# PERSONALIZED FEDERATED LEARNING ON LONG-TAILED DATA VIA ADVERSARIAL FEATURE AUGMENTATION

*Yang Lu[1], Pinxin Qian[1], Gang Huang[2], Hanzi Wang[1,*]*

[1]Fujian Key Laboratory of Sensing and Computing for Smart City, Xiamen University, Xiamen, China
[2]Zhejiang Lab, Hangzhou, Zhejiang, China
pxqian@stu.xmu.edu.cn, huanggang@zju.edu.cn, {luyang, hanzi.wang}@xmu.edu.cn

## ABSTRACT

Personalized Federated Learning (PFL) aims to learn personalized models for each client based on the knowledge across all clients in a privacy-preserving manner. Existing PFL methods generally assume that the underlying global data across all clients are uniformly distributed without considering the long-tail distribution. The joint problem of data heterogeneity and long-tail distribution in the FL environment is more challenging and severely affects the performance of personalized models. In this paper, we propose a PFL method called Federated Learning with Adversarial Feature Augmentation (FedAFA) to address this joint problem in PFL. FedAFA optimizes the personalized model for each client by producing a balanced feature set to enhance the local minority classes. The local minority class features are generated by transferring the knowledge from the local majority class features extracted by the global model in an adversarial example learning manner. The experimental results on benchmarks under different settings of data heterogeneity and long-tail distribution demonstrate that FedAFA significantly improves the personalized performance of each client compared with the state-of-the-art PFL algorithm. The code is available at https://github.com/pxqian/FedAFA.

***Index Terms***— Federated Learning, Long-Tail, Adversarial Sample, Feature Augmentation

## 1. INTRODUCTION

A common way to build a deep learning model is to collect the training data and train the model on a server, called centralized training. However, with the increasing data security awareness, such a centralized training paradigm is not applicable in some applications when sensitive data is stored in each data holder. Federated Learning (FL) develops a paradigm to train models without transmitting private data from each data holder (called client in FL) to a centralized server to address the privacy issue in deep learning [1, 2, 3, 4]. However, one major problem in FL is that the aggregated global model is usually not guaranteed to generalize overall clients well because each client's data distribution is not-independent and identically distributed (non-IID), which is usually defined as the *data heterogeneity problem* in FL. Therefore, targeting the model generalization ability on each client, Personalized Federated Learning (PFL) aims to obtain a 'tailored' local model that utilizes the global model's generalization ability and simultaneously fits the client's local data distribution.

Existing PFL methods [5, 6, 7, 8] have generally achieved promising personalized performance of each client on heterogeneous data. However, the data usually exhibits long-tail distribution in real-world scenarios, where a few head classes contain most samples while a large number of tail classes contain only a few samples. Therefore, it is reasonable to assume that the global data (data across all clients) is in long-tail distribution. In the PFL environment with data heterogeneity, the global tail class samples will be sporadically distributed on only a few clients, and each client's data distribution is also locally imbalanced. In addition, the local data distributions are likely to differ from the global data distribution, which yields different data imbalance degrees among clients. In this case, the generalization ability of PFL models will further deteriorate because the global tail classes are underrepresented by the aggregated global model, which mainly provides the knowledge of the global head classes. In addition, the PFL models are prone to overfitting the local minority classes with only a few samples in personalized optimization on each client. Some existing long-tail learning methods [9, 10, 11] are also not applicable in the PFL environment because the global data distribution is unknown to both the server and clients due to privacy issues.

To address the joint problem of data heterogeneity and long-tail distribution in PFL, we propose a PFL method called Federated Learning with Adversarial Feature Augmentation (FedAFA) that utilizes the global model learned across clients to rebalance the local feature set for robust personalized training. Specifically, inspired by the targeted adversarial attack [12, 13], we generate new local minority class features by adding specific small perturbations to the local majority class features. To not affect the performance of the original

---

local majority classes, we also propose a new optimization objective during personalized training. Experimental results on various heterogeneous and long-tailed benchmarks show that FedAFA surpasses the state-of-the-art PFL and long-tail learning methods.

