

# Local Learning Matters: Rethinking Data Heterogeneity in Federated Learning

Matias Mendieta<sup>1</sup>, Taojiannan Yang<sup>1</sup>, Pu Wang<sup>2</sup>, Minwoo Lee<sup>2</sup>, Zhengming Ding<sup>3</sup>, Chen Chen<sup>1</sup>

<sup>1</sup>Center for Research in Computer Vision, University of Central Florida, USA

<sup>2</sup>Department of Computer Science, University of North Carolina at Charlotte, USA

<sup>3</sup>Department of Computer Science, Tulane University, USA

{mendieta,taoyang1122}@knights.ucf.edu; chen.chen@crcv.ucf.edu

{pu.wang,minwoo.lee}@uncc.edu; zding1@tulane.edu

## Abstract

*Federated learning (FL) is a promising strategy for performing privacy-preserving, distributed learning with a network of clients (i.e., edge devices). However, the data distribution among clients is often non-IID in nature, making efficient optimization difficult. To alleviate this issue, many FL algorithms focus on mitigating the effects of data heterogeneity across clients by introducing a variety of proximal terms, some incurring considerable compute and/or memory overheads, to restrain local updates with respect to the global model. Instead, we consider rethinking solutions to data heterogeneity in FL with a focus on local learning generality rather than proximal restriction. To this end, we first present a systematic study informed by second-order indicators to better understand algorithm effectiveness in FL. Interestingly, we find that standard regularization methods are surprisingly strong performers in mitigating data heterogeneity effects. Based on our findings, we further propose a simple and effective method, FedAlign, to overcome data heterogeneity and the pitfalls of previous methods. FedAlign achieves competitive accuracy with state-of-the-art FL methods across a variety of settings while minimizing computation and memory overhead. Code is available at <https://github.com/mmendiet/FedAlign>.*

## 1. Introduction

Federated learning (FL) [17] enables a large number of clients to perform collaborative training of machine learning models without compromising data privacy. In the FL setting, participating clients are typically deployed in a variety of environments or owned by a diverse set of users. Therefore, the distribution of each client’s local data can vary considerably (i.e., data heterogeneity). This non-IID data distribution among participating devices in FL makes optimization particularly challenging. As each client trains

locally on their own data, they step towards their respective local minimum. However, this local convergence point may not be well aligned with the objective of the global model (that is, the model being learned through aggregation at the central server). Therefore, the client model often drifts away from the ideal global optimization point and overfits to its local objective. When such client drifting occurs, the performance of the central aggregated model is hindered [9, 14].

One straight-forward solution to this phenomenon is to simply limit the number of local training epochs performed between central aggregation steps. However, this severely hinders the convergence speed of the FL system, and many communication rounds are required to achieve adequate performance. The time to convergence and immense communication overhead incurred by such an approach are often not tolerable for real-world distributed systems. Therefore, effectively addressing data heterogeneity is of paramount concern in federated learning.

Many algorithmic solutions to this problem have been proposed in the literature [1, 10, 15, 23]. These strategies typically focus on mitigating the effects of data heterogeneity across clients by introducing a variety of *proximal terms* to restrain local updates with respect to the global model. However, by restraining the drift, they also inherently limit the local convergence potential; less novel information is gathered per communication round. Consequently, many current FL algorithms do not provide stable performance improvements across different non-IID settings in comparison to classic baselines [14, 15], especially on vision tasks beyond the difficulty of MNIST [13]. Furthermore, existing methods have paid little attention to the resource constraints of the client, typically scarce for deployed FL edge devices, and in some cases incur considerable compute and/or memory overheads on the client in their effort to alleviate client drift. For example, the state-of-the-art (SOTA) method MOON performs well on federated image tasks, but to do so incurs a  $\sim 3x$  overhead in both memory and com-

pute compared to the standard FedAvg baseline [17].

**Motivation.** In the centralized training paradigm, network generalization capability has been well studied to combat overfitting. Even in standard settings where the training and test data are drawn from a similar distribution, models still overfit on the training data if no precautions are taken. This effect is further intensified when the training and test data are of different distributions. Various regularization techniques are introduced to enforce learning generality during training and preserve suitable test performance. Similarly, overfitting to the local training data of each device in FL is detrimental to overall network performance, as the client drifting effect creates conflicting objectives among local models. *Thus, a focus on improving model generality should be of primary concern in the presence of data heterogeneity.* Improving local learning generality during training would inherently position the objective of the clients closer to the overall global objective. However, despite its intuitive motivations, this perspective has been overlooked by the bulk of current FL literature.

Therefore, in this paper, we propose rethinking approaches to data heterogeneity in terms of **local learning generality** rather than proximal restriction. Specifically, we carefully analyze the effectiveness of various data and structural regularization methods at reducing client drift and improving FL performance (Section 3). Utilizing second-order information and insights from out-of-distribution generality literature [19, 21], we identify theoretical indicators for successful FL optimization, and evaluate across a variety of FL settings for empirical validation.

