# FedBABU: Towards Enhanced Representation for Federated Image Classification

**Jaehoon Oh**
Graduate School of KSE
KAIST
jhoon.oh@kaist.ac.kr

**Sangmook Kim**
Graduate School of AI
KAIST
sangmook.kim@kaist.ac.kr

**Se-Young Yun**
Graduate School of AI
KAIST
yunseyoung@kaist.ac.kr

## Abstract

Federated learning has evolved to improve a single global model under data heterogeneity (as a curse) or to develop multiple personalized models using data heterogeneity (as a blessing). However, there has been little research considering both directions simultaneously. In this paper, we first investigate the relationship between them by analyzing Federated Averaging [31] at the client level and determine that a better federated global model performance does not constantly improve personalization. To elucidate the cause of this personalization performance degradation problem, we decompose the entire network into the body (i.e., extractor), related to universality, and the head (i.e., classifier), related to personalization. We then point out that this problem stems from training the head. Based on this observation, we propose a novel federated learning algorithm, coined as FedBABU, which updates only the body of the model during federated training (i.e., the head is randomly initialized and *never* updated), and the head is fine-tuned for personalization during the evaluation process. Extensive experiments show consistent performance improvements and an efficient personalization of FedBABU.

## 1 Introduction

Federated learning (FL) [31], a distributed learning framework with personalized data, has become an attractive field of research. From the early days of this field, improving *a single global model* across devices has been the main objective [48, 13, 28, 1], where the global model suffers from data heterogeneity. Many researchers have recently focused on *multiple personalized models* by leveraging data heterogeneity across devices as a blessing in disguise [4, 12, 47, 14, 37, 38]. Although many studies have been conducted on each research line, there remains a lack of research on how to train a good global model for personalization purposes [22, 23]. In this study, for personalized training, each local client model starts from a global model that learns information from all clients and leverages the global model to fit its own data distribution.

In [23], personalization methods that adapt the global model through fine-tuning on each device were analyzed. They observed that the effects of fine-tuning are encouraging and that training in a central location increases the initial accuracy (of a single global model) while decreasing the personalized accuracy (of on-device fine-tuned models). We focus on *why* opposite changes in the two performances appear. This suggests that the factors for universality and personalization must be dealt with separately, inspiring us to decouple the entire network into the body related to generality and the head related to specialty, as in [24, 45, 44, 10] for advanced analysis. Note that popular networks have one linear layer (e.g., MobileNet [21] and ResNet [20]), and the head is defined as this linear layer and the body is defined as all the layers except the head. The body of the model is related to representation learning, and the head of the model is related to linear decision boundary learning. We shed light on the cause of the personalization performance degradation problem by decoupling parameters, pointing out that such a problem stems from training the head.
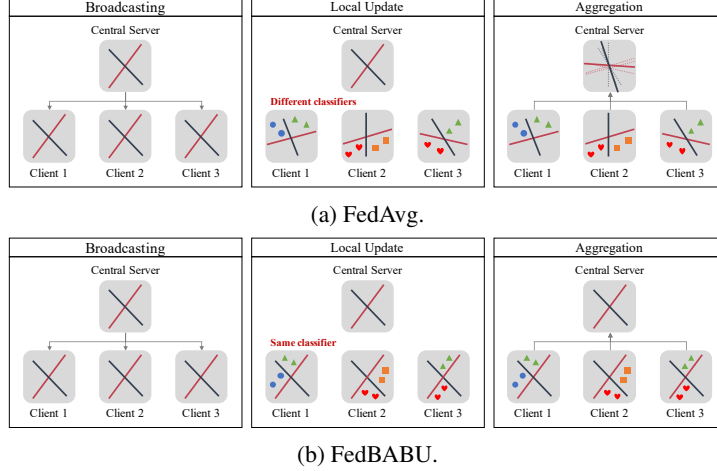
(a) FedAvg.



(b) FedBABU.

Figure 1: **Difference in the local update and aggregation stages between FedAvg and FedBABU**. In the figure, the lines represent the decision boundaries defined by the head (i.e., the last linear classifier) of the network. Different shapes indicate different classes. It is assumed that each client has two classes. (a) FedAvg updates the entire network during local updates on each client, and then the local networks are aggregated entirely. Therefore, the heads of all clients and the head of the server are different. Whereas, (b) FedBABU only updates the body (i.e., all the layers except the head) during local updates on each client, and then the local networks are aggregated body-partially. Therefore, the heads of all clients and the head of the server are the same.

Inspired by the above observations, we propose an algorithm to learn a single global model that can be efficiently personalized by simply changing the Federated Averaging (FedAvg) algorithm. FedAvg consists of four stages: client sampling, broadcasting, local update, and aggregation. In the client sampling stage, clients are sampled because the number of clients is so large that not all clients can participate per round. In the broadcasting stage, the server sends a global model (i.e., an initial random model at the first broadcast stage or an aggregated model afterward) to the participating clients. In the local update stage, the broadcast model of each device is trained based on its own data set. In the aggregation stage, locally updated models are sent to the server and are aggregated by averaging. Among the four stages, we focus on the *local update* stage from both the universality and personalization perspectives. *Here, we only update the body of the model in the local update stage, i.e., the head is never updated during federated training.* From this, we propose FedBABU, **Fed**erated Averaging with **B**ody **A**ggregation and **B**ody **U**pdate. Figure 1 describes the difference during the local update and aggregation stages between FedAvg and FedBABU. *Intuitively, our approach is a type of representation learning based on the same fixed criteria across all clients.* This simple change improves the representation power of a single global model and enables the trained single global model to be personalized more accurately and rapidly than FedAvg.

