

Towards Diverse Device Heterogeneous Federated Learning via Task Arithmetic Knowledge Integration

Mahdi Morafah¹, Vyacheslav Kungurtsev², Hojin Chang¹, Chen Chen³,
Bill Lin¹

¹University of California, San Diego

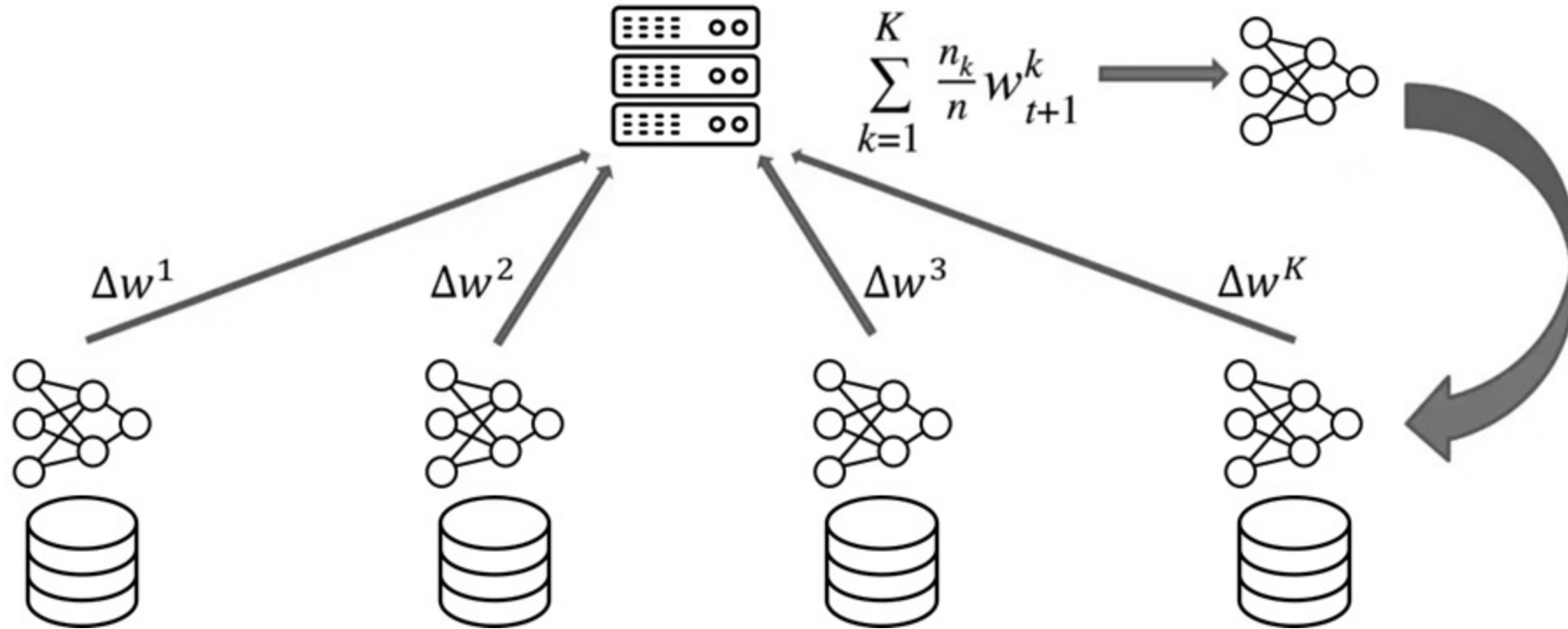
²Czech Technical University in Prague

³University of Central Florida

NeurIPS 2024



Standard Federated Learning (FedAvg)



- Assumes clients can train an identical model
- Clients train on local data, weights averaged at the server
- Global shared model benefits from all client data

