

PERSONALIZING FEDERATED LEARNING WITH OVER-THE-AIR COMPUTATIONS

Zihan Chen^{†*} Zeshen Li^{†*} Howard H. Yang[‡] Tony Q.S. Quek[†]

[†] Singapore University of Technology and Design, Singapore 487372

[‡] ZJU-UIUC Institute, Zhejiang University, Haining 314400, China

ABSTRACT

Federated edge learning is a promising technology to deploy intelligence at the edge of wireless networks in a privacy-preserving manner. Under such a setting, multiple clients collaboratively train a global generic model under the coordination of an edge server. But the training efficiency is often throttled by challenges arising from limited communication and data heterogeneity. This paper presents a distributed training paradigm that employs analog over-the-air computation to address the communication bottleneck. Additionally, we leverage a bi-level optimization framework to personalize the federated learning model so as to cope with the data heterogeneity issue. As a result, it enhances the generalization and robustness of each client’s local model. We elaborate on the model training procedure and its advantages over conventional frameworks. We provide a convergence analysis that theoretically demonstrates the training efficiency. We also conduct extensive experiments to validate the efficacy of the proposed framework.

Index Terms— Federated learning, personalization, wireless edge network, over-the-air computation, robustness.

1. INTRODUCTION

With the increasing concerns on data privacy as well as the rapid growing capability of edge devices, deploying the federated learning (FL) [1] at the edge of wireless network, commonly coined as *federated edge learning* (FEEL), is attracting arising attentions [2, 3], where the computation tasks could be decoupled from the cloud to the edge of the network in a privacy-preserving paradigm.

However, in real-world implementations of the FEEL system, a typical training process of a generic global model requires hundreds of communication rounds among the massively distributed clients. The iterative gradient exchange would bring hefty communication overhead [1, 4]. Hence, for a digital communication based-FEEL system run over the resource-constrained network, the limited communication bandwidth would inevitably constrain the scalability,

since every selected client in each round requires an assigned orthogonal sub-channel to perform the update [5, 6].

To combat the communication bottleneck, an array of recent studies [6–11] suggest incorporate *analog over-the-air* (A-OTA) computations into the design of FEEL systems, exploiting the superposition property of the multi-access channels for fast and scalable model aggregations. The adoption of the A-OTA computations with FEEL, termed as *A-OTA-FEEL*, have been demonstrated to have high spectral efficiency, low access latency, enhanced privacy protection, and reduced communication costs [5, 6, 12], all benefiting from the automatic “one-shot” gradient aggregation for model update [7, 13]. Nevertheless, A-OTA computations inevitably introduce the random channel fading and interference into the aggregated gradients, leading to performance degradations such as the slower convergence and instability [5, 10]. Hence, robust training techniques could be adopted to enhance the performance with channel imperfections.

In addition to the inherent channel fading and interference, current approaches for A-OTA-FEEL system design have not addressed the existing discrepancies in both local data statistics and qualities (i.e., data heterogeneity and label noise) due to the diverse preferences, bias, and hardware capabilities of different clients [11, 14]. Such discrepancies in clients’ datasets can significantly degrade the FL performance. More crucially, these discrepancies would even make the single generic global model fail to achieve good generalization and robustness performance on diverse local data [15–17]. On the other hand, the future intelligent network is envisioned to be able to provide customized services to the clients [18, 19]. It is necessary to address the individuality of the clients in the design of the A-OTA-FEEL system with personalized intelligent services.

In view of the above challenges, we propose a personalized training framework in the context of the A-OTA-FEEL. The proposed framework provides personalized model training services while still enjoying the benefits of analog over-the-air computations, in which each client would maintain two models (i.e., generic and personalized models) at the local via two different global and local objectives. We also provide a convergence analysis of the proposed personalized A-OTA-FEEL framework. Both the theoretical and numerical results validate the gain from the personalization design.