Our contributions are summarized as follows:

- We investigate the connection between a single global model and fine-tuned personalized models by analyzing the FedAvg algorithm at the client level and show that training the head using centralized data has a negative impact on personalization.
- We demonstrate that a fixed random classifier can have comparable performance as a learned classifier under a centralized setting. From this observation, we suggest that sharing a fixed random classifier across all clients can be more potent for matched aggregation than averaging each learned classifier in federated settings.
- We propose a novel algorithm, **FedBABU**, that reduces the update and aggregation parts from the entire model to the body of the model during federated training.
- We show that FedBABU is efficient, particularly under more significant data heterogeneity. Furthermore, a single global model trained with the FedBABU algorithm can be personalized rapidly (even with one fine-tuning epoch), and the personalization performance of FedBABU overwhelms that of other existing personalization FL algorithms.
- We adapt the body update and body aggregation idea to the regularization-based federated learning algorithm (such as FedProx [28]). We show that regularization reduces the personalization capabilities and that this problem is mitigated through BABU.

## 2 Related Works

**FL for a Single Global Model** Canonical federated learning, FedAvg proposed by [31], aims to learn a single global model that collects the benefits of affluent data without storing the raw data of the clients in a central server, reducing the communication costs through local updates. However, it is difficult to develop a globally optimal model for non-independent and identically distributed (non-IID) data derived from various clients. To solve this problem, studies have been conducted that make the data distribution of the clients IID-like or add regularization to the distance from the global model during local updates. [48] suggested that all clients share a subset of public data, and [13] augmented data for balancing the label distribution of clients. Recently, [28, 1] penalized local models that have a large divergence from the global model, adding a regularization term to the local optimization process and allowing the global model to converge more reliably. However, it should be noted that a single global model trained using the above methods is not optimized for each client.

**Personalized FL** Personalized federated learning aims to learn personalized local models stylized to each client. Although local models can be developed without federation, this method suffers from a limited amount of data. Therefore, to maintain the benefits of the federation and personalized models, many other methods have been applied to FL: clustering, multi-task learning, transfer learning, and meta-learning. [3, 30] clustered similar clients to match the data distribution within a cluster and learned separate models for each cluster without inter-cluster federation. Similar to clustering, multi-task learning aims to learn models for related clients simultaneously. Note that clustering ties related clients into a single cluster, whereas multi-task learning does not. The generalization of each model can be improved by sharing representations between the related clients. [38, 12] showed that multi-task learning is an appropriate learning scheme for personalized FL. Transfer learning is also a recommended learning scheme because it aims to deliver knowledge among clients. In addition, [43, 7] utilized transfer learning to enhance local models by transferring knowledge between related clients. Unlike the aforementioned methods that develop local models during training, [4, 14] attempted to develop a good initialized shared global model using bi-level optimization through a Model-Agnostic Meta-Learning (MAML) approach [15]. A well-initialized model can be personalized through updates on each client (such as inner updates in MAML). [23] argued that the FedAvg algorithm could be interpreted as a meta-learning algorithm, and a personalized local model with high accuracy can be obtained through fine-tuning from a global model learned using FedAvg. In addition, various technologies and algorithms for personalized FL are presented, and it will be helpful to read [40, 25] for more details.

**Decoupling the Body (Extractor) and the Head (Classifier) for Personalized FL** The training scheme through decoupling the entire network into the body and the head has been used in various fields, including long-tail recognition [24, 45], noisy label learning [46], and meta-learning [32, 34][1]. For personalized FL, there have been attempts to use this decoupling approach. For a consistent explanation, we describe each algorithm from the perspective of local update and aggregation parts. FedPer [2] learns the entire network jointly during local updates and aggregates the bottom layers only. When the bottom layers are matched with the body, the body is shared on all clients and the head is personalized to each client. LG-FedAvg [29] learns the entire network jointly during local updates and aggregates the top layers only based on the pre-trained global network via FedAvg. When the top layers are matched with the head, the body is personalized to each client and the head is shared on all clients. FedRep [9] learns the entire network sequentially during local updates and aggregates the body only. In the local update stage, each client first learns a classifier only with the aggregated representation, and then learns an extractor only with its own classifier with a single epoch. Unlike the above decoupling methods, we propose a FedBABU algorithm that learns only the body with a randomly initialized classifier during local updates and aggregates only the body. It is thought that fixing a classifier during the entire federated training provides the same guidelines on learning representations across all clients. Personalized local models are then obtained by fine-tuning the head.