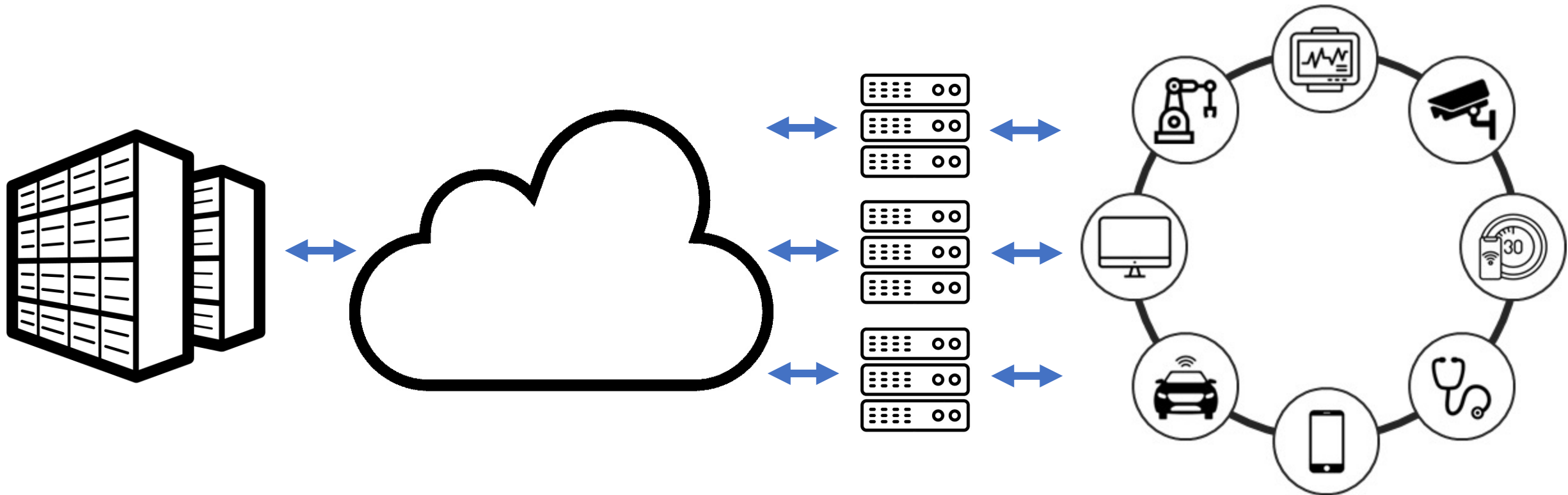
Heterogeneous Device Prototypes

DATA CENTERS

CLOUD

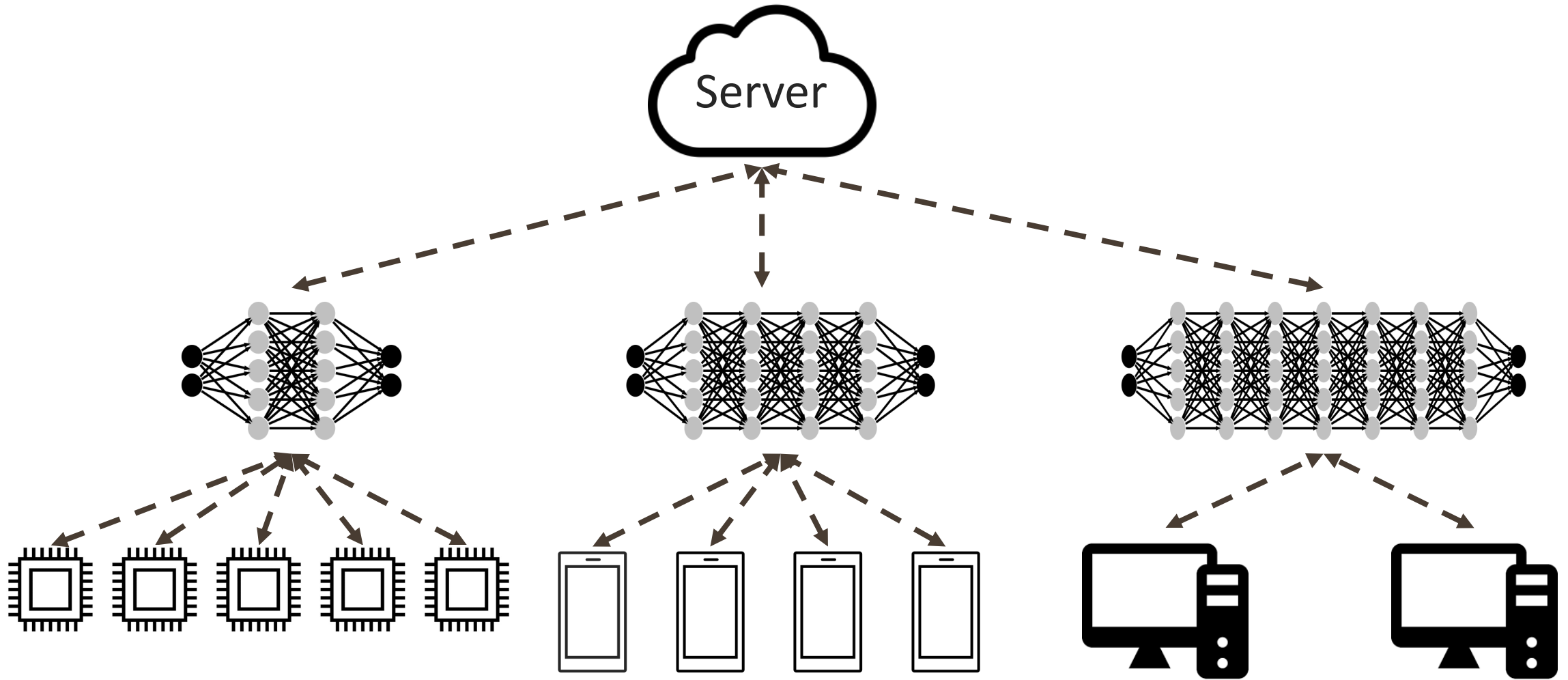
EDGE NODES

EDGE DEVICES