Although some of the regularization methods perform well at mitigating client drift, *significant resource overheads* are still incurred to achieve the best performance (see Section 4). Therefore, we propose **FedAlign**, a distillation-based regularization method that promotes local learning generality while maintaining excellent resource efficiency. Specifically, FedAlign focuses on regularizing the Lipschitz constants of the final block in a network with respect to its representations. By focusing solely on the last block, we effectively regularize the portion of the network most prone to overfitting and keep additional resource needs to a minimum. Therefore, FedAlign achieves state-of-the-art accuracy on multiple datasets across a variety of FL settings, while requiring significantly less computation and memory overhead in comparison to other state-of-the-art methods.

Our contributions are as follows:

- We approach one of the most troublesome FL challenges (*i.e.* client drift caused by data heterogeneity) from a unique angle than any other previous work. We do not focus on reparameterization tricks to maintain closeness to the central model, or adjust the aggregation scheme to mitigate the effects of non-IID data distributions. *Rather, we propose the rethinking of this problem from fundamen-*

*tal machine learning training principles.* In this way, we analyze the performance of standard regularization methods on FL and their effectiveness against data heterogeneity.

- Not only do we empirically analyze the performance of regularization methods in FL, we also propose to take a deeper look. Specifically, we inform our analysis with theoretical indicators of learning generality to provide insight into which methods are best and why. We find that Hessian eigenvalue/trace measurements and Hessian matching across clients to be meaningful indicators for optimal FL methods. Additionally, we perform a thorough ablation study across a variety of FL settings to understand the empirical effects of different methods. Our aim is to provide this valuable knowledge to the FL community to inspire new, productive research directions.
- Informed by our analysis and examining the pitfalls of previous methods, we propose FedAlign, which achieves competitive state-of-the-art accuracy while maintaining memory and computational efficiency.

## 2. Related Work

**Federated Learning.** In general, federated learning algorithms aim to obtain a collective model which minimizes the training loss across all clients. This objective can be expressed as

$$\min_w F(w) = \sum_{c=1}^C \alpha_c F_c(w), \quad (1)$$

where  $F_c(w)$  is the local loss of device  $c$ , and  $\alpha_c$  is an arbitrary weight parameter with  $\sum_{c=1}^C \alpha_c = 1$ . One of the earliest algorithms proposed in FL is Federated Averaging, or FedAvg [17]. This approach simply optimizes the local training loss with standard SGD training, and aggregates using a weighted average approach with  $a_c = \frac{n_c}{n}$ , where  $n_c$  is equal to the number of training samples on client  $c$ , with a total of  $n$  training samples partitioned across all  $C$  clients.

Recent works attempt to improve over this baseline with two distinct focuses: improvements to the local training at the client, or improvements to the global aggregation process at the server. In this work, we focus on local training and client drift, and therefore we will first discuss methods of this nature. To mitigate data heterogeneity complications, a common approach is to introduce proximal terms to the local training loss. For instance, FedProx [23] adds the proximal term  $\frac{\mu}{2} \|w - w^t\|^2$ , where  $\mu$  is a hyperparameter,  $w$  is the current local model weights, and  $w^t$  is the global model weights from round  $t$ . The goal of this reparameterization is to minimize client drift by limiting the impact of local updates from becoming extreme. More recently, MOON [15] proposes a similar reparameterization idea inspired by contrastive learning. Specifically, the authors form a local model constrastive loss comparing rep-

representations of three models: the global model, the current local model, and a copy of the local model from the previous round. The goals of this term are similar to that of Fed-Prox but in feature representation space; to push the current local representation closer to the global representation. At the same time, the current local model is being pushed away from the representations of the local model copy of the previous round. Other methods [1, 10] follow similar ideas; they aim to limit the impact of the local update or shift the update with a correction term.

However, these approaches have two main downsides. First, by restraining the drift, they also inherently limit the local convergence potential. With this, not as much new information is gathered per communication round. Second, many of these methods incur substantial overheads in memory and/or computation. For instance, because of its model contrastive loss, MOON [15] requires the storage of three full-size models in memory simultaneously during training, and forward passing through each of these every iteration. This requires a great deal of additional resources, which are often already scarce in FL client settings.

Other works focus on the server side of the system, aiming to improve the aggregation algorithm. [34] propose a Bayesian nonparametric method for matching neurons across local models at aggregation rather than naively averaging. However, the presented framework is limited in application to fully-connected networks, and therefore [27] extend it to CNNs and LSTMs. FedNova [28] presents a normalized averaging method as an alternative to the simple FedAvg update. As we focus on the local training, these works are orthogonal to our work. A few approaches [18, 25, 32] propose federated schemes inspired by the data augmentation method Mixup, using similar averaging techniques on the local data and sharing the augmented data with the global model or other devices. However, even though the data is augmented in some way prior to distribution, the sharing of private data from the client is less than ideal for privacy preservation. Furthermore, sharing additional data worsens the communication burden on the system, which is a principal concern in FL.