## 3 Preliminaries

**FL training procedure** We summarize the training procedure of FL following the aforementioned four stages with formal notations. Let $\{1, \cdots, N\}$ be the set of all clients. Then, the participating

---

[1]Although there are more studies related to decoupling parameters [35, 26, 16, 8], we focus on decoupling the entire network into the body and head.

clients in the communication round $k$ with client fraction ratio $f$ is $C^k = \{C_i^k\}_{i=1}^{\lfloor Nf \rfloor}$. By broadcasting, the local parameters of the participating clients $\{\theta_i^k(0)\}_{i=1}^{\lfloor Nf \rfloor}$ are initialized by the global parameter $\theta_G^{k-1}$, that is, $\theta_i^k(0) \leftarrow \theta_G^{k-1}$ for all $i \in [1, \lfloor Nf \rfloor]$ and $k \in [1, \mathrm{K}]$. Note that $\theta_G^0$ is randomly initialized first. On its own device, each local model is updated using a locally kept data set. After local epochs $\tau$, the locally updated models become $\{\theta_i^k(\tau I_i^k)\}_{i=1}^{\lfloor Nf \rfloor}$, where $I_i^k$ is the number of iterations of one epoch on client $C_i^k$ (i.e., $\lceil \frac{n_{C_i^k}}{B} \rceil$), $n_{C_i^k}$ is the number of data samples for client $C_i^k$, and $B$ is the batch size. Therefore, $\tau I_i^k$ is the total number of updates. Note that our research deals with a balanced environment in which all clients have the same size data set (i.e., $I_i^k$ is a constant for all $k$ and $i$). Finally, the global parameter $\theta_G^k$ is aggregated by $\sum_{i=1}^{\lfloor Nf \rfloor} \frac{n_{C_i^k}}{n} \theta_i^k(\tau I_i^k)$, where $n = \sum_{i=1}^{\lfloor Nf \rfloor} n_{C_i^k}$, at the server. For our algorithm, the parameters $\theta$ are decoupled into the extractor parameters $\theta_{ext}$ and the classifier parameters $\theta_{cls}$.

**Experimental setup** We mainly use MobileNet on CIFAR100.[2] We set the number of clients to 100, and then each client has 500 training data and 100 test data. The classes in the training and test data sets are the same. For the heterogeneous distribution of client data, we refer to the experimental setting in [31]. We sort the data by label and divide the data into the same-sized shards. Because there is no overlapping data between shards, the size of a shard is defined by $\frac{|D|}{N \times s}$, where $|D|$ is the data set size, $N$ is the total number of clients, and $s$ is the number of shards per user. We control FL environments with three hyperparameters: client fraction ratio $f$, local epochs $\tau$, and shards per user $s$. $f$ is the ratio of participating clients in the total number of clients in every round, and a small $f$ is natural in the FL settings because the total number of clients is numerous. Local epochs $\tau$ are equal to the interval between two consecutive communication rounds. To fix the number of total updates for the consistency in all experiments, we fix the product of communication rounds and local epochs to 320 (e.g., if local epochs are 4, then the total number of communication rounds is 80). $\tau$ is closely related to the trade-off between accuracy and communication costs. A small $\tau$ provides an accurate federation but requires considerable communication costs. $s$ is related to the maximum number of classes each user can have; hence, as $s$ decreases, the degree of data heterogeneity increases.

**Evaluation** We calculate the *initial accuracy* and *personalized accuracy* of FedAvg and FedBABU following the federated personalization evaluation procedure proposed in [42] to analyze the algorithms at the client level: (1) the learned global model is broadcast to all clients, and is then evaluated on the test data set of each client $D_i^{ts}$ (referred to as the *initial accuracy*), (2) the learned global model is personalized using the training data set of each client $D_i^{tr}$ by fine-tuning with the fine-tuning epochs of $\tau_f$, and the personalized models are then evaluated on the test data set of each client $D_i^{ts}$ (referred to as the *personalized accuracy*). In addition, we calculate the *personalized accuracy* of other personalized FL algorithms (such as FedPer, LG-FedAvg, and FedRep). The algorithm 2 in Appendix A describes the evaluation procedure. The values (X$\pm$Y) in all tables indicate the mean$\pm$standard deviation of the accuracies across all clients. Here, reducing the variance over the clients could be an interesting topic, but is beyond the scope of this study.

## 4 Personalization of a Single Global Model

Table 1: Initial and personalized accuracy of FedAvg on CIFAR100 under various FL settings with 100 clients. MobileNet is used. The initial and personalized accuracy indicate the evaluated performance without fine-tuning and after five fine-tuning epochs on each client, respectively.

| FL settings | | $s$=100 (heterogeneity $\downarrow$) | | $s$=50 | | $s$=10 (heterogeneity $\uparrow$) | |
|---|---|---|---|---|---|---|---|
| $f$ | $\tau$ | Initial | Personalized | Initial | Personalized | Initial | Personalized |
| | 1 | 46.93$\pm$5.47 | 51.93$\pm$5.19 | 45.68$\pm$5.50 | 57.84$\pm$5.08 | 37.27$\pm$6.97 | 77.46$\pm$5.78 |
| 1.0 | 4 | 37.44$\pm$4.98 | 42.66$\pm$5.09 | 36.05$\pm$4.04 | 47.17$\pm$4.26 | 24.17$\pm$5.50 | 70.41$\pm$6.83 |
| | 10 | 29.58$\pm$4.87 | 34.62$\pm$4.97 | 29.57$\pm$4.29 | 40.59$\pm$5.23 | 17.85$\pm$7.38 | 63.51$\pm$7.38 |
| | 1 | 39.07$\pm$5.22 | 43.92$\pm$5.55 | 38.20$\pm$5.73 | 49.55$\pm$5.36 | 29.12$\pm$7.11 | 71.24$\pm$7.82 |
| 0.1 | 4 | 35.39$\pm$4.58 | 39.67$\pm$5.21 | 33.49$\pm$4.72 | 43.63$\pm$4.77 | 21.14$\pm$6.86 | 67.14$\pm$6.72 |
| | 10 | 28.18$\pm$4.83 | 33.13$\pm$5.22 | 27.34$\pm$4.96 | 38.09$\pm$5.17 | 14.40$\pm$5.64 | 62.67$\pm$6.52 |