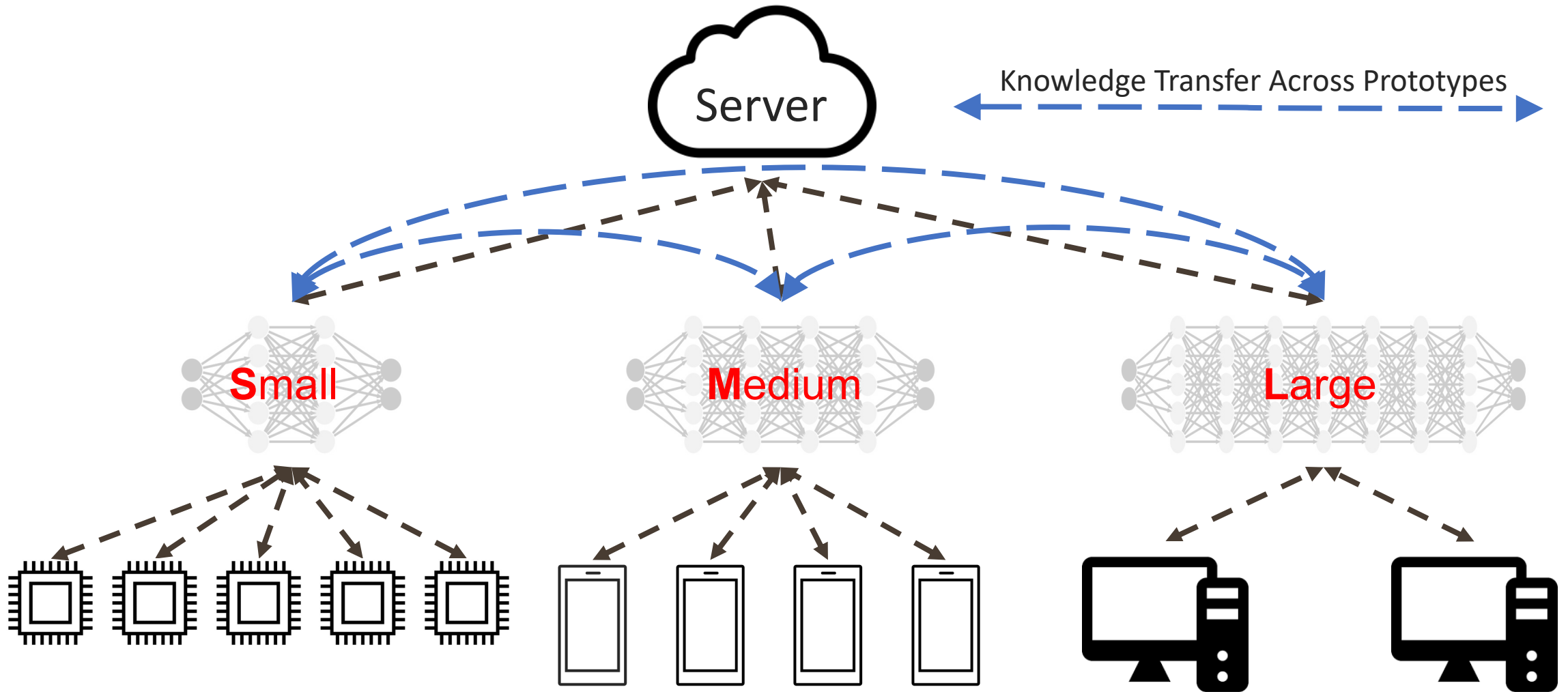


- Data for ML is distributed across diverse devices from small to large
- Learn ML models from diverse devices while maintaining data privacy