*Equal contribution. This work was supported in part by the National Natural Science Foundation of China under Grant 62271513. (*Corresponding Author: Howard H. Yang*)

2. SYSTEM MODEL

We consider a wireless system consisting of one edge server that is attached to an access point and K clients, where the i -th device has a local dataset \mathcal{D}_i . In this system, communications between the clients and the server are taken place over the spectrum. Each client's goal is to (a) learn a statistical model based on its own dataset and (b) exploit information from the dataset of other clients and, aided by the orchestration of the server, attain an improvement toward its locally learned model while preserving privacy. Such tasks can be achieved via a bi-level optimization based PFL framework. More precisely, every client k aims to find a local model $\mathbf{v}_k \in \mathbb{R}^d$ that solves the following *personal objective* function

$$\min_{\mathbf{v}_k} f_k(\mathbf{v}_k; \mathbf{w}^*) = F_k(\mathbf{v}_k) + \frac{\lambda}{2} \|\mathbf{v}_k - \mathbf{w}^*\|^2 \quad (1)$$

$$\text{s.t.} \quad \mathbf{w}^* \in \arg \min_{\mathbf{w}} \frac{1}{K} \sum_{i=1}^K F_i(\mathbf{w}) \quad (2)$$

in which $F_i(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes the loss function of client i , $\mathbf{w} \in \mathbb{R}^d$ is a globally trained generic model, and λ is a hyper-parameter that controls the level of personalization of the clients' locally trained personal models. We use η_l to denote the learning rate in the optimization of personal objective. Notably, a large value of λ indicates that the clients' local models $\{\mathbf{v}_i\}_{i=1}^K$ need to well align with the global model \mathbf{w}^* , promoting commonality across the local models. In contrast, a small λ improves personalization. Moreover, benefiting from such a bi-level optimization design, the personalized local models $\{\mathbf{v}_i\}_{i=1}^K$ would have better generalization and robustness performance on the limited local data.

To solve the above optimization problem, the clients need to not just train their local models through (1), but more importantly, jointly minimize a global objective function as per (2). Due to privacy concerns, the clients will carry out the minimization problem (2) without sharing data in an FL manner. The following section presents a model training approach that capitalizes on the properties of analog transmissions for low-latency and high-privacy federated computing.

3. MODEL TRAINING PROCEDURE

This section details the PFL model training process based on over-the-air computing schemes. (See Fig. 1 for an overview.) More precisely, we employ A-OTA computations for fast (and highly scalable) gradient aggregation that significantly improves the training efficiency of the global model. The detailed training procedure is elaborated on below.

1) *Local Model Training*: Without loss of generality, we assume the system has progressed to the t -th round of global training, where the clients just received the global model parameters \mathbf{w}^t from the edge server.¹ Then, each client k up-

¹Because of the high transmit power of the access point, we assume the global model can be successfully received by all the clients.

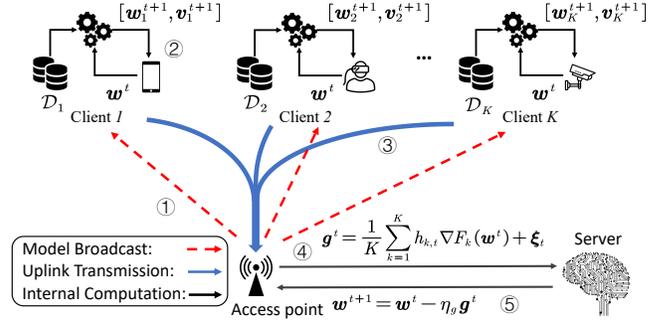


Fig. 1. An overview of personalized analog over-the-air federated edge learning, in which each client maintains a common global model and local personalized model.

dates its personalized local model \mathbf{v}_k^t by optimizing the local personal objective function $f_k(\mathbf{v}_k; \mathbf{w}^t)$. (For simplicity, we use \mathbf{v}_k to denote personal model.) Each client k also computes its local gradient $\nabla F_k(\mathbf{w}^t)$ for global model update.