---

[2]We also use 4convNet on CIFAR10 and ResNet on CIFAR100. The details of the architectures used are presented in Appendix A. The results of 4convNet and ResNet are presented in Appendix E and Appendix I, respectively.

Table 2: Initial and personalized accuracy of FedAvg on CIFAR100 under a realistic FL setting ($N$=100, $f$=0.1, $\tau$=10) according to $p$, which is the percentage of all client data that the server also has. Here, the entire network (F) or body (B) is updated on the server using the available data.

| $p$ | $s$=100 (heterogeneity $\downarrow$) | | $s$=50 | | $s$=10 (heterogeneity $\uparrow$) | |
|---|---|---|---|---|---|---|
| | Initial | Personalized | Initial | Personalized | Initial | Personalized |
| 0.00 | 28.18±4.83 | 33.13±5.22 | 27.34±4.96 | 38.09±5.17 | 14.40±5.64 | 62.67±6.52 |
| 0.05 (F) | 29.23±5.03 | 32.59±5.24 | 27.13±4.34 | 34.34±4.78 | 18.22±5.64 | 54.68±6.77 |
| 0.10 (F) | 30.59±4.93 | 33.34±5.30 | 29.62±4.27 | 35.50±4.84 | 19.24±5.15 | 49.62±7.48 |
| 0.05 (B) | 28.50±4.93 | 33.03±5.36 | 27.96±4.86 | 39.10±5.55 | 14.78±5.59 | 60.19±6.46 |
| 0.10 (B) | 32.90±4.77 | 36.82±4.66 | 32.81±4.97 | 40.80±5.62 | 18.35±6.75 | 60.94±7.30 |

We first investigate personalization of a single global model using the FedAvg algorithm, following [42, 23], to connect a single global model to multiple personalized models. We evaluate both the initial accuracy and personalized accuracy, assuming that the test data sets are not gathered in the server but scattered on the clients. Table 1 describes the accuracy of FedAvg on CIFAR100 according to different FL settings ($f$, $\tau$, and $s$) with 100 clients. The initial and personalized accuracy indicate the evaluated performance without fine-tuning and with five fine-tuning epochs for each client, respectively. As previous studies have shown, the more realistic the setting is (i.e., a smaller $f$, larger $\tau$, and smaller $s$), the lower the initial accuracy. Moreover, the tendency of the personalized models to converge higher than the global model observed in [23] is the same. More interestingly, it is shown that the gap between the initial accuracy and the personalized accuracy increases significantly as the data become more heterogeneous. It is thought that a small $s$ makes local tasks to be easy because the label distribution owned by each client is limited and the number of samples per class increases.

In addition, we conduct an intriguing experiment in which the initial accuracy increases but the personalized accuracy decreases, maintaining the FL training procedures. In [23], the authors showed that centralized trained models are more difficult to personalize. Similarly, we design an experiment where the federated trained models are more difficult to personalize. We assume that the server has a small portion of $p$ of the non-privacy data of the clients[3] such that the non-privacy data can be used in the server to mitigate the degradation derived from data heterogeneity. We update the global model using the non-privacy data after every aggregation with only one epoch. Table 2 describes the result of this experiment. We update the entire network (F in Table 2) on the server. As $p$ increases, the initial accuracy increases, as expected. However, the personalized accuracy decreases under significant data heterogeneity ($s$=10). This result implies that boosting a single global model can hurt personalization, which may be considered more important than the initial performance. Therefore, we agree that the development of an excellent global model should even consider the ability to be adequately fine-tuned or personalized.

To investigate the cause of personalization performance degradation, we hypothesize that unnecessary and confusing information for personalization is injected into a model, particularly a classifier, when a global model is trained on the server. To capture this, we only update the body (B in Table 2) on the server by zeroing the learning rate corresponding to a classifier. By narrowing the update parts, the personalization degradation problem is remedied significantly without affecting the initial accuracy. From this observation, we argue that training the head using harmonized data has a negative impact on personalization.[4]

## 5   FedBABU: Federated Averaging with Body Aggregation and Body Update

In this section, we propose a novel algorithm that learns a better single global model to be personalized efficiently. Inspired by prior studies on long-tailed recognition [45, 24], fine-tuning [10], and self-supervised learning [18], as well as our data sharing experiment, we decouple the entire network into the body (i.e., extractor) and the head (i.e., classifier). Extractors are trained for generalization, and classifiers are then trained for specialization. We apply this idea to federated learning by *never* training classifiers in the federated training phase (i.e., developing a single global model) and by fine-tuning classifiers for personalization (in the evaluation process).