FL with Heterogeneous Device Prototypes



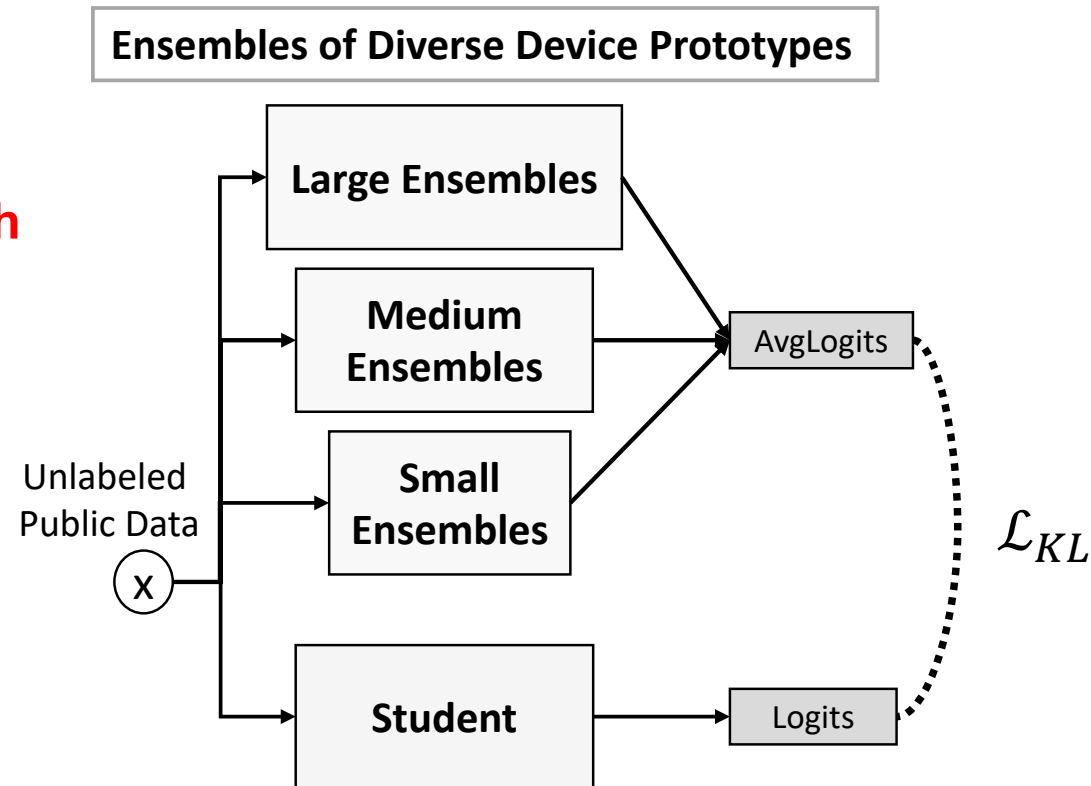
FL with Heterogeneous Device Prototypes



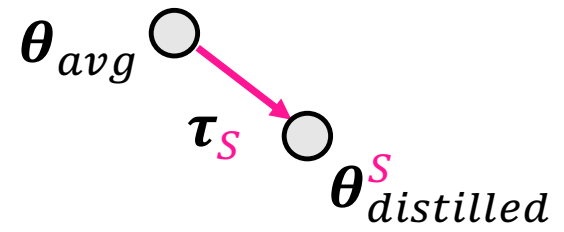
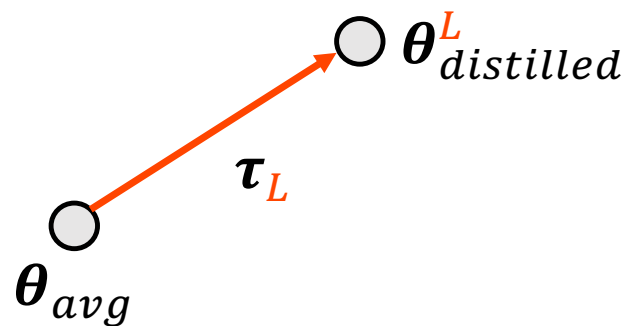
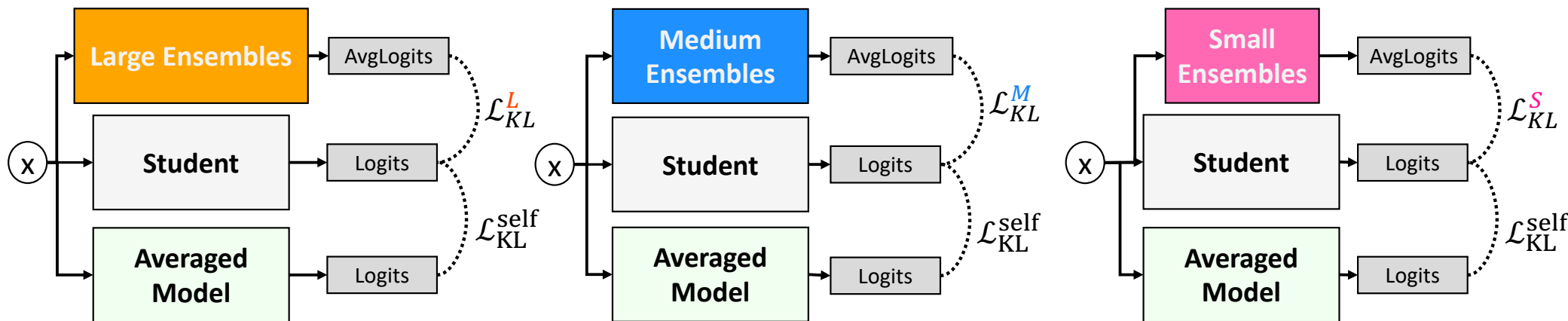
Limitations of Existing KD-based Methods

Information Dilution

one-size-fits-all approach



Proposed Method: TAKFL



Distillation process for Large Ensembles: $\theta_{avg} \xrightarrow{\tau_L} \theta_{distilled}^L$

Distillation process for Medium Ensembles: $\theta_{avg} \xrightarrow{\tau_M} \theta_{distilled}^M$

Distillation process for Small Ensembles: $\theta_{avg} \xrightarrow{\tau_S} \theta_{distilled}^S$

Distillation process for Large Ensembles: $\theta_{avg} \xrightarrow{\lambda_L \tau_L} \theta_{merged}$

Distillation process for Medium Ensembles: $\theta_{avg} \xrightarrow{\lambda_M \tau_M} \theta_{merged}$