2) *Analog Gradient Aggregation*: We consider the clients adopt analog transmissions to upload their locally trained parameters. Specifically, once $\nabla F_k(\mathbf{w}^t)$ is computed, client k modulates it entry-by-entry onto the magnitudes of a common set of orthogonal baseband waveforms [5], forming the following analog signal

$$x_k(s) = \langle \mathbf{u}(s), \nabla F_k(\mathbf{w}^t) \rangle \quad (3)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product between two vectors and $\mathbf{u}(s) = (u_1(s), \dots, u_d(s))$, $s \in [0, \tau]$ has its entries satisfying

$$\int_0^\tau u_i^2(s) ds = 1, \quad i = 1, 2, \dots, d \quad (4)$$

$$\int_0^\tau u_i(s) u_j(s) ds = 0, \quad i \neq j. \quad (5)$$

τ represents the total time of signal duration. Once the analog waveforms $\{x_k(s)\}_{k=1}^K$ are available, the clients transmit them concurrently to the access point. Owing to the superposition property of electromagnetic waves, the signal received at the radio front end of the access point can be expressed as:

$$y(s) = \sum_{k=1}^K h_{k,t} P_k x_k(s) + \xi(s), \quad (6)$$

where $h_{k,t}$ is the channel fading experienced by client k , P_k the corresponding transmit power, and $\xi(s)$ denotes the additive noise. In this work, we assume the channel fading is i.i.d. across clients, with mean μ_h and variance σ_h^2 . Besides, the transmit power of each client is set to compensate for the large-scale path loss and we use P to denote the average power for all clients. This received signal will be passed through a bank of match filters, with each branch tuning to

Algorithm 1 Personalized A-OTA FEEL framework

Input: Initial global model \mathbf{w}^0 , initial personal local models

$$\{\mathbf{v}_i\}_{k=1}^K, T, \lambda, \eta_g$$

Output: Global model \mathbf{w}^T , personal model $\{\mathbf{v}_i\}_{i=1}^K$

- 1: **for** $t = 0, 1, 2$ **to** $T - 1$ **do**
 - 2: **for** $k = 1, 2,$ **to** K **in parallel do**
 # global generic model update
 - 3: $\nabla F_i(\mathbf{w}^t) \leftarrow \text{CLIENTUPDATE}(k, \mathbf{w}^t)$
 # local personalized model update
 - 4: Update \mathbf{v}_k via solving $f_k(\mathbf{v}_k; \mathbf{w}^t)$
 - 5: Transmit local gradient $\nabla F_i(\mathbf{w}^t)$ to edge server
 # Noisy aggregation via analog OTA computations
 - 6: $\mathbf{g}^t = \frac{1}{K} \sum_{k=1}^K h_{k,t} \nabla F_k(\mathbf{w}^t) + \boldsymbol{\xi}_t$
 - 7: $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_g \mathbf{g}^t$ # Global model update
 - 8: **return** $\mathbf{w}^T, \{\mathbf{v}_i\}_{k=1}^K$
-

$u_i(s), i = 1, 2, \dots, d$. On the output side, the server obtains:

$$\mathbf{g}^t = \frac{1}{K} \sum_{k=1}^K h_{k,t} \nabla F_k(\mathbf{w}^t) + \boldsymbol{\xi}_t, \quad (7)$$

in which $\boldsymbol{\xi}_t$ is a d -dimensional random vector with each entry being i.i.d. and follows a zero-mean Gaussian distribution with variance σ^2 . It is noteworthy that the vector given in (7) is a distorted version of the globally aggregated gradient.

3) *Global Generic Model Update:* Using \mathbf{g}^t , the server updates the global model as follows:

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta_g \mathbf{g}^t, \quad (8)$$

where η_g is the learning rate for generic global model update. After this, the server broadcasts the \mathbf{w}^{t+1} to all clients for the next round local computing. Such a process will be iterated through multiple rounds until the global model converges.

Notably, the bi-level optimization in the personal model mitigates impacts from the random channel fading and noise introduced by analog over-the-air computations to the globally aggregated gradient, thus improving the robustness of the analog over-the-air federated edge learning system. Consequently, personalization enhances both the generalization and robustness of the FL system in the presence of data heterogeneity and noisy model aggregation.²

We summarize the proposed framework in Algorithm 1. It is worthwhile to highlight several advantages of the presented framework, including high scalability, low access latency, enhanced privacy, better generalization as well as robustness, brought together by analog over-the-air computations and personalized training. We would also like to address that we make no assumption about the generic model

²This paper does not consider the architecture-based PFL methods in which each client maintains a personal model with unique architecture via techniques such as sparsification or model weight decoupling [15]. It would increase the cost of the synchronizations for signal transmission to achieve automatic signal aggregations in the context of A-OTA computations.

training, as well as the OTA communication, which indicates that the performance could be further enhanced by advanced federated optimization [4] and OTA techniques [8, 12].