---

[3]Non-privacy data are sampled randomly in our experiment, which can violate the FL environments. However, this experiment is conducted simply as motivation for our study.

[4]We also investigate personalization of centralized trained models such as [23], and the results are reported in Appendix K. In the centralized training case, training the head also has a negative impact on personalization.

## 5.1 Frozen Classifier in the Centralized Setting

Before proposing our algorithm, we empirically demonstrate that a model with an initialized and fixed classifier (body in Figure 2) has a performance comparable to a model that jointly learns an extractor and a classifier (full in Figure 2). Figure 2 depicts the test accuracy curve of MobileNet on CIFAR100 for various training scenarios in the centralized setting. The blue line represents the accuracy when all layers are trained, the red line represents the accuracy when only the body of the model is trained, and the green line represents the accuracy when only the head of the model is trained. It turns out that the fully trained model and the fixed classifier have almost the same



Figure 2: Test accuracy curves according to the update part in the centralized setting.

performance, whereas the fixed extractor performs poorly. Thus, we claim that *randomly initialized fixed classifiers are acceptable, whereas randomly initialized fixed extractors are unacceptable. Initialized classifiers are thought to serve as guidelines*. This characteristic is derived from the orthogonal initialization [36] on the head, which is explained in Appendix B.
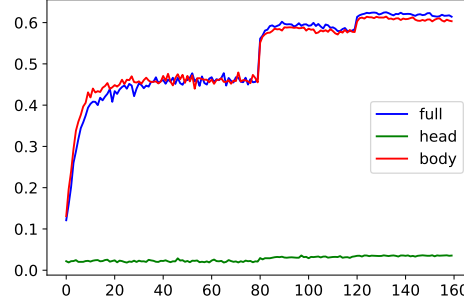
## 5.2 FedBABU Algorithm

Based on the insights and results in Section 5.1, we propose a new FL algorithm, called FedBABU (**Fed**erated Averaging with **B**ody **A**ggregation and **B**ody **U**pdate). Extractors are trained only, whereas classifiers are never trained during federated training. Therefore, there is no need to aggregate the head. Formally, the model parameters $\theta$ are decoupled into the extractor parameters $\theta_{ext}$ and the classifier parameters $\theta_{cls}$. Note that $\theta_{cls}$ on any client is fixed with the classifier parameters of a randomly initialized global parameters $\theta_G^0$ until a single global model converges. This is implemented by zeroing the learning rate corresponding to the classifier. *Intuitively, it is thought that the same fixed classifier on all clients serves as the same criteria on learning representations across all clients despite the passage of training time.* The FedBABU algorithm is described in Appendix A. In this section, we demonstrate the ability related to personalization of FedBABU (from Section 5.2.1 to Section 5.2.3). We further show our algorithm's applicability (Section 5.2.4).

### 5.2.1 Personalization of FedBABU

To investigate the dominant factor for personalization, we compare the performance according to the fine-tuned part. Table 3 describes the results of this experiment. Global models are fine-tuned with five epochs based on the training data set of each client. It is shown that fine-tuning including the head (i.e., Head or Full in Table 3) is better for personalization than body-partially fine-tuning (i.e., Body in Table 3). For consistency of the evaluation, for FedBABU, we fine-tune the entire model in the remainder of this paper. Note that the computational costs can be reduced by fine-tuning only the head in the case of FedBABU

Table 3: Personalized accuracy of MobileNet on CIFAR100 according to the fine-tuned part. The fine-tuning epochs is 5, and $f$ is 0.1.

| Hyperparameter | | Update part for personalization | | |
|---|---|---|---|---|
| $s$ | $\tau$ | Body | Head | Full |
| 100 | 1 | 44.26±5.12 | 49.76±5.03 | 49.67±4.92 |
| | 4 | 39.61±4.68 | 44.74±5.02 | 44.74±5.10 |
| | 10 | 32.45±5.42 | 36.48±5.04 | 35.94±5.06 |
| 50 | 1 | 48.54±5.23 | 56.76±5.68 | 56.69±5.16 |
| | 4 | 41.27±5.04 | 49.45±5.41 | 49.55±5.58 |
| | 10 | 35.42±5.60 | 42.55±5.70 | 42.63±5.59 |
| 10 | 1 | 72.81±7.32 | 75.97±6.29 | 76.02±6.29 |
| | 4 | 69.12±6.70 | 70.74±6.47 | 71.00±6.63 |
| | 10 | 64.77±7.14 | 66.28±6.77 | 66.32±7.02 |

without performance degradation; however, in the case of FedAvg, a performance gap appears (Appendix D). The personalization of FedBABU looks similar to the linear evaluation protocol, which adds a linear layer to the well-trained extractor and makes it suitable for multiple tasks.