**Learning Generality.** In traditional centralized training, the practice of regularization of various forms is common practice for improving generality. Data-level regularization, including basic data augmentations and other more advanced techniques [33, 36], are known to be quite effective. Other methods introduce a level of noise to the training process via structural modification; for instance, random or deliberate modifications to the network connectivity [3, 6, 26]. [29] proposes a hybrid approach that introduces self-guided gradient perturbation to the training process through the use of sub-network representations, knowledge distillation, and input transformations. As part of this work, we employ a variety of regularization methods in many FL settings and

analyze their performance in comparison to state-of-the-art FL algorithms.

### 3. Empirical Study

**We wish to assess the data heterogeneity challenge of FL from a simple yet unique perspective of local learning generality.** Specifically, we first study the effectiveness of standard regularization techniques as solutions to this FL challenge in comparison to state-of-the-art methods.

#### 3.1. Preliminaries

We employ three FL algorithms, namely FedAvg, Fed-Prox, and MOON. These works represent both classic baselines and current state-of-the-art, and are described in Section 2. For comparison, we employ three state-of-the-art regularization methods: Mixup [36], Stochastic Depth [6], and GradAug [29]. Specifically, these regularization methods are applied to the local optimization within a standard FedAvg setup, and their operations are described as follows.

Mixup is a data-level augmentation technique that performs linear interpolation between two samples. Specifically, given two sample-label pairs  $(x_i, y_i)$  and  $(x_j, y_j)$ , they are combined as  $\tilde{x} = \beta x_i + (1 - \beta)x_j$  and  $\tilde{y} = \beta y_i + (1 - \beta)y_j$ , where  $\beta \sim \text{Beta}(\gamma, \gamma)$ .

Stochastic depth (StochDepth) is a structural-based method that drops layers during training, thereby creating an implicit network ensemble of different effective lengths. Specifically, the output of layer (or residual block)  $\ell$  is given by  $\zeta_\ell = \sigma(\lambda \mathcal{F}_{\theta_\ell}(\zeta_{\ell-1}) + \mathcal{I}(\zeta_{\ell-1}))$ , where  $\lambda$  is a Bernoulli random variable,  $\mathcal{F}_{\theta_\ell}$  is the operation within the network with parameter  $\theta$  at layer  $\ell$ ,  $\mathcal{I}$  is the identity mapping operation of residual connections, and  $\sigma$  is a non-linear activation function. The keep probability is defined as  $\rho = P(\lambda = 1)$ , where in practice each layer has its own keep probability set with a linear decay rule  $\rho_\ell = 1 - \frac{\ell}{L}(1 - \rho_L)$ , with  $L$  denoting the total number of layers (or blocks) in the network.

GradAug is a recent regularization approach that combines data-level and structural techniques in a distillation-based framework. Its training loss is defined as

$$L_{GA} = L_{CE}(\mathcal{F}_\theta(x), y) + \mu \sum_{i=1}^n L_{KD}(\mathcal{F}_{\theta^{\omega_i}}(T^i(x)), \mathcal{F}_\theta(x)), \quad (2)$$

where  $\mathcal{F}_{\theta^{\omega_i}}$  denotes a slimmed sub-network of fractional width  $\omega_i$ ,  $T^i$  is a transformation performed on the input (e.g. resolution scaling), and  $\mu$  is a balancing parameter between the cross-entropy loss  $L_{CE}$  and the summed Kullback–Leibler divergence ( $L_{KD}$ ) loss on  $n$  sub-networks. The  $\omega_i$  fractional width for each sub-network is sampled from a uniform distribution between a lower bound  $\omega^b$  and 1.0 (full-width).

### 3.2. Experimental Setup

To begin our analysis, we test the accuracy of several state-of-the-art FL algorithms with several regularization methods in a common FL setting. We perform experiments using CIFAR-100 [12], an image recognition dataset with 50,000 training images across 100 categories, and employ ResNet56 [5] (as implemented in FedML [4] with PyTorch [20]) as the model. As common in the literature [1, 4, 15], the dataset is partitioned into  $K$  unbalanced subsets using a Dirichlet distribution ( $Dir(\alpha)$ ), with the default being  $\alpha = 0.5$ . With this data partitioning scheme, it is possible for a client to have no samples for one or multiple classes. Therefore, many clients will only see a portion of the total class instances. This makes the setting more realistic and challenging. For all methods and experiments we use an SGD optimizer with momentum, and a fixed learning rate of 0.01. In our basic setting, training is conducted for 25 rounds, with 16 clients and 20 local epochs per round. Any modifications to this setting in subsequent results will be stated clearly.