Distillation process for Small Ensembles: $\theta_{avg} \xrightarrow{\lambda_S \tau_S} \theta_{merged}$

Distillation process for Large Ensembles: $\tau = \theta_{distilled} - \theta_{avg}$

Distillation process for Large Ensembles: $\theta_{merged} = \theta_{avg} + (\lambda_L \tau_L + \lambda_M \tau_M + \lambda_S \tau_S)$

Experiment Results (CV Task)

Table 4: **Performance Results for CV task on CIFAR-10 and CIFAR-100.** Training data is distributed among S, M, and L device prototypes in a 1:3:6 ratio, subdivided among clients using Dirichlet distribution. Public datasets are CIFAR-100 for CIFAR-10 and ImageNet-100 for CIFAR-100. Client configurations include 100, 20, and 4 clients for S, M, and L, with sampling rates of 0.1, 0.2, and 0.5. In homo-family settings, architectures are ResNet8, ResNet14, and ResNet18; in hetero-family settings, they are ViT-S, ResNet14, and VGG-16. All models are trained from scratch for 60 rounds. See Appendix F.1 for more details.

Homo-family Architecture Setting									
Dataset	Baseline	Low Data Heterogeneity				High Data Heterogeneity			
		S	M	L	Average	S	M	L	Average
CIFAR-10	FedAvg	36.21 \pm 2.24	46.41 \pm 2.33	59.46 \pm 6.17	47.36	22.01 \pm 0.78	25.26 \pm 3.89	51.51 \pm 3.52	32.93
	FedDF	49.31 \pm 0.15	50.63 \pm 0.73	49.82 \pm 0.98	49.92	34.71 \pm 1.48	35.27 \pm 4.74	51.08 \pm 4.04	40.35
	FedET	49.21 \pm 0.72	55.01 \pm 1.81	53.60 \pm 6.47	52.61	29.58 \pm 3.00	30.96 \pm 4.70	45.53 \pm 6.46	35.36
	TAKFL	55.90 \pm 1.70	57.93 \pm 3.49	60.58 \pm 2.35	58.14	37.40 \pm 1.68	38.96 \pm 0.17	51.49 \pm 6.15	42.61
	TAKFL+Reg	56.37\pm0.46	58.60\pm0.43	65.69\pm1.28	60.22	40.51\pm1.05	40.12\pm1.24	53.24\pm2.51	44.62
CIFAR-100	FedAvg	13.22 \pm 0.14	21.39 \pm 1.11	29.47 \pm 0.86	21.36	11.86 \pm 0.08	14.63 \pm 0.65	26.25 \pm 1.64	17.58
	FedDF	19.54 \pm 0.20	24.32 \pm 0.45	29.29 \pm 1.45	24.38	16.09 \pm 0.32	19.80 \pm 0.17	26.59 \pm 0.25	20.83
	FedET	19.67 \pm 0.35	25.27 \pm 0.66	31.10 \pm 1.53	25.35	11.18 \pm 1.68	18.22 \pm 0.35	26.40 \pm 0.65	18.60
	TAKFL	24.48 \pm 0.42	27.60 \pm 0.25	29.84 \pm 0.94	27.31	22.90\pm0.18	23.63 \pm 0.72	26.98 \pm 0.13	24.50
	TAKFL+Reg	27.18\pm0.27	29.14\pm0.20	31.15\pm0.97	29.16	22.88 \pm 0.37	23.92\pm0.57	28.01\pm0.34	24.94
Hetero-family Architecture Setting									
Dataset	Baseline	Low Data Heterogeneity				High Data Heterogeneity			
		S	M	L	Average	S	M	L	Average
CIFAR-10	FedAvg	27.53 \pm 0.83	47.30 \pm 3.17	55.10 \pm 8.60	43.31	20.93 \pm 1.54	25.62 \pm 6.04	36.80 \pm 5.47	27.78
	FedDF	34.15 \pm 0.87	54.06 \pm 1.06	69.07 \pm 4.99	52.43	24.20 \pm 0.74	34.07 \pm 3.08	39.81 \pm 5.45	32.69
	FedET	33.24 \pm 1.27	58.86 \pm 0.94	65.56 \pm 3.49	52.55	24.37 \pm 1.26	37.77 \pm 4.71	43.64 \pm 3.36	35.26
	TAKFL	33.29 \pm 0.15	57.64 \pm 0.19	68.44 \pm 0.66	53.12	24.92 \pm 1.32	38.07 \pm 3.19	48.01 \pm 3.99	37.00
	TAKFL+Reg	33.34\pm3.36	59.01\pm3.12	70.22\pm4.40	54.19	25.10\pm1.87	38.81\pm5.36	50.26\pm6.42	38.06
CIFAR-100	FedAvg	8.51 \pm 0.37	22.11 \pm 0.58	37.91 \pm 2.60	22.84	7.01 \pm 0.47	14.94 \pm 0.96	28.51 \pm 1.46	16.82
	FedDF	10.46 \pm 0.17	23.46 \pm 0.65	36.81 \pm 0.82	23.58	7.76 \pm 0.40	18.92 \pm 0.39	29.81 \pm 1.09	18.83
	FedET	11.16 \pm 0.18	25.40 \pm 0.30	37.38 \pm 0.60	24.65	8.20 \pm 0.54	20.66 \pm 0.50	28.95 \pm 1.79	19.27
	TAKFL	10.29 \pm 0.11	27.14 \pm 0.89	39.15\pm0.88	25.53	7.88 \pm 0.68	21.41 \pm 0.37	31.31 \pm 0.66	20.20
	TAKFL+Reg	11.25\pm0.37	27.86\pm0.86	38.68 \pm 0.45	25.93	8.45\pm0.20	22.16\pm0.87	31.95\pm1.13	20.85