### 5.2.2 Personalization Performance Comparison

We compare `FedBABU` with existing methods from the perspective of personalization. Details of evaluation procedure and implementations of each algorithm are presented in Appendix A. Table 4 describes the personalized accuracy of various algorithms. Interestingly, `FedAvg` overwhelms other recent personalized FL algorithms on CIFAR100 in most cases.[5] These results are similar to recent trends in the field of meta-learning, where fine-tuning based on well-trained representations overwhelms advanced few-shot classification algorithms for heterogeneous tasks [5, 6, 11, 41]. `FedBABU` (ours) further overwhelms `FedAvg`. It is believed that enhanced representation by freezing the head during federated training improve performance, explained in Appendix O.

---

[5]The reason why FedAvg+Fine-tuning is effective itself needs to be discussed more.

Table 4: Personalized accuracy comparison on CIFAR100 under various settings with 100 clients. MobileNet is used. **Bold** indicates the best accuracy.

| Hyperparameter | | | Personalized accuracy | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $s$ | $f$ | $\tau$ | FedBABU (Ours) | FedAvg [31] | FedPer [2] | LG-FedAvg [29] | FedRep [9] | Per-FedAvg [14] | Local-only |
| 100 | 1.0 | 1 | **55.79±4.57** | 51.93±5.19 | 51.95±5.30 | 53.01±5.26 | 18.29±3.59 | 47.09±7.35 | 20.6±3.15 |
| | | 4 | **44.49±4.91** | 42.66±5.09 | 40.87±5.05 | 43.09±4.74 | 15.32±3.79 | 39.07±7.59 | |
| | | 10 | 33.20±4.54 | **34.62±4.97** | 32.91±4.97 | 34.64±5.03 | 13.45±3.26 | 30.22±6.59 | |
| | 0.1 | 1 | **49.67±4.92** | 43.92±5.55 | 45.17±4.70 | 40.91±5.50 | 23.84±3.92 | 48.10±7.42 | |
| | | 4 | **44.74±5.10** | 39.67±5.21 | 39.30±4.92 | 37.87±4.99 | 16.01±3.48 | 33.70±7.04 | |
| | | 10 | **35.94±5.06** | 33.13±5.22 | 32.08±4.97 | 30.08±5.34 | 11.11±3.13 | 25.82±5.83 | |
| 50 | 1.0 | 1 | **61.09±4.91** | 57.84±5.08 | 57.16±5.26 | 58.44±5.53 | 24.75±5.02 | 43.75±7.94 | 28.02±4.01 |
| | | 4 | **51.56±5.04** | 47.17±4.26 | 48.89±5.40 | 47.78±4.72 | 21.55±4.36 | 37.59±7.87 | |
| | | 10 | **42.09±5.12** | 40.59±5.23 | 39.90±5.54 | 40.32±4.70 | 19.92±4.50 | 28.75±6.40 | |
| | 0.1 | 1 | **56.69±5.16** | 49.55±5.36 | 51.63±5.27 | 42.64±5.55 | 32.88±5.09 | 43.96±7.40 | |
| | | 4 | **49.55±5.58** | 43.63±4.77 | 46.31±5.63 | 38.54±4.71 | 21.13±3.96 | 28.67±6.98 | |
| | | 10 | **42.63±5.59** | 38.09±5.17 | 39.81±4.88 | 30.79±6.12 | 15.15±4.01 | 21.64±6.16 | |
| 10 | 1.0 | 1 | **79.17±6.51** | 77.46±5.78 | 74.71±6.35 | 77.49±5.60 | 61.28±8.27 | 36.59±8.98 | 61.52±7.22 |
| | | 4 | **74.60±6.69** | 70.41±6.83 | 65.61±7.13 | 69.97±6.42 | 50.59±7.94 | 18.31±10.57 | |
| | | 10 | **66.64±6.84** | 63.51±7.38 | 59.71±7.35 | 61.50±7.28 | 42.13±7.53 | 11.54±8.87 | |
| | 0.1 | 1 | **76.02±6.29** | 71.24±7.82 | 69.36±6.77 | 51.75±9.32 | 60.13±7.72 | 31.21±11.66 | |
| | | 4 | **71.00±6.63** | 67.14±6.72 | 62.62±7.63 | 35.80±10.55 | 45.91±7.68 | 14.34±9.51 | |
| | | 10 | **66.32±7.02** | 62.67±6.52 | 59.50±7.33 | 25.04±12.02 | 34.30±7.84 | 9.17±6.95 | |

### 5.2.3 Personalization Speed of FedAvg and FedBABU

Table 5: Performance according to the fine-tune epochs (FL setting: $f$=0.1, and $\tau$=10).

| $s$ | Algorithm | Fine-tune epochs ($\tau_f$) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 0 (Initial) | 1 | 2 | 3 | 4 | 5 | 10 | 15 | 20 |
| 50 | FedAvg | 27.34±4.96 | 29.17±5.01 | 32.39±4.77 | 34.97±5.13 | 36.78±5.13 | 38.09±5.17 | 40.56±5.43 | 41.20±5.51 | 40.86±5.13 |
| | FedBABU | 27.91±5.27 | 35.20±5.58 | 40.60±5.47 | 42.12±5.61 | 42.74±5.60 | 42.63±5.59 | 41.94±5.68 | 41.19±5.52 | 40.61±5.28 |
| 10 | FedAvg | 14.40±5.64 | 27.43±6.46 | 48.63±7.30 | 58.08±6.11 | 61.27±6.15 | 62.67±6.52 | 63.91±6.49 | 64.56±6.45 | 64.89±6.53 |
| | FedBABU | 18.50±7.82 | 63.29±7.55 | 66.05±6.93 | 66.10±6.54 | 66.40±7.24 | 66.32±7.02 | 66.07±7.57 | 66.24±7.67 | 66.32±7.71 |