We compare the previously described FL algorithms and regularization methods. FedProx, MOON, and GradAug all have a hyperparameter  $\mu$  to balance their additional loss terms. We report all results with the optimal  $\mu$  for all approaches, being 0.0001, 1.0, and 1.75 for FedProx, MOON, and GradAug respectively. For Mixup and Stochastic Depth,  $\gamma$  and  $\rho_L$  are set to 0.1 and 0.9 respectively. For GradAug specifically, the number of sub-networks  $n = 2$ ,  $\omega^b = 0.8$ , and the applied transformation  $T$  is random resolution scaling. A two-layer projection layer is added to the model for MOON and the default temperature parameter  $\tau = 0.5$  as specified in the original paper. Basic data augmentations (random crop, horizontal flip, and normalization) are kept consistent across all methods.

Table 1. Results for accuracy (%) on CIFAR-100 and second-order metrics indicating the smoothness of the loss space ( $\lambda_{max}$ ,  $H_T$ ) and cross-client consistency ( $H_N$ ,  $H_D$ ) for each method.

Method	Acc. $\uparrow$	$\lambda_{max} \downarrow$	$H_T \downarrow$	$H_N \downarrow$	$H_D \uparrow$
FedAvg	52.9	297	6240	11360	0.98
FedProx	53.0	270	6132	6522	0.98
MOON	55.3	252	5520	5712	0.97
Mixup	54.0	216	5468	15434	0.99
StochDepth	55.5	215	3970	8267	0.97
GradAug	<b>57.1</b>	<b>167</b>	<b>2597</b>	2924	0.96

### 3.3. Results Comparison

The accuracy results are shown in Table 1. Within the current state-of-the-art FL algorithms (upper portion of Table 1), MOON achieves the best accuracy. This is expected, as MOON is the most intricate of the FL methods, requiring the usage of three individual models for its contrastive

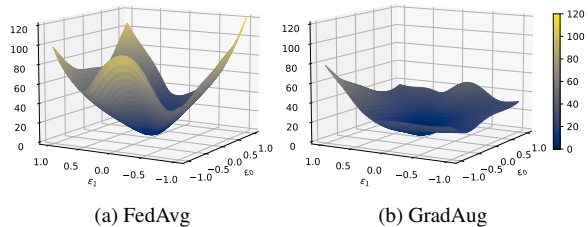


Figure 1. Visualization of the parametric loss landscape with Hessian eigenvectors  $\epsilon_0$  and  $\epsilon_1$  for each resulting global model.

learning technique. However, when we compare with standard regularization techniques (Mixup, StochDepth and GradAug in the lower portion of Table 1), we see that these perform similarly or substantially better. GradAug particularly stands out, achieving an accuracy  $\sim 2\%$  higher than MOON and  $\sim 4\%$  higher than FedAvg and FedProx. StochDepth also achieves similar accuracy to MOON. Furthermore, these regularization methods bring the same or better performance than MOON, with less memory and/or compute requirements. *We find that regularization methods appear to have an advantage in this situation; however, we wish to further investigate why this could be the case.* Next, we present our in-depth analysis based on second-order information in Section 3.4.

### 3.4. Algorithm Analysis based on Second-order Information

Recent works in the Neural Architecture Search domain [2, 35], as well as in network generalization [8, 11, 31], have noted the importance of the top Hessian eigenvalue ( $\lambda_{max}$ ) and Hessian trace ( $H_T$ ) as a predictor of performance and indicator of network generality. Having a lower  $\lambda_{max}$  and  $H_T$  typically yields a network that is less sensitive to small perturbations in the networks weights. This has the beneficial effects of smoothing the loss space during training, reaching a flatter minima, and easing convergence. These properties are particularly advantageous in federated learning, where extreme non-IID distributions and limited local data often make convergence difficult.

Motivated by these insights, we analyze the top Hessian eigenvalue and Hessian trace of the global models trained with each FL scheme to provide insight into the effectiveness of each method. As described in [30], the top Hessian eigenvalues can be approximated with the Power Iteration [31] method using a simple inner product and standard backpropagation. Furthermore, [30] also find a similar approximation for the trace utilizing the Hutchinson method [7]. We conduct our analysis with the top Hessian eigenvalues and trace of the final averaged models using these methods.

In Table 1, we include the results of the Hessian analysis. First, we find that FedAvg has the highest  $\lambda_{max}$  and



$H_T$ . FedProx and MOON each result in lower values, indicating some degree of improved generalization. However, interestingly, we find that regularization methods are most effective at reducing the  $\lambda_{max}$  and  $H_T$ , with GradAug having by far the lowest in both values. We visualize the effect of this reduction in  $\lambda_{max}$  and  $H_T$  in Fig. 1, where it can be seen that GradAug is able to smooth out the loss landscape considerably in comparison to FedAvg.

