shared parameters do not encode generalized knowledge well compared to full global model



This Work: PerAda

- **Framework:** parameter-efficient personalized FL using Adapters and Knowledge Distillation (KD)
- **Benefits:**
 - Reduce communication and computation costs with a pretrained model and adapters
 - Achieve personalization while maintaining generalization to test-time distribution shifts with regularization and KD
- **Theoretical justification** of PerAda (in paper)
 - Convergence analysis for global & personalized models
 - generalization guarantees for global & personalized models

Method



Personalized objective:

- train personalized adapter with **regularization** towards a global adapter to prevent overfitting
- $$\min_{v_m} P_m(v_m, w) := \mathcal{L}_m(u, v_m) + \frac{\lambda}{2} \|v_m - w\|^2, \quad (\text{Personal Obj})$$

Global objective:

- leverage server-side **ensemble distillation** to enrich the global adapter with ensemble knowledge from clients' local models
 - avoid directly averaging clients' models
- $$\min_w \mathcal{R}_{KD}(u, \{\theta_m\}_{m=1}^M, w) \quad (\text{Global Obj})$$
- where $\theta_m = \arg \min_{\theta} \mathcal{L}_m(u, \theta)$, initialized with w .
- \mathcal{R}_{KD} is average distillation loss (between the averaged logits of local models and logits of the global model) on an **auxiliary (unlabeled) dataset**

Experiments

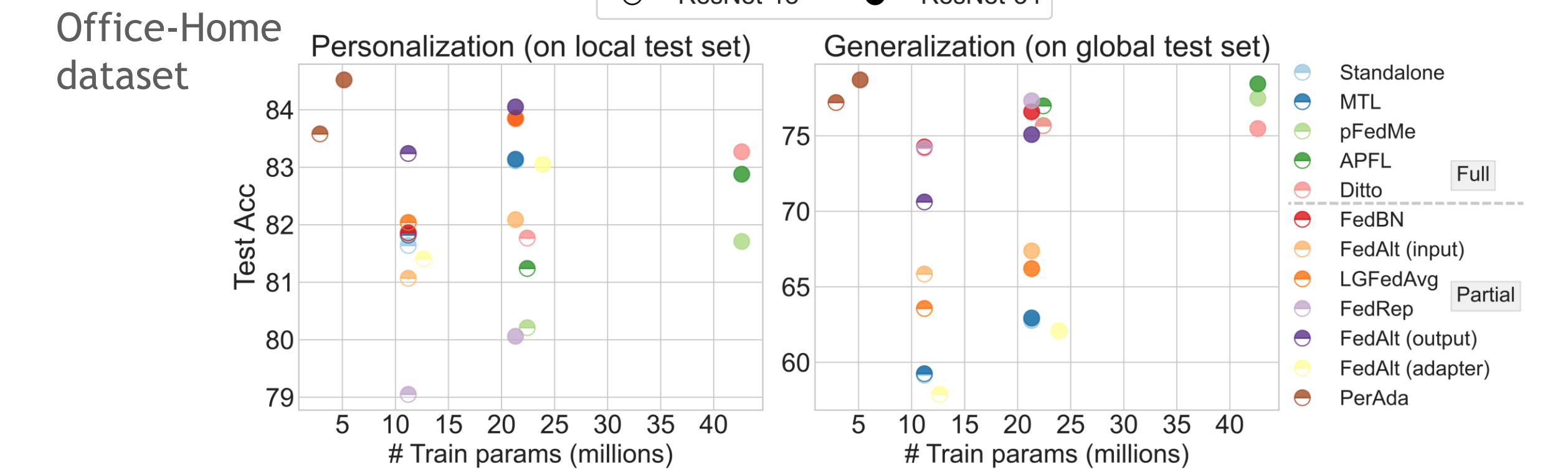
Parameter-efficiency and averaged test accuracy across all clients' personalized models

- Local-test: clients' corresponding local test data \rightarrow personalization
- Global-test: the union of clients' local test data \rightarrow generalization