We investigate the personalization speed of FedAvg and FedBABU by controlling the fine-tuning epochs $\tau_f$ in the evaluation process. Table 5 describes the initial (when $\tau_f$ is 0) and personalized (otherwise) accuracy of FedAvg and FedBABU. Here, 1 epoch is equal to 10 updates in our case because each client has 500 training samples and the batch size is 50. It is shown that FedAvg requires sufficient epochs for fine-tuning, as reported in [23]. Notably, FedBABU achieves better accuracy with a small number of epochs. It means that FedBABU can personalize accurately and rapidly, especially when fine-tuning is either costly or restricted.

### 5.2.4 Body Aggregation and Body Update on the FedProx

Table 6: Initial and personalized accuracy of FedProx and FedProx+BABU with $\mu$ of 0.01 on CIFAR100 with 100 clients and $f$ of 0.1.

| FL settings | | $s$=100 (heterogeneity ↓) | | $s$=50 | | $s$=10 (heterogeneity ↑) | |
|---|---|---|---|---|---|---|---|
| Algorithm | $\tau$ | Initial | Personalized | Initial | Personalized | Initial | Personalized |
| FedProx | 1 | 46.52±4.56 | 50.95±4.65 | 42.20±4.90 | 51.29±5.20 | 28.16±9.00 | 66.39±7.79 |
| | 4 | 36.54±4.74 | 39.83±4.71 | 33.59±4.80 | 40.17±5.11 | 18.20±7.62 | 41.56±9.34 |
| | 10 | 28.63±4.40 | 31.90±4.16 | 26.88±4.59 | 32.92±5.00 | 13.62±7.73 | 43.48±9.32 |
| FedProx +BABU | 1 | 48.53±5.15 | 57.44±4.72 | 46.25±5.31 | 63.12±5.25 | 33.13±8.11 | 78.86±5.70 |
| | 4 | 37.17±4.41 | 45.26±4.76 | 33.86±5.44 | 50.18±5.14 | 22.94±9.90 | 75.71±5.33 |
| | 10 | 27.79±3.95 | 35.68±4.34 | 27.48±5.22 | 42.37±6.10 | 15.66±8.29 | 67.15±7.10 |

FedProx [28] regularizes the distance between a global model and local models to prevent local models from deviating. The degree of regularization is controlled by $\mu$. Note that when $\mu$ is 0.0, FedProx is reduced to FedAvg. Table 6 describes the initial and personalized accuracy of FedProx and FedProx+BABU with $\mu$ of 0.01 when $f$=0.1, and the results when $f$=1.0 are reported in Appendix L. The global models trained by FedProx reduce the personalization capabilities compared to FedAvg (refer to Table 1 when $f$=0.1), particularly under realistic FL settings. We adapt the body aggregation and body update idea to the FedProx algorithm, referred to as FedProx+BABU, which performs better than personalization of FedAvg. Furthermore, FedProx+BABU also has an improved personalized performance compared to FedBABU (refer to Table 4 when $f$=0.1). This means that the regularization of the body is still meaningful. Our algorithm and various experiments suggest future directions of federated learning: *Which parts should be federated and enhanced? Representation!*

## 6 Conclusion

In this study, we focused on how to train a good federated global model for personalization purposes. Based on parameter decoupling, we proposed the FedBABU algorithm, which learns a shared global model that can rapidly adapt to heterogeneous data on each client. This global model can be efficiently personalized by fine-tuning each client's model using its own data set. Extensive experimental results showed that FedBABU outperforms various personalized FL algorithms. Our improvement emphasizes the importance of federating and enhancing the representation for FL.

## Acknowledgments and Disclosure of Funding

## References

[1] Durmus Alp Emre Acar, Yue Zhao, Ramon Matas Navarro, Matthew Mattina, Paul N Whatmough, and Venkatesh Saligrama. Federated learning based on dynamic regularization. In *International Conference on Learning Representations*, 2021.

[2] Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. *arXiv preprint arXiv:1912.00818*, 2019.

[3] Christopher Briggs, Zhong Fan, and Peter Andras. Federated learning with hierarchical clustering of local updates to improve training on non-iid data. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9. IEEE, 2020.

[4] Fei Chen, Mi Luo, Zhenhua Dong, Zhenguo Li, and Xiuqiang He. Federated meta-learning with fast convergence and efficient communication. *arXiv preprint arXiv:1802.07876*, 2018.

[5] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *International Conference on Learning Representations*, 2019.

[6] Yinbo Chen, Xiaolong Wang, Zhuang Liu, Huijuan Xu, and Trevor Darrell. A new meta-baseline for few-shot learning. *arXiv preprint arXiv:2003.04390*, 2020.

[7] Yiqiang Chen, Xin Qin, Jindong Wang, Chaohui Yu, and Wen Gao. Fedhealth: A federated transfer learning framework for wearable healthcare. *IEEE Intelligent Systems*, 35(4):83–93, 2020.