In the separate field of out-of-distribution (O.O.D.) generalization for centralized training, second-order information is being found quite useful as a theoretical indicator. Recent works [19, 21] find that forming representations that are “hard to vary” seem to result in better O.O.D. performance. More specifically, they show that the resulting loss landscapes across domains for the learned model should be consistent with each other. In terms of theoretical indicators, this translates to matching domain-level Hessians, as the Hessian provides an approximation of local curvature. Similarly, in federated learning, each client is essentially a separate domain. Therefore, matching Hessians in norm and direction across clients reveals additional detail and reasoning behind the effectiveness of each method. In light of these findings in O.O.D. literature, we analyze the difference in Hessian norm ( $H_N$ ) and the Hessian direction across clients ( $H_D$ ), where

$$H_N^{k,j} = (\|\text{Diag}(\mathbf{H}_k)\|_F - \|\text{Diag}(\mathbf{H}_j)\|_F)^2 \text{ and} \quad (3)$$

$$H_D^{k,j} = \frac{\text{Diag}(\mathbf{H}_k) \odot \text{Diag}(\mathbf{H}_j)}{\|\text{Diag}(\mathbf{H}_k)\|_F \cdot \|\text{Diag}(\mathbf{H}_j)\|_F}. \quad (4)$$

Here,  $\odot$  is the dot product,  $\mathbf{H}_k$  and  $\mathbf{H}_j$  are the Hessian matrices of clients  $k$  and  $j$ , and  $\|\cdot\|_F$  is the Frobenius norm.  $H_N^{k,j}$  and  $H_D^{k,j}$  are averaged across all pairs of clients and reported as simply  $H_N$  and  $H_D$  in Table 1. For these Hessian matching criteria, a lower  $H_N$  (less difference) and a higher  $H_D$  (essentially the cosine similarity) are desired.

As seen on the right side of Table 1,  $H_D$  is fairly consistent across all methods. In terms of  $\lambda_{max}$ ,  $H_T$ , and  $H_D$ , most methods seem to correlate decently well between these values and performance. However, there are a few cases which require more information. First, Mixup has a similar  $H_T$  value as MOON, but lower accuracy.  $H_N$  provides another detail; the Hessian norms of Mixup are not nearly as similar across clients as those of MOON. Between MOON and StochDepth, we see that MOON has both a higher  $\lambda_{max}$  and  $H_T$ , but StochDepth has a higher  $H_N$ . In the end, MOON and StochDepth result in similar performance, with perhaps a slight edge towards the latter.

**Key Insight.** It appears that both the eigenvalue/trace analysis and Hessian matching criteria can serve as a guiding indicator for optimal FL methods. Particularly, they provide insight into the facilitation of convergence and aggregation thorough landscape smoothness and consistency. To

understand how these differences will play out empirically, we conduct a variety of ablations in Section 3.5.

### 3.5. Ablation Study under Various FL Settings

**Data Heterogeneity.** Federated systems can be deployed with many different setups and diverse environments. We conduct further analysis across a variety of FL settings to ensure the generality of our findings. First, we examine the effect of varying the degree of heterogeneity in the client data distributions. The results are shown in Table 2. We report the mean accuracy  $\pm$  the standard deviation across three runs. All other settings are maintained from Section 3.2; only the data distribution  $Dir(\alpha)$  is varied. A lower  $\alpha$  value indicates a more heterogeneous distribution.

Table 2. Ablation results for varying degrees of data heterogeneity.

Method	$\alpha = 0.1$	$\alpha = 0.5$	$\alpha = 2.5$	homog
FedAvg	45.0 $\pm$ 0.2	52.9 $\pm$ 0.1	54.4 $\pm$ 0.2	54.9 $\pm$ 0.4
FedProx	45.2 $\pm$ 0.3	53.1 $\pm$ 0.3	54.5 $\pm$ 0.3	54.8 $\pm$ 0.5
MOON	46.5 $\pm$ 0.5	55.0 $\pm$ 0.5	56.3 $\pm$ 0.6	56.3 $\pm$ 0.5
Mixup	44.3 $\pm$ 0.1	54.0 $\pm$ 0.1	55.5 $\pm$ 0.4	56.7 $\pm$ 0.4
StochDepth	48.2 $\pm$ 0.3	55.5 $\pm$ 0.2	57.6 $\pm$ 0.2	58.1 $\pm$ 0.6
GradAug	<b>48.6<math>\pm</math>0.4</b>	<b>57.0<math>\pm</math>0.1</b>	<b>59.6<math>\pm</math>0.2</b>	<b>60.5<math>\pm</math>0.2</b>

As the degree of data heterogeneity decreases, the effect of client drift should become less significant. Therefore, we expect that the accuracy for each method will increase, with peak performance in the homogeneous setting. All regularization methods, as well as FedAvg, perform as expected, and find consistent improvement across the degrees of data distribution. However, we see that the accuracy improvement of FedProx and MOON slows as the data approaches homogeneity, with accuracy in the purely homogeneous setting (“homog” in Table 2) remaining quite low. In their attempt to mitigate client drift and keep local updates close to the global model, it appears that they also hinder their ability to fully learn on minorly heterogeneous or even homogeneous data. This is not ideal for deployable FL systems, as the degree of heterogeneity is not known ahead of time. Moreover, even in the most heterogeneous cases, the structural regularization methods perform better than the standard FL algorithms. For instance, StochDepth achieves a  $\sim 1.7\%$  improvement over MOON at  $\alpha = 0.1$ , while also having improvement in more homogeneous situations. In all settings, GradAug performs the best.