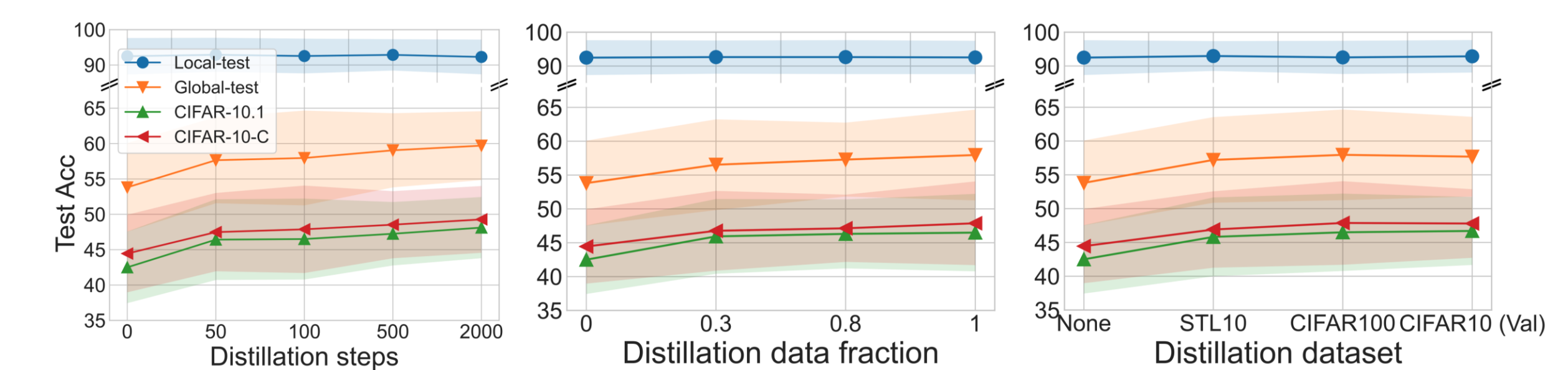
Algorithm	Personalized Params	# Trained Params	# Comm. Params	CIFAR-10				Office-Home		CheXpert	
				Local-test	Global-test	CIFAR-10.1	CIFAR-10-C	Local-test	Global-test	Local-test	Global-test
STANDALONE	Full model	11.18M	0M	85.94 \pm 8.82	29.77 \pm 8.09	25.82 \pm 6.27	26.67 \pm 7.07	81.64 \pm 6.08	59.15 \pm 3.32	65.06 \pm 1.88	65.45 \pm 2.3
MTL [57]	Full model	11.18M	11.18M	86.24 \pm 8.46	29.46 \pm 8.33	25.64 \pm 6.42	26.4 \pm 7.29	81.82 \pm 5.53	59.25 \pm 2.84	65.15 \pm 1.95	65.48 \pm 2.3
FEDAVG+FT [65]	Full model	11.18M	11.18M*	88.91 \pm 6.71	43.99 \pm 9.57	35.49 \pm 8.02	36.51 \pm 8.36	79.42 \pm 5.62	77.19 \pm 0.56	70.16 \pm 0.78	70.6 \pm 0.31
PFEDME [59]	Full model	22.36M	11.18M	90.73 \pm 4.67	45.06 \pm 8.65	36.51 \pm 7.2	37.65 \pm 7.6	80.21 \pm 5.32	75.69 \pm 0.69	65.07 \pm 1.2	64.86 \pm 1.22
APFL [10]	Full model	22.36M	11.18M	90.74 \pm 4.76	43.92 \pm 9.18	35.83 \pm 7.5	36.51 \pm 7.94	81.24 \pm 4.51	76.98 \pm 1.39	68.98 \pm 1.04	68.96 \pm 1.1
DITTO [33]	Full model	22.36M	11.18M	90.21 \pm 4.61	53.82 \pm 6.35	42.72 \pm 5.68	44.32 \pm 5.73	81.77 \pm 4.31	75.66 \pm 1.01	68.79 \pm 1.4	68.86 \pm 1.22
FEDBN [36]	Batch norm.	11.18M	11.17M	90.37 \pm 5.19	43.18 \pm 8.67	35.01 \pm 7.24	36.29 \pm 7.43	81.86 \pm 5.13	74.26 \pm 0.52	68.74 \pm 1.17	68.83 \pm 1.08
FEDALT [48]	Input layer	11.18M	6.45M	87.07 \pm 6.54	32.23 \pm 8.23	27.49 \pm 6.41	28.51 \pm 7.11	81.07 \pm 5.59	65.85 \pm 0.9	67.63 \pm 1.18	67.74 \pm 1.1
FEDSIM [48]	Input layer	11.18M	6.45M	87.93 \pm 6.26	33.07 \pm 8.16	28.21 \pm 6.41	29.15 \pm 7.16	82.45 \pm 5.03	67.66 \pm 0.82	67.49 \pm 1.32	67.54 \pm 1.24
LG-FEDAVG [38]	Feat. extractor	11.18M	0.005M	86.7 \pm 8.01	29.96 \pm 8	25.97 \pm 6.21	26.83 \pm 6.95	82.04 \pm 5.96	63.57 \pm 2.32	65.78 \pm 1.62	66.23 \pm 1.75
FEDREP [9]	Output layer	11.18M	11.17M	87.76 \pm 6.46	35.19 \pm 6.97	30.15 \pm 5.89	30.68 \pm 6.31	79.05 \pm 5.88	74.17 \pm 2.02	66.66 \pm 1.82	66.52 \pm 1.47
FEDALT [48]	Output layer	11.18M	11.17M	89.68 \pm 5.4	40.68 \pm 7.3	33.61 \pm 6.12	34.3 \pm 6.5	83.24 \pm 3.96	70.62 \pm 1.46	68.27 \pm 1.3	68.36 \pm 1.31
FEDSIM [48]	Output layer	11.18M	11.17M	89.75 \pm 6.51	41.98 \pm 7.66	34.21 \pm 6.22	35.31 \pm 6.79	82.91 \pm 4.46	72.34 \pm 0.51	68.22 \pm 1.34	68.12 \pm 1.24
FEDALT [48]	Adapter	12.59M	1.41M	87.26 \pm 7.78	31.51 \pm 8.55	27.38 \pm 6.65	27.77 \pm 7.19	81.41 \pm 6.5	57.88 \pm 3.57	72.13 \pm 1.34	74.67 \pm 1.57
FEDSIM [48]	Adapter	12.59M	1.41M	87.76 \pm 7.57	31.97 \pm 7.44	27.76 \pm 5.78	28.1 \pm 6.46	82.14 \pm 5.46	58.62 \pm 3.24	71.75 \pm 1.4	74.09 \pm 1.65
PERADA w/o KD	Adapter	2.82M	1.41M	91.27 \pm 5.15	53.81 \pm 6.27	42.5 \pm 5.06	44.45 \pm 5.48	83.31 \pm 4.31	76.55 \pm 2.47	76.77 \pm 2.24	77.59 \pm 2.18
PERADA	Adapter	2.82M	1.41M	91.82 \pm 4.43	59.05 \pm 5.24	47.25 \pm 4.48	48.53 \pm 4.74	83.58 \pm 4.74	77.2 \pm 1.63	76.98 \pm 3.87	77.88 \pm 1.55