## 2. PROPOSED METHOD

### 2.1. Basics of Personalized Federated Learning

In a typical PFL setting, there are $K$ clients participating in the training process and each client $k$ has a private dataset $\mathcal{D}_k = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{n_k}$, where $\mathbf{x}_i$ denotes the $i$th sample in $\mathcal{D}_k$, $\mathbf{y}_i \in \{0,1\}^C$ is the corresponding label over $C$ classes, and $n_k$ is the number of samples in $\mathcal{D}_k$. PFL aims to look for good local models (personalized models) for all $K$ clients, which are usually adapted from the global model $\mathbf{w}$. FedAvg-FT [14] is the most straightforward PFL method. It is a locally adaptive algorithm based on fine-tuning the global model by $\mathbf{w}_k = \mathbf{w} - \eta \nabla \mathcal{L}_k(\mathcal{D}_k; \mathbf{w})$ on $\mathcal{D}_k$, where $\mathcal{L}_k(\mathcal{D}_k; \mathbf{w})$ is the local training loss of the model $\mathbf{w}$ on $\mathcal{D}_k$. $\mathcal{L}_k(\mathcal{D}_k; \mathbf{w}) = \frac{1}{|\mathcal{D}_k|} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{D}_k} \ell(\mathbf{y}, f(\mathbf{x}; \mathbf{w}))$ is calculated by averaging sample losses in $\mathcal{D}_k$, where $\ell(\cdot, \cdot)$ is the sample loss and $f(\mathbf{x}; \mathbf{w})$ is the prediction result of sample $\mathbf{x}$ by model $\mathbf{w}$. However, this method slightly improves heterogeneous data because local model $\mathbf{w}_k$ is prone to overfit the local data $\mathcal{D}_k$, which results in poor local generalization performance.

### 2.2. FedAFA Framework

FedAFA is based on the following intuitions. First, the global model of FL is usually more robust than local models because it obtains the information from each client, although it is in an indirect manner. Second, it has been empirically shown that the feature extractor in a neural network is less affected by the data distribution compared with its classifier [15, 16, 17]. It means that the features extracted by the global model are still of high quality despite the negative influence of data heterogeneity and long-tail distribution. Based on the above observations, we propose to utilize the feature extractor of the global model to transfer information from the local majority classes to the local minority classes in the feature space.

We randomly select pairs of classes as the *source majority class* $y_s$ and the *target minority class* $y_t$ depending on their corresponding number of samples denoted as $n_s$ and $n_t$. We use the idea of adversarial samples [12, 13] to generate the features of the target minority class by adding specific perturbations to the features of source majority class samples. The perturbations are obtained from the gradients of the loss of predicting the source majority class sample into the target minority class. However, unlike generating adversarial samples, the input and output of FedAFA are in the feature space. After generating a certain number of features to achieve a balanced feature set for all local classes, we combine them with the original local samples with class-balanced sampling to train the local personalized model. Thus, the PFL model of
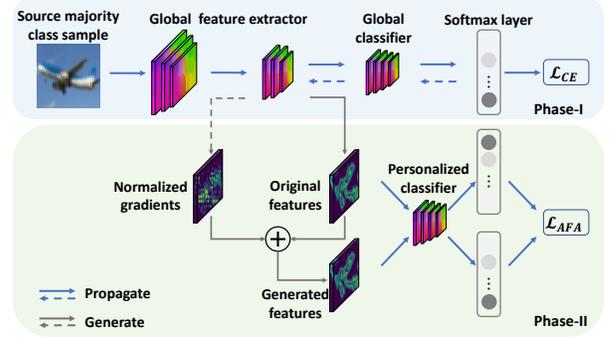


**Fig. 1.** An overview of personalized training of proposed FedAFA.

FedAFA is constructed by the feature extractor of the global model and the personalized classifier trained on the rebalanced feature set, which is illustrated in Figure 1. In order to explain FedAFA with better notations, we split the model into two parts: a feature extractor $g$ parameterized by $\mathbf{u}$ and a classifier $f$ parameterized by $\mathbf{v}$. Thus, the features of a sample $\mathbf{x}$ is generated by $\mathbf{h} = g(\mathbf{x}; \mathbf{u})$, and the prediction result is given by $f(\mathbf{h}; \mathbf{v})$. In FedAFA, client $k$ uses the feature extractor of the global model $\mathbf{u}$ and has its own personalized classifier with parameters $\mathbf{v}_k$. Next, two major steps in FedAFA, i.e., local feature augmentation and personalized model optimization, are described in detail.