Experiment Results (NLP Task)

Table 6: **Performance Results for NLP Task on 4 Datasets.** Training data is distributed among S, M, and L device prototypes in a 1:3:6 ratio, subdivided among clients using Dir(0.5). Client configurations are 8, 4, and 2 clients for S, M, and L, with sample rates of 0.3, 0.5, and 1.0, respectively. Architectures include Bert-Tiny, Bert-Mini, and Bert-Small for S, M, and L, initialized from pre-trained parameters and fine-tuned for 20 communication rounds. See Appendix F.2 for more details.

Private	Public	Baseline	S	M	L	Average
MNLI	SNLI	FedAvg	36.15 \pm 0.46	54.47 \pm 2.48	57.51 \pm 2.79	49.37
		FedDF	54.21 \pm 0.15	60.44 \pm 1.91	66.71 \pm 1.09	60.45
		FedET	48.03 \pm 6.32	50.33 \pm 7.87	53.80 \pm 6.18	50.72
		TAKFL	57.43 \pm 0.21	63.58 \pm 0.31	68.74 \pm 0.12	63.25
		TAKFL+Reg	57.61\pm0.89	63.91\pm1.05	68.96\pm1.10	63.49
SST2	Sent140	FedAvg	54.98 \pm 1.81	74.71 \pm 8.22	86.69 \pm 0.06	72.13
		FedDF	74.41 \pm 2.62	80.71 \pm 1.63	84.35 \pm 1.66	79.82
		FedET	66.63 \pm 9.14	65.89 \pm 16.35	70.05 \pm 15.83	67.52
		TAKFL	74.73 \pm 0.55	82.17 \pm 0.31	86.93 \pm 0.42	81.28
		TAKFL+Reg	74.88\pm0.43	82.40\pm0.83	87.33\pm0.63	81.54
MARC	Yelp	FedAvg	33.76 \pm 1.13	49.08 \pm 1.28	59.26 \pm 1.43	47.36
		FedDF	53.01 \pm 1.24	55.37 \pm 0.87	56.81 \pm 0.99	55.06
		FedET	52.63 \pm 2.29	54.28 \pm 2.31	56.11 \pm 2.61	54.34
		TAKFL	55.70 \pm 2.08	58.64 \pm 1.75	59.39 \pm 1.16	57.91
		TAKFL+Reg	55.96\pm1.66	59.18\pm1.13	59.61\pm1.89	58.25
AG-News	DBPedia	FedAvg	83.64 \pm 3.51	83.47 \pm 2.35	91.48 \pm 2.22	86.20
		FedDF	85.97 \pm 2.45	89.10 \pm 1.85	91.37 \pm 1.10	88.81
		FedET	75.27 \pm 3.85	81.13 \pm 3.21	83.19 \pm 4.58	79.86
		TAKFL	87.37 \pm 1.31	90.11 \pm 1.56	92.48 \pm 1.12	89.99
		TAKFL+Reg	87.66\pm1.83	90.30\pm2.05	92.61\pm1.72	90.19