**Number of Local Training Epochs.** The main purpose for adequately handling data heterogeneity is to allow for more productive training on the client each round, therefore reducing the time to convergence and required communication cost. Therefore, to examine the training productivity of each method, we examine their accuracy with various allotted local training epochs per round ( $E$ ) in Table 3.

Ideally methods should continue to improve in accuracy with more allotted local training epochs. In Table 3, we see

Table 3. Ablation results for number of local training epochs.

Method	$E = 10$	$E = 20$	$E = 30$
FedAvg	50.6±0.1	52.9±0.1	53.2±0.3
FedProx	50.7±0.5	53.1±0.3	52.8±0.1
MOON	50.7±0.4	55.0±0.5	55.2±0.4
Mixup	50.5±0.4	54.0±0.1	54.4±0.3
StochDepth	50.9±0.6	55.5±0.2	56.4±0.3
GradAug	<b>53.5±0.3</b>	<b>57.0±0.1</b>	<b>57.7±0.3</b>

Table 4. Ablation results for varying number of clients  $C$  in synchronous and client sampling cases.

Method	$C = 16$	$C = 32$	$C = 64$	$C = 64 \times 0.25$	$C = 64 \times 0.25 (100)$
FedAvg	52.9±0.1	44.5±0.3	34.6±0.2	32.7±0.5	46.5±0.6
FedProx	53.1±0.3	44.5±0.6	34.8±0.2	32.5±0.4	46.2±0.1
MOON	55.0±0.5	45.8±0.3	35.2±0.8	34.2±0.2	49.5±0.7
Mixup	54.0±0.1	46.0±0.1	36.0±0.2	33.6±0.6	49.1±0.2
StochDepth	55.5±0.2	47.5±0.2	35.5±0.6	34.6±0.1	51.4±0.1
GradAug	<b>57.0±0.1</b>	<b>50.4±0.1</b>	<b>40.2±0.1</b>	<b>38.1±0.3</b>	<b>53.3±0.5</b>

that all methods steadily improve from 10 epochs per round to 20. However, from 20 to 30, the trends vary considerably. As a baseline, FedAvg slightly improves by  $\sim 0.3\%$ . Surprisingly, FedProx and MOON stay relatively stagnant from 20 to 30 epochs. Meanwhile, the standard (particularly structural) regularization methods continue to increase in accuracy. Therefore, these methods illustrate the ability to maintain productive training, even across a wide range of allotted local epochs.

**Number of Clients.** In real-world FL settings, the number of participating clients can vary widely. Moreover, only a portion of clients are potentially sampled per round, whether for connectivity reasons or other capacity restrictions of the central system. Therefore, it is crucial that an FL method can converge under such conditions. We study the affect of client number and client sampling in Table 4.  $C = 64 \times 0.25$  indicates that there are 64 total clients in the system, but only a fraction (0.25) are sampled each round. The rest of the presented results in Table 4 sample all  $K$  clients each round.  $C = 64 \times 0.25 (100)$  is run for 100 rounds, and all other settings for the default 25 rounds.

The trends of most methods are similar with increasing clients. However, FedProx struggles to keep up with the FedAvg baseline, especially in the client sampling cases. These scenarios are particularly important; when a small percentage of clients are sampled, only a portion of the dataset is effectively trained on each round. Therefore, learning efficiency becomes paramount for maintaining suitable convergence. The standard regularization methods maintain better accuracy than FedAvg in all settings, often by a significant margin, and even in the client sampling scenario. Overall, GradAug performs the best in all cases. *Therefore, even though these regularization methods were not designed for the FL setting and partial client sampling, they still perform on par with or improve over current state-of-the-art FL algorithms.*

## 4. Proposed Method – FedAlign

Overall, we find that GradAug is particularly effective in the FL setting, having the highest accuracy in all tested scenarios along with the lowest  $\lambda_{max}$ ,  $H_T$ , and  $H_N$ . However, while this method is quite memory efficient in comparison to many FL methods (only requires a single stored model during training), it does incur a substantial increase in training time and local computation over the FedAvg baseline. This is because GradAug requires multiple forward passes through slimmed sub-networks for the distillation loss. It is possible to reduce the computation burden to some extent by using a smaller number of sub-networks during the knowledge distillation process, as seen in Table 5. Here, the  $\mu$  in GradAug is adjusted to 2.0, 1.5, and 1.25 for  $n = 1, 3$ , and 4, respectively. Nonetheless, a considerable gap still remains between GradAug and vanilla FedAvg in local compute requirements and subsequent wall-clock time. **Therefore, the question is, can we devise a method which provides similar effect and performance as GradAug in FL, but with substantially less computational overhead?** This is particularly important in the FL setting, where clients are typically deployed devices with minimal memory and computational resources.