#### 2.2.1. Local Feature Augmentation

First, we introduce how to adversarially generate features of the target minority class by utilizing the information of the source majority class. Given a target minority class $\mathbf{y}_t$, the probability $p(\mathbf{y}_s|\mathbf{y}_t)$ of selecting the source majority class $\mathbf{y}_s$ is proportional to $Bernoulli(\frac{n_s - n_t}{n_k})$. Once a source majority class $\mathbf{y}_s$ is determined, a sample $(\mathbf{x}_s, \mathbf{y}_s) \in \mathcal{D}_k$ is randomly selected. Then, the sample is fed into the feature extractor of the global model to obtain its original features $\mathbf{h}_s = g(\mathbf{x}_s; \mathbf{u})$. The original features $\mathbf{h}_s$ is going to be transformed into generated features $\widehat{\mathbf{h}}_t$ in the target class $\mathbf{y}_t$. The objective of feature transformation is formulated as:

$$\widehat{\mathbf{h}}_t = \underset{\widehat{\mathbf{h}}_t := \mathbf{h}_s + \delta}{\arg\min}\ \ell(\mathbf{y}_t, f(\widehat{\mathbf{h}}_t; \mathbf{v})), \quad (1)$$

where $\delta$ denotes the perturbation added to the original features $\mathbf{h}_s$, such that $\mathbf{h}_s + \delta$ can be predicted as the target minority class $\mathbf{y}_t$ by minimizing the loss $\ell(\mathbf{y}_t, f(\widehat{\mathbf{h}}_t; \mathbf{v}))$. An intuitive strategy to optimize Eq. (1) is to obtain the perturbations $\delta$ by the negative normalized gradients of classifier $\mathbf{v}$ [12]:

$$\nabla_{\widehat{\mathbf{h}}_t} = \frac{\partial \ell(\mathbf{y}_t, f(\widehat{\mathbf{h}}_t; \mathbf{v}))}{\partial \widehat{\mathbf{h}}_t}, \quad (2)$$

$$\delta = -\frac{\nabla_{\widehat{\mathbf{h}}_t}}{\|\nabla_{\widehat{\mathbf{h}}_t}\|_2}. \quad (3)$$

The gradient $\nabla_{\widehat{\mathbf{h}}_t}$ carries the information of how to predict $\widehat{\mathbf{h}}_t$ into class $\mathbf{y}_t$. Therefore, iteratively adding its negative

normalized value $\delta$ to $\widehat{\mathbf{h}}_t$ is towards the direction of loss decreasing of the next optimization iteration.

After obtaining the features $\widehat{\mathbf{h}}_t$, we cannot guarantee that $\widehat{\mathbf{h}}_t$ is certainly classified into the target minority class $\mathbf{y}_t$. Therefore, we check out its prediction confidence of $\mathbf{y}_t$ and consider if it can be used for personalized training. We accept all transformed features whose prediction confidence higher than a threshold $p_d$ called *drop probability* in FedAFA. We also empirically validate the influence of $p_d$ in Section 3.3. After selection by the drop probability, we put the selected generated features $(\widehat{\mathbf{h}}_t, \mathbf{y}_t)$, as well as the source majority class features $(\mathbf{h}_s, \mathbf{y}_s)$ into a set $\mathcal{G}_k$ for personalized training.

*2.2.2. Personalized Model Optimization*

One potential risk of using the above feature augmentation method is that it may damage the performance of the source majority classes. While the high-quality generated features of the target minority class $\mathbf{y}_t$ enhance the performance of the local minority classes, some original samples in the source majority class $\mathbf{y}_s$, which are used for augmentation, are likely to be misclassified to the target minority class $\mathbf{y}_t$. This is because the original and generated features $\mathbf{h}_s$ and $\widehat{\mathbf{h}}_t$ are still close in the feature space, although they can be classified into different classes. Therefore, to make the classification boundary separate the original and generated features, we proposed a new objective of the local personalized training $\mathcal{L}_{AFA}$ consisting of two parts: $\mathcal{L}_{gen}$ and $\mathcal{L}_{ori}$. They are the averaged training loss on the generated balanced feature set $\mathcal{G}_k^{bal}$ and the class-balanced local dataset $\mathcal{D}_k^{bal}$, respectively:

$$\mathcal{L}_{gen} = \frac{1}{|\mathcal{G}_k^{bal}|} \sum_{(\mathbf{h},\mathbf{y}) \in \mathcal{G}_k^{bal}} \ell(\mathbf{y}, f(\mathbf{h}; \mathbf{v}_k)), \quad (4)$$

$$\mathcal{L}_{ori} = \frac{1}{|\mathcal{D}_k^{bal}|} \sum_{(\mathbf{x},\mathbf{y}) \in \mathcal{D}_k^{bal}} \ell(\mathbf{y}, f(g(\mathbf{x}; \mathbf{u}); \mathbf{v}_k)), \quad (5)$$

$$\mathcal{L}_{AFA} = \lambda \mathcal{L}_{gen} + (1 - \lambda)\mathcal{L}_{ori}, \quad (6)$$

where $\lambda$ is a hyperparameter called *balance factor* to control the balance of adjusting the decision boundary between the original and generated features.

## 3. EXPERIMENTS

### 3.1. Experiment Setup

We evaluate FedAFA on two image classification datasets: CIFAR-10-LT and CIFAR-100-LT [18] [1]. Specifically, the number of samples in class $k$ decay exponentially by $\rho^k n_c$ [9], where $\rho^k \in (0, 1)$ controls the degree of long-tail and $n_c$ is the number of samples in each class of the original balanced dataset. In all experiments, we set the degree of data imbalance at 100, calculated by the number of samples in the largest class divided by the ones in the smallest class. We

---

[1]We refer to the long-tailed versions of these three benchmarks as CIFAR-10-LT and CIFAR-100-LT.
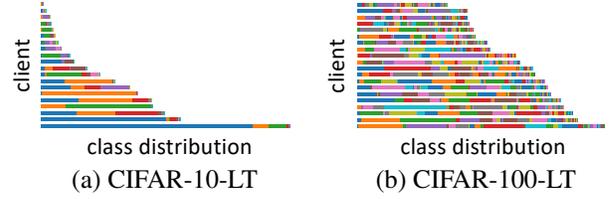


**Fig. 2**. Class distributions of each client on CIFAR-10/100-LT with $\alpha = 0.2$. Different colors represent different classes, and the length of the color block represents the number of samples in this class.

adopt the Dirichlet Distribution with the hyperparameter $\alpha$ to simulate different degrees of data heterogeneity [19]. A larger value of $\alpha$ means higher similarity between data distributions. We choose $\alpha = [0.5, 0.2]$ to make heterogeneous data across clients. Figure 2 shows the data distributions of each client for two long-tailed datasets used in our experiments with $\alpha = 0.2$. We choose ResNet32 as the learning model. The features of the second block of ResNet32 are selected to do feature augmentation of FedAFA. SGD is the optimizer with a learning rate of 0.005, a momentum of 0.9, and a weight decay of $5 \times 10^{-4}$ for local model optimization of all methods. We set the total number of clients at 20, the number of clients selected in each round at 10, the batch size at 64, the local epoch for personalized training at 1, and the global communication round at 500.

### 3.2. Experimental Results

In this subsection, we first compare FedAFA with two groups of state-of-the-art PFL methods. The methods in the first group are all PFL methods but do not specifically design for imbalanced data. The methods in the other group apply the long-tail learning methods to the PFL framework, including random oversampling (ROS), M2m [20], and cRT [16]. To verify that the features obtained by the global model are more robust, we also compare FedAFA and FedAFA_Loc, whose balanced feature sets are generated by the global model and the local personalization model, respectively.

As shown in Table 1, the performance of local training is poor because the model is only trained locally without any communication with other clients. Other PFL methods that only take data heterogeneity into account without considering the long-tail distribution also underperform compared with the methods that specifically address the long-tail distribution. Although these two groups of PFL methods improve FedAvg-FT from different perspectives, FedAFA is the only method consistently achieving promising personalized performance. In addition, by comparing the test accuracy of FedAFA_Loc and FedAFA, it is verified that the features extracted by the global model are more robust because more local noise is introduced when selecting the local personalized model to generate features.