- PerAda achieves the highest personalized performance and generalization by updating the smallest number of model parameters
- Existing partial model personalization methods have poor generalization to distribution shifts.
- Adapter-based personalization methods are generally effective on CheXpert.

Comparison with SOTA

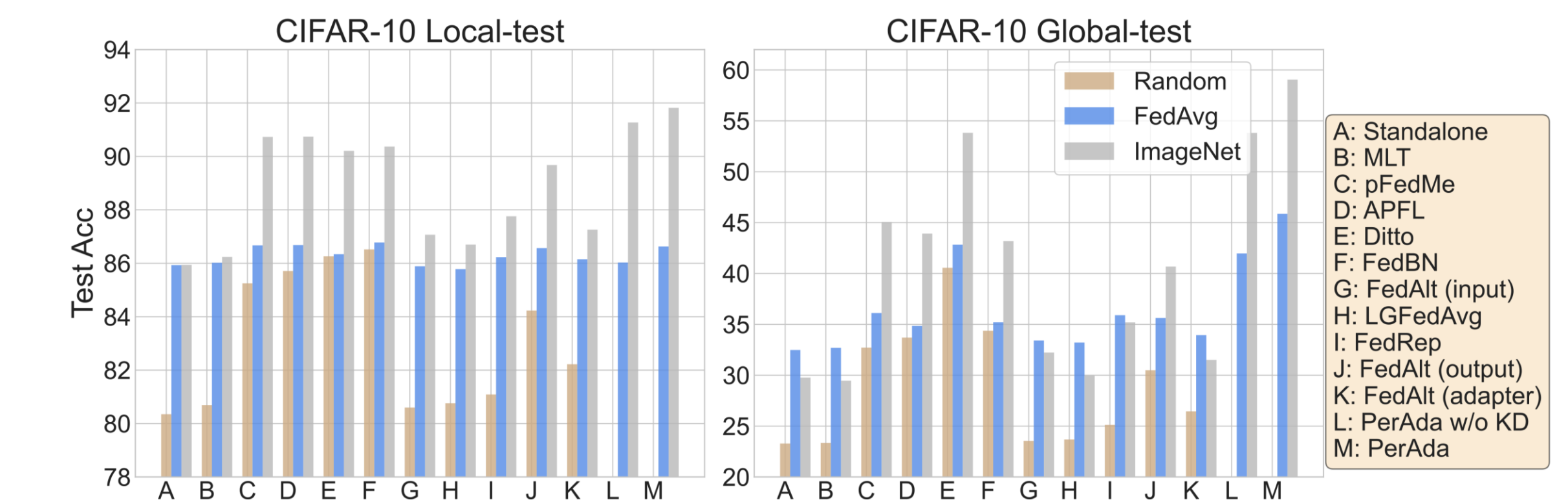


Effect of knowledge distillation



- More distillation steps & data help
- Out-of-domain distillation data achieve similar performance as in-domain data

Effect of pretrained models



- ImageNet-pretraining leads to better performance than Fed-Avg pretraining

Utility under differential privacy (DP) guarantees

CIFAR-10 dataset with ViT-S/16-224 model

Algorithm	Personalization	$\epsilon = \infty$	$\epsilon = 5.99 \pm 3.03$	$\epsilon = 3.7 \pm 2.12$	$\epsilon = 1.81 \pm 1.12$
Ditto	Full	98.59 \pm 1.63	76.76 \pm 24.14	76.75 \pm 24.13	76.67 \pm 24.12
PERADA w/o KD	Adapter	97.69 \pm 1.79	77.49 \pm 21.21	77.32 \pm 21.16	76.68 \pm 21
PERADA	Adapter	98.08 \pm 1.28	80.33 \pm 20.76	79.79 \pm 20.45	77.83 \pm 19.58

- Perform local training with DP-SGD for personalized & global model
- PerAda achieves higher utility than full model personalization under DP