4. CONVERGENCE ANALYSIS

This section provides the convergence analysis of our proposed framework from the perspective of both the global model and the local personalized model.

To facilitate the analysis, we assume that each client's loss function is μ -strongly convex and the local gradient $\nabla F_k(\mathbf{w}^t)$ is Lipschitz continuous with constant $L_k > 0$. We use \bar{L} to denote the maximal constant among all clients and L is the Lipschitz gradient constant of global objective. δ is the diameter of the compact convex parameter set that all model parameters lie in. We consider that if the global model converges, its convergence rate is denoted by $g(t)$, i.e., there exists $g(t)$ that $\lim_{t \rightarrow \infty} g(t) = 0$ and $E[|\mathbf{w}^t - \mathbf{w}^*|^2] \leq g(t)$. In this work, we denote by \mathbf{v}_k^* and \mathbf{z}_k^* as $\mathbf{v}_k^* = \arg \min_{\mathbf{v}} f_k(\mathbf{v}; \mathbf{w}^*)$ and $\mathbf{z}_k^* = \arg \min_{\mathbf{z}} F_k(\mathbf{z})$, respectively. We assume the l_2 distance between the optimal local and global model is bounded, i.e., for any $k \in [K]$, $\|\mathbf{z}_k^* - \mathbf{w}^*\| \leq M$.

We now present the main theoretical finding of this paper. First of all, the following theorem provides the convergence rate of the global generic model.

Theorem 4.1. *Under the considered A-OTA FEEL system, let $r_0^2 \triangleq \|\mathbf{w}^0 - \mathbf{w}^*\|^2$ be the squared distance between the initial estimate \mathbf{w}_0 and \mathbf{w}^* . If the learning rate η_g satisfies*

$$0 < \eta_g < \min \left\{ \frac{2}{\mu_h(\mu + L)}, \frac{2\mu_h\mu LK}{\sigma_h^2 \bar{L}^2(1 + 2\delta)(\mu + L)} \right\}, \quad (9)$$

then the error of \mathbf{w}^t can be bounded as:

$$E[|\mathbf{w}^t - \mathbf{w}^*|^2] \leq c^t r_0^2 + \frac{\eta_g^2}{(1-c)} \left(\frac{\sigma_h^2 \delta \bar{L}^2(2+\delta)}{K} + \frac{d\sigma^2}{P^2 K^2} \right) \quad (10)$$

where $0 < c \triangleq 1 - \frac{2\eta_g\mu_h\mu L}{\mu + L} + \frac{\eta_g^2\sigma_h^2 \bar{L}^2(1+2\delta)}{K} < 1$.

Proof. Please refer to [10] for a detailed proof. \square

Next, we employ the following lemma to characterize the convergence rate of the local personalized models.

Lemma 4.2. *Under the considered system, let local learning rate satisfy condition (9), then the local model of client k converges as:*

$$\begin{aligned} E[|\mathbf{v}_k^{t+1} - \mathbf{v}_k^*|^2] &\leq (1 - \mu\eta_l) E[|\mathbf{v}_k^t - \mathbf{v}_k^*|^2] + \eta_l^2 \lambda^2 M^2 \\ &+ \eta_l^2 \lambda^2 E[|\mathbf{w}^t - \mathbf{w}^*|^2] + 2\eta_l^2 \lambda^2 M \sqrt{E[|\mathbf{w}^t - \mathbf{w}^*|^2]} \\ &+ 2\eta_l \lambda \sqrt{E[|\mathbf{v}_k^t - \mathbf{v}_k^*|^2]} E[|\mathbf{w}^t - \mathbf{w}^*|^2]. \end{aligned} \quad (11)$$

Proof. Please refer to [16] for detailed proof. \square

Aided by the above result, we obtain the convergence rate of the global model as the following.

Theorem 4.3. *Under the considered A-OTA FEEL system, if there exists a variable A satisfying $\frac{g(t+1)}{g(t)} \geq 1 - \frac{g(t)}{A}$, then, there is a constant $C < \infty$ such that for any client k , $E \left[\|\mathbf{v}_k^t - \mathbf{v}_k^*\|^2 \right] \leq Cg(t)$ with a local learning rate given by $\eta = \frac{2g(t)}{A\mu}$.*

Proof. We omit the proof due to the space limit. \square

To this end, we can see that via A-OTA computing, both the global and local personalized models attain linear convergence rates, while addressing the non-ideal gradient updates.