Table 5. Analysis of local compute, stored parameters, and wall-clock time. FLOPs are calculated for the compute needs for the forward pass of the training process. Parameters include the total number of stored parameters needed for each method during training. Wall-clock time is measured as a per-round average on CIFAR-100 with  $C=16$  and  $E=20$  across 4 RTX-2080Ti GPUs.

Method	Acc (%) $\uparrow$	MFLOPs $\downarrow$	Param (M) $\downarrow$	Time (s)
FedAvg	52.9±0.1	87.3	<b>0.61</b>	137.2
FedProx	53.1±0.3	87.3	1.21	161.9
MOON	55.0±0.5	262.2	2.21	414.2
Mixup	54.0±0.1	87.3	<b>0.61</b>	137.8
StochDepth	55.5±0.2	<b>82.4</b>	<b>0.61</b>	<b>136.7</b>
GradAug ( $n = 1$ )	56.7±0.3	133.9	<b>0.61</b>	229.2
GradAug ( $n = 2$ )	<b>57.0±0.1</b>	170.7	<b>0.61</b>	323.9
GradAug ( $n = 3$ )	56.8±0.3	217.4	<b>0.61</b>	417.7
GradAug ( $n = 4$ )	56.9±0.3	264.1	<b>0.61</b>	514.4
<b>FedAlign</b>	56.9±0.5	89.1	<b>0.61</b>	166.2

To do so, we first take note of the following insights gathered during our analysis: 1) Second-order information is insightful for understanding the learning generality of neural networks. Particularly, we find that flatness and consistency in this realm are desirable traits. 2) In practice, we find that structural regularization, and especially distillation-based like GradAug, is quite effective. Furthermore, the weight sharing mechanisms of such approaches are memory efficient compared to other methods that rely on global model or previous model storage. Therefore, we combine these insights into a novel algorithm to optimize for performance and resource needs in FL.

We propose **FedAlign**, a distillation-based regulariza-

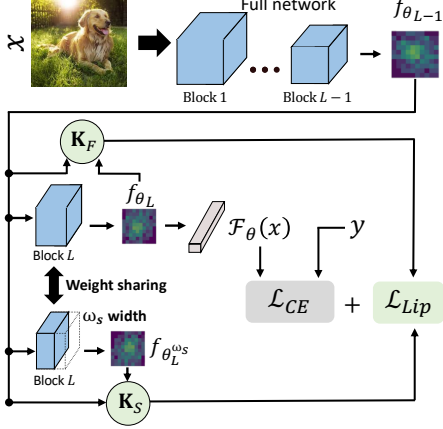


Figure 2. The proposed FedAlign for local client training in FL. Features  $f_{\theta_{L-1}}$  are run through Block  $L$  as normal. The only additional inference in FedAlign is through Block  $L$  at a reduced width (i.e. sub-block), reusing features  $f_{\theta_{L-1}}$  as input. The channels throughout the layers in the sub-block are a  $\omega_S$  fraction of the original number. This is accomplished via temporary uniform pruning of Block  $L$ .

tion method that aligns the Lipschitz constants (i.e. top Hessian eigenvalues) of the most critical network components through the use of slimmed sub-blocks. Fig. 2 shows an overview of FedAlign, whose design is based on two key principles. First, motivated by the insights of Section 3.4, we internally regularize the Lipschitz constants of network blocks to *promote smooth optimization and consistency within the model*. Recent work [24] presents a quick approximation of the Lipschitz constants for neural network layers in a differentiable manner. This enables the use of second-order information in the distillation process, traditionally between a fully trained teacher and a learning student. We adapt this technique for distillation-based regularization with a *single untrained network* in place of the traditional logit-based loss. Second, in order to reduce computation in a purposeful manner, we take note of certain network properties. Particularly, it has been shown that the final layers of a neural network are most prone to overfit to the client distribution [16]. Therefore, we design FedAlign with a focus on these critical points in the network. *The question we raise is, when aiming to concentrate our regularization efforts on the final layers, why should we run all sub-networks for distillation from start to finish?* Instead, we propose to reuse the intermediate features of the full network as input to just the final block at a reduced width, and therefore significantly reduce computation. In this way, we harness the benefits of distillation-based regularization in performance and memory footprint, while effectively mitigating computational overhead.