Experiment Results (Scalability)

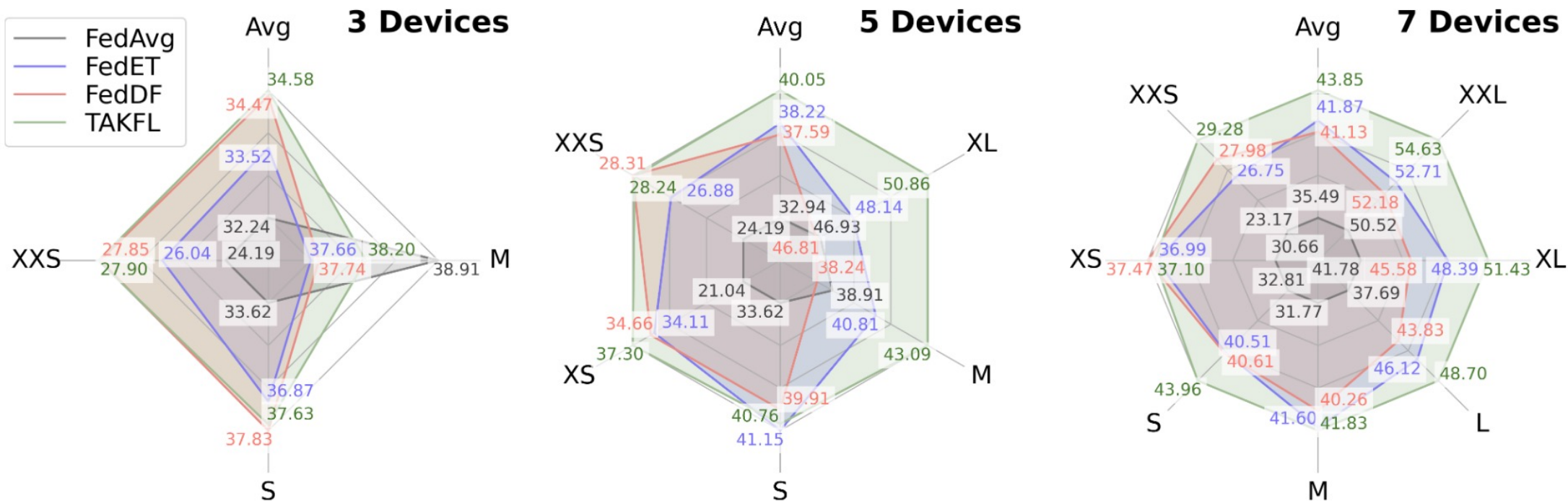


Figure 3: **Scalability Evaluation of TAKFL.** Image classification on CINIC-10 [9] dataset is used to evaluate TAKFL’s scalability across device prototypes ranging from XXS to XXL. Training data is distributed among prototypes in a 1:2:3:4:5:6:7 ratio, further subdivided using Dir(0.5). Client configurations range from 35 for XXS to 5 for XXL. Architectures span from ResNet10-XXS for XXS to ResNet50 for XXL prototype, all initialized from scratch and trained over 30 communication rounds. The public dataset is CIFAR-100 [24]. See Appendix D.4 for more details.

Thank you for listening!

Feel free to reach out if you have any questions!