5. NUMERICAL RESULTS

This section evaluates the performance of our proposed framework. Particularly, we examine the performance of the personalized local training in terms of generalization power and robustness compared to conventional settings and baselines. We also explore the robustness performance of the framework in the context of the noisy local data (i.e., part of the local training data are annotated with wrong labels).

5.1. Experiment setup

We evaluate our framework on image classification tasks on CIFAR-10/100 [20] with ResNet-18 and ResNet-34 [21], respectively. Both IID and non-IID data settings are considered, in which the non-IID data partitions are implemented with Dirichlet distribution and the identicalness of the distributions could be controlled by the parameter α . Unless otherwise specified, we use $K = 100$ for CIFAR-10, $K = 50$ for CIFAR-100, and Rayleigh fading with average channel gain $\mu_h = 1$. We select λ from comparison experiments. The federated label noise setting is the same as the [14]³.

5.2. Performance evaluation

We first compare the personalized models performance of our proposed framework with generic global model from conventional FL setup in Fig. 2 with IID local data partition, using the same configurations of A-OTA-FEEL system. The two sub-figures demonstrate the consistent outperformance of personalization training scheme. Specifically, increasing total number of clients in the system (i.e., a larger K) would improve the system performance for both two settings in A-OTA, and personalized training presents a more robust generalization with diverse local data quality.

³For all label noise settings, we use lower bound 0.5 for local label noise level. Details can be found in [14].

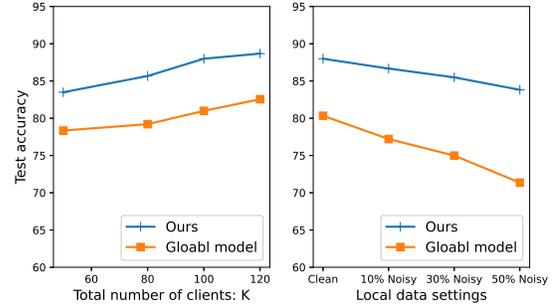


Fig. 2. Performance comparison of the best test accuracy on CIFAR-10 with IID data settings. (Left): Performance with different total number of clients K . (Right): Performance with different ratios of clients containing noisy local data.

Table 1. Average (3 trails) of the best test accuracy comparison on CIFAR-10/100 with real-world data settings. The highest accuracy for each setting is boldfaced.

	Methods	CIFAR-10		CIFAR-100
		$\alpha = 10$	$\alpha = 1$	$\alpha = 1$
Clean	OTA-FedAvg	76.32	72.71	65.12
	OTA-FedProx	76.45	72.90	66.35
	OTA-FedRep	82.44	79.93	-
	Ours	83.57	81.05	69.33
Noisy	OTA-FedAvg	70.51	67.15	58.81
	OTA-FedProx	72.06	69.62	59.29
	OTA-FedRep	77.09	74.23	-
	Ours	78.74	75.31	63.30

To further demonstrate the outperformance of the proposed framework, we provide the detailed best test accuracy comparison in Tab. 1 on CIFAR-10/100 with non-IID data, compared with FedAvg [1], FedProx [4] and FedRep [22] with same OTA setup. In such context, our proposed personalized training method achieves best test accuracies across all settings, which shows the superiority with respect to the generalization and robustness.

6. CONCLUSION

In this paper, we proposed a personalized A-OTA-FEEL framework that utilizes bi-level optimization and analog transmissions to address the data heterogeneity and communication efficiency challenges. Both the theoretical and empirical results were provided to demonstrate the effectiveness of the proposed framework. We highlighted the robustness performance of the PFL in edge learning. To the best of our knowledge, this is the first work that explores the PFL model in A-OTA FEEL systems. We envision that PFL could be a potential technique to provide customized services in future intelligent networks.