[8] Yutian Chen, Abram L Friesen, Feryal Behbahani, Arnaud Doucet, David Budden, Matthew W Hoffman, and Nando de Freitas. Modular meta-learning with shrinkage. *arXiv preprint arXiv:1909.05557*, 2019.

[9] Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. Exploiting shared representations for personalized federated learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 2089–2099. PMLR, 2021.

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.

[11] Guneet S Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. *arXiv preprint arXiv:1909.02729*, 2019.

[12] Canh T Dinh, Tung T Vu, Nguyen H Tran, Minh N Dao, and Hongyu Zhang. Fedu: A unified framework for federated multi-task learning with laplacian regularization. *arXiv preprint arXiv:2102.07148*, 2021.

[13] Moming Duan, Duo Liu, Xianzhang Chen, Yujuan Tan, Jinting Ren, Lei Qiao, and Liang Liang. Astraea: Self-balancing federated learning for improving classification accuracy of mobile deep learning applications. In *2019 IEEE 37th International Conference on Computer Design (ICCD)*, pages 246–254. IEEE, 2019.

[14] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*, 2020.

[15] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR, 2017.

[16] Sebastian Flennerhag, Andrei A Rusu, Razvan Pascanu, Francesco Visin, Hujun Yin, and Raia Hadsell. Meta-learning with warped gradient descent. *arXiv preprint arXiv:1909.00025*, 2019.

[17] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.

[18] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[21] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[22] Shaoxiong Ji, Teemu Saravirta, Shirui Pan, Guodong Long, and Anwar Walid. Emerging trends in federated learning: From model fusion to federated x learning, 2021.

[23] Yihan Jiang, Jakub Konečnỳ, Keith Rush, and Sreeram Kannan. Improving federated learning personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019.

[24] 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*, 2019.

[25] Viraj Kulkarni, Milind Kulkarni, and Aniruddha Pant. Survey of personalization techniques for federated learning. In *2020 Fourth World Conference on Smart Trends in Systems, Security and Sustainability (WorldS4)*, pages 794–797. IEEE, 2020.

[26] Yoonho Lee and Seungjin Choi. Gradient-based meta-learning with learned layerwise metric and subspace. In *International Conference on Machine Learning*, pages 2927–2936. PMLR, 2018.

[27] José Lezama, Qiang Qiu, Pablo Musé, and Guillermo Sapiro. Ole: Orthogonal low-rank embedding-a plug and play geometric loss for deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8109–8118, 2018.

[28] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.

[29] Paul Pu Liang, Terrance Liu, Liu Ziyin, Nicholas B Allen, Randy P Auerbach, David Brent, Ruslan Salakhutdinov, and Louis-Philippe Morency. Think locally, act globally: Federated learning with local and global representations. *arXiv preprint arXiv:2001.01523*, 2020.

[30] Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*, 2020.

[31] H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. Communication-efficient learning of deep networks from decentralized data. In *Proc. 20th Int'l Conf. Artificial Intelligence and Statistics (AISTATS)*, pages 1273–1282, 2017.

[32] Jaehoon Oh, Hyungjun Yoo, ChangHwan Kim, and Se-Young Yun. Boil: Towards representation change for few-shot learning. In *International Conference on Learning Representations*, 2021.

[33] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32:8026–8037, 2019.

[34] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. In *International Conference on Learning Representations*, 2019.

[35] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. *arXiv preprint arXiv:1807.05960*, 2018.

[36] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.

[37] Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. Personalized federated learning using hypernetworks. *arXiv preprint arXiv:2103.04628*, 2021.

[38] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet Talwalkar. Federated multi-task learning. *arXiv preprint arXiv:1705.10467*, 2017.

[39] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4080–4090, 2017.

[40] Alysa Ziying Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *arXiv preprint arXiv:2103.00710*, 2021.

[41] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? *arXiv preprint arXiv:2003.11539*, 2020.

[42] Kangkang Wang, Rajiv Mathews, Chloé Kiddon, Hubert Eichner, Françoise Beaufays, and Daniel Ramage. Federated evaluation of on-device personalization. *arXiv preprint arXiv:1910.10252*, 2019.

[43] Hongwei Yang, Hui He, Weizhe Zhang, and Xiaochun Cao. Fedsteg: A federated transfer learning framework for secure image steganalysis. *IEEE Transactions on Network Science and Engineering*, 2020.

[44] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *arXiv preprint arXiv:1411.1792*, 2014.

[45] Haiyang Yu, Ningyu Zhang, Shumin Deng, Zonggang Yuan, Yantao Jia, and Huajun Chen. The devil is the classifier: Investigating long tail relation classification with decoupling analysis. *arXiv preprint arXiv:2009.07022*, 2020.

[46] Hui Zhang and Quanming Yao. Decoupling representation and classifier for noisy label learning. *arXiv preprint arXiv:2011.08145*, 2020.

[47] Michael Zhang, Karan Sapra, Sanja Fidler, Serena Yeung, and Jose M Alvarez. Personalized federated learning with first order model optimization. *arXiv preprint arXiv:2012.08565*, 2020.

[48] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.