To show that FedAFA produces better personalized mod-

| Dataset | CIFAR-10-LT | | CIFAR-100-LT | |
|---|---|---|---|---|
| $\alpha$ | 0.5 | 0.2 | 0.5 | 0.2 |
| Local training | 28.06 | 27.88 | 9.82 | 8.37 |
| FedAvg-FT [14] | 48.21 | 48.11 | 30.66 | 26.32 |
| FedProx [21] | 51.07 | 50.26 | 31.40 | 30.58 |
| LG-FedAvg [5] | 51.13 | 50.93 | 32.65 | 31.18 |
| Per-FedAvg [8] | 49.96 | 49.81 | 32.27 | 30.60 |
| pFedMe [6] | 50.05 | 49.59 | 32.82 | 30.05 |
| Ditto [7] | 50.84 | 50.26 | 31.98 | 31.11 |
| FedBN [22] | 50.98 | 50.17 | 31.62 | 30.37 |
| FedAvg+ROS | 59.29 | 56.51 | 34.04 | 32.02 |
| FedAvg+M2m | 62.85 | 56.30 | 32.63 | 32.69 |
| FedAvg+cRT | 63.33 | 59.55 | 34.42 | 33.11 |
| FedAFA_Loc | 62.24 | 62.15 | 34.12 | 33.06 |
| FedAFA | **64.52** | **63.70** | **36.57** | **35.44** |

**Table 1**. Test accuracy (%) of FedAFA compared with different methods. The best results are shown in bold.

els for each client, we evaluate the performance boost by FedAFA on each class against the baseline FedAvg-FT from the global perspective in Figure 3(a). It can be observed that FedAFA improves the performance for almost every class. In addition, we randomly select one client from all clients to illustrate the improvement of FedAFA on local minority classes. As shown in Figure 3(b), FedAFA can significantly improve the accuracy of local minority classes with only a few samples while keeping the same accuracy on the majority classes. These two experiments validate that FedAFA can improve the tail classes' generalization ability while guaranteeing the head classes' performance.

### 3.3. Effects of Hyperparameters

In FedAFA, there are three hyperparameters: the balance factor $\lambda$, the drop probability $p_d$, and the selected layer for feature augmentation. We conduct various experiments on CIFAR-10-LT to evaluate their influences.

We first compare the performance of FedAFA with $\lambda \in [0, 1]$ in Figure 4(a). FedAFA degenerates into the case of only training on $\mathcal{D}_k^{bal}$ when $\lambda = 0$, and FedAFA is equivalent to the case of only training on $\mathcal{G}_k$ when $\lambda = 1$. It can be observed that $\lambda \in [0.6, 0.8]$ generally performs better than either extreme case, which validates the effectiveness of the proposed loss function $L_{AFA}$.

The performance of FedAFA with $p_d \in (0, 1)$ is also shown in Figure 4(b). It can be observed that the performance increases as $p_d$ reaches 0.5 and then decreases when $p_d$ reaches 1. When $p_d$ is close to 1, only a few generated features can be selected for personalized training, which may only provide limited help. On the other hand, when $p_d$ is close to 0, generated features with any confidence are used for personalized training. They cannot focus on improving the decision boundary between the local majority and minority
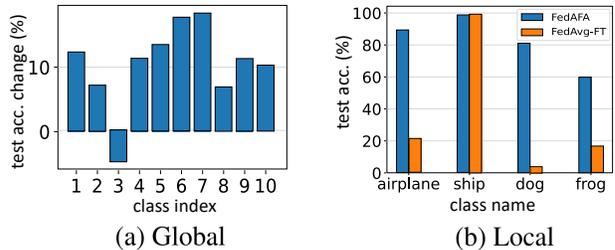


(a) Global  (b) Local

**Fig. 3**. Visualization of test accuracy change on each class with $\alpha = 0.2$. (a) The class index is sorted from head classes to tail classes from the global perspective. (b) The numbers of samples of each class from left to right are $\{8, 933, 9, 6\}$.
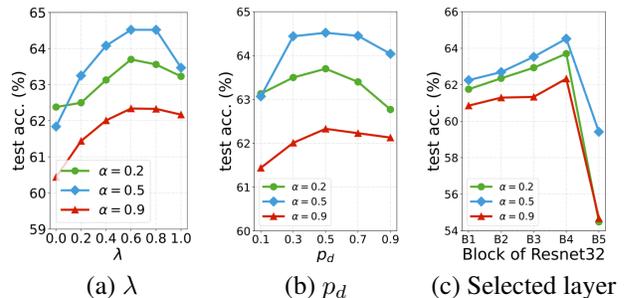


(a) $\lambda$  (b) $p_d$  (c) Selected layer

**Fig. 4**. Effects of hyperparameters of FedAFA.

classes. Therefore, this experiment validates that drop probability can effectively select generated features helpful for personalized training.