7. REFERENCES

- [1] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proc. Int. Conf. Artif. Intell. Stat.*, Fort Lauderdale, USA, Apr. 2017, pp. 1273–1282.
- [2] Howard H. Yang, Zuozhu Liu, Tony Q. S. Quek, and H. Vincent Poor, “Scheduling policies for federated learning in wireless networks,” *IEEE Trans. Commun.*, vol. 68, no. 1, pp. 317–333, 2020.
- [3] Walid Saad, Mehdi Bennis, and Mingzhe Chen, “A vision of 6g wireless systems: Applications, trends, technologies, and open research problems,” *IEEE Netw.*, vol. 34, no. 3, pp. 134–142, 2019.
- [4] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith, “Federated learning: Challenges, methods, and future directions,” *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [5] Howard H Yang, Zihan Chen, Tony Q.S. Quek, and H Vincent Poor, “Revisiting analog over-the-air machine learning: The blessing and curse of interference,” *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 3, pp. 406–419, 2021.
- [6] Kai Yang, Tao Jiang, Yuanming Shi, and Zhi Ding, “Federated learning via over-the-air computation,” *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2022–2035, 2020.
- [7] Huayan Guo, Yifan Zhu, Haoyu Ma, Vincent KN Lau, Kaibin Huang, Xiaofan Li, Huabin Nong, and Mingyu Zhou, “Over-the-air aggregation for federated learning: Waveform superposition and prototype validation,” *J. of Commun. and Inf. Netw.*, vol. 6, no. 4, pp. 429–442, 2021.
- [8] Li Chen, Nan Zhao, Yunfei Chen, F Richard Yu, and Guo Wei, “Over-the-air computation for iot networks: Computing multiple functions with antenna arrays,” *IEEE Internet Things J.*, vol. 5, no. 6, pp. 5296–5306, 2018.
- [9] Mohammad Mohammadi Amiri and Deniz Gündüz, “Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air,” *IEEE Trans. Signal Process.*, vol. 68, pp. 2155–2169, Mar. 2020.
- [10] Tomer Sery and Kobi Cohen, “On analog gradient descent learning over multiple access fading channels,” *IEEE Trans. Signal Process.*, vol. 68, pp. 2897–2911, 2020.
- [11] Zihan Chen, Zeshen Li, and Jingyi Xu, “Analog over-the-air federated learning with real-world data,” in *IEEE Int. Conf. on Sensing, Commun, and Netw. (SECON Workshops)*. IEEE, 2022, pp. 31–36.
- [12] Dongzhu Liu and Osvaldo Simeone, “Privacy for free: Wireless federated learning via uncoded transmission with adaptive power control,” *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 170–185, 2020.
- [13] Guangxu Zhu, Jie Xu, Kaibin Huang, and Shuguang Cui, “Over-the-air computing for wireless data aggregation in massive iot,” *IEEE Wireless Commun.*, vol. 28, no. 4, pp. 57–65, 2021.
- [14] Jingyi Xu, Zihan Chen, Tony Q.S. Quek, and Kai Fong Ernest Chong, “Fedcorr: Multi-stage federated learning for label noise correction,” in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2022.
- [15] Alysa Ziyang Tan, Han Yu, Lizhen Cui, and Qiang Yang, “Towards personalized federated learning,” *IEEE Trans. Neural Netw.*, 2022.
- [16] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith, “Ditto: Fair and robust federated learning through personalization,” in *Proc. Int. Conf. Mach. Learn.* PMLR, 2021, pp. 6357–6368.
- [17] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar, “Federated multi-task learning,” *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017.
- [18] Yiqing Zhou, Ling Liu, Lu Wang, Ning Hui, Xinyu Cui, Jie Wu, Yan Peng, Yanli Qi, and Chengwen Xing, “Service-aware 6g: An intelligent and open network based on the convergence of communication, computing and caching,” *Digit. Commun. Netw.*, vol. 6, no. 3, pp. 253–260, 2020.
- [19] Hao Tang, Zechao Li, Zhimao Peng, and Jinhui Tang, “Blockmix: Meta regularization and self-calibrated inference for metric-based meta-learning,” in *ACM Multimedia*, 2020, pp. 610–618.
- [20] Alex Krizhevsky and Geoffrey Hinton, “Learning multiple layers of features from tiny images,” Tech. Rep. 0, University of Toronto, Toronto, Ontario, 2009.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [22] Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai, “Exploiting shared representations for personalized federated learning,” in *Proc. Int. Conf. Mach. Learn.* PMLR, 2021, pp. 2089–2099.