Finally, Figure 4(c) shows the performance of FedAFA by augmenting features on different blocks of ResNet32. It can be seen that FedAFA performs best when we select the features of Block 3 or Block 4. This observation confirms that selecting the feature space for augmentation is more effective than the sample space.

### 4. CONCLUSION

In this paper, FedAFA is proposed to solve the problem of PFL with heterogeneous and long-tailed data. FedAFA first transfers the knowledge of the local majority class to the local minority class to improve the performance of the local minority class. At the same time, FedAFA proposes a new optimization objective to maintain the performance of local majority classes. Therefore, FedAFA can enhance the generalization ability of the local minority classes while preserving the robust performance of the local majority classes. Experimental results show the superiority of FedAFA compared to other state-of-the-art PFL methods under different settings.

# 5. REFERENCES

[1] Dhruv Guliani, Françoise Beaufays, and Giovanni Motta, "Training speech recognition models with federated learning: A quality/cost framework," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021, pp. 3080–3084.

[2] Xiaodong Cui, Songtao Lu, and Brian Kingsbury, "Federated acoustic modeling for automatic speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021, pp. 6748–6752.

[3] Mehdi Salehi Heydar Abad, Emre Ozfatura, Deniz Gunduz, and Ozgur Ercetin, "Hierarchical federated learning across heterogeneous cellular networks," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 8866–8870.

[4] Hasin Us Sami and Başak Güler, "Over-the-air personalized federated learning," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2022, pp. 8777–8781.

[5] Filip Hanzely and Peter Richtárik, "Federated learning of a mixture of global and local models," *arXiv*, 2020.

[6] Canh T. Dinh, Nguyen Tran, and Josh Nguyen, "Personalized federated learning with moreau envelopes," in *Advances in Neural Information Processing Systems*, 2020, pp. 21394–21405.

[7] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith, "Ditto: Fair and robust federated learning through personalization," in *International Conference on Machine Learning*, 2021, pp. 6357–6368.

[8] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar, "Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach," in *Advances in Neural Information Processing Systems*, 2020, pp. 3557–3568.

[9] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie, "Class-balanced loss based on effective number of samples," in *IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9268–9277.

[10] Seulki Park, Jongin Lim, Younghan Jeon, and Jin Young Choi, "Influence-balanced loss for imbalanced visual classification," in *IEEE/CVF International Conference on Computer Vision*, 2021, pp. 735–744.

[11] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *International Conference on Neural Information Processing Systems*, 2019, pp. 1567–1578.

[12] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry, "Adversarial examples are not bugs, they are features," in *International Conference on Neural Information Processing Systems*, 2019, pp. 125–136.

[13] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, "Explaining and harnessing adversarial examples," *arXiv*, 2014.

[14] Kangkang Wang, Rajiv Mathews, Chloé Kiddon, Hubert Eichner, Françoise Beaufays, and Daniel Ramage, "Federated evaluation of on-device personalization," *arXiv*, 2019.

[15] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen, "Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9719–9728.

[16] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis, "Decoupling representation and classifier for long-tailed recognition," in *International Conference on Learning Representations*, 2020, pp. 1–16.

[17] Mi Luo, Fei Chen, Dapeng Hu, Yifan Zhang, Jian Liang, and Jiashi Feng, "No fear of heterogeneity: Classifier calibration for federated learning with non-iid data," *Advances in Neural Information Processing Systems*, vol. 34, pp. 5972–5984, 2021.

[18] Alex Krizhevsky, Geoffrey Hinton, et al., "Learning multiple layers of features from tiny images," *Technical Reports*, 2009.

[19] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown, "Measuring the effects of non-identical data distribution for federated visual classification," *arXiv*, 2019.

[20] Jaehyung Kim, Jongheon Jeong, and Jinwoo Shin, "M2m: Imbalanced classification via major-to-minor translation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13896–13905.

[21] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith, "Federated optimization in heterogeneous networks," *arXiv*, 2018.

[22] Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou, "Fedbn: Federated learning on non-iid features via local batch normalization," in *International Conference on Learning Representations*, 2021, pp. 1–27.