Combining these two key principles, we form the FedAlign local objective as

$$\mathcal{L}_{FA} = \mathcal{L}_{CE}(\mathcal{F}_{\theta}(x), y) + \mu \mathcal{L}_{Lip}(\mathbf{K}_S, \mathbf{K}_F), \quad (5)$$

## Algorithm 1 FedAlign

### SERVER OPERATIONS

**Inputs:** Round number  $R$ , Set of clients  $S$   
**Output:** Final global model weights  $\theta_{global}^R$

```

Initialize model weights  $\theta_{global}^0$ 
for  $r = 0, 1, \dots, R - 1$  do
  Sample available clients  $C$  from  $S$ 
  for client  $c \in C$  in parallel do
     $\theta_c^r \leftarrow \text{CLIENTOPERATIONS}(\theta_{global}^r)$ 
  end for
   $\theta_{global}^{r+1} \leftarrow \sum_{c=1}^C \frac{n_c}{n} \theta_c^r$ 
end for

```

### CLIENT OPERATIONS

**Input:** Model weights  $\theta_{global}$   
**Output:** Updated local model weights  $\theta$

```

Load received weights  $\theta_{global}$  to local model  $\mathcal{F}_{\theta}$ 
for epoch  $e = 0, 1, \dots, E - 1$  do
  for batch  $\{x, y\} \in D$  do
     $f_{\theta_{L-1}}, f_{\theta_L}, pred = \mathcal{F}_{\theta}(x)$ 
     $f_{\theta_L^{\omega_S}} = \mathcal{F}_{\theta_L^{\omega_S}}(f_{\theta_{L-1}})$ 
     $\mathbf{X}_S, \mathbf{X}_F = TM(f_{\theta_L^{\omega_S}}, f_{\theta_{L-1}}, f_{\theta_L})$ 
     $\mathbf{K}_S, \mathbf{K}_F = \|\mathbf{X}_S\|_{SN}, \|\mathbf{X}_F\|_{SN}$ 
     $\mathcal{L}_{FA} = \mathcal{L}_{CE}(pred, y) + \mu \mathcal{L}_{Lip}(\mathbf{K}_S, \mathbf{K}_F)$ 
     $\theta \leftarrow \text{update}(\theta, \mathcal{L}_{FA})$ 
  end for
end for
Send updated local model weights  $\theta$  to server

```

where  $\mu$  is a balancing constant,  $\mathcal{L}_{CE}$  is the cross-entropy loss, and  $\mathcal{L}_{Lip}$  is the mean squared error between the approximated Lipschitz constant vectors  $\mathbf{K}_S$  and  $\mathbf{K}_F$  for the reduced width (i.e. sub-block) and full width block  $L$ , respectively. Specifically, the Lipschitz approximations are calculated via the spectral norm of a transmitting matrix using feature maps as in [24], which bypasses the need for singular value decomposition. Therefore, we use the intermediate features for these transmitting matrices  $\mathbf{X}_F$  and  $\mathbf{X}_S$ , where

$$\mathbf{X}_F = (f_{\theta_{L-1}})^{\top} f_{\theta_L}, \text{ and} \quad (6)$$

$$\mathbf{X}_S = (f_{\theta_{L-1}})^{\top} f_{\theta_L^{\omega_S}}. \quad (7)$$

$f_{\theta_L}$  and  $f_{\theta_{L-1}}$  are the feature maps outputted by the last and prior-to-last blocks of the full network  $\mathcal{F}_{\theta}(x)$ ;  $f_{\theta_L^{\omega_S}}$  is the output feature map of the final block  $L$  at reduced width  $\omega_S$  (see Fig. 2). Finally, the spectral norm (SN) of  $\mathbf{X}_F$  and  $\mathbf{X}_S$  are approximated using the Power Iteration method [31], and therefore  $\mathbf{K}_F = \|\mathbf{X}_F\|_{SN}$  and  $\mathbf{K}_S = \|\mathbf{X}_S\|_{SN}$ . A pseudocode implementation of FedAlign is presented in Alg. 1. Looking back to Eq. 5, one could view  $\mathcal{L}_{Lip}$  as a correction term; however, there is a key distinction between this form of regularization and that of traditional FL algorithms. *Our correction term promotes the local client models to learn well-generalized representations based on their own data, instead of forcing the local models to be close to the global model.*

As seen in Table 5, FedAlign achieves state-of-the-art accuracy in a resource-efficient manner. With just a

Table 6. FedAlign ablation results on CIFAR-100.

Method	$\alpha = 0.1$	$\alpha = 2.5$	homog	$E = 10$	$E = 30$	$C = 32$	$C = 64$	$C = 64 \times 0.25$	$C = 64 \times 0.25 (100)$
FedAlign	48.7±0.2	57.6±0.6	58.2±0.1	51.2±0.3	57.9±0.6	47.8±0.3	36.5±0.1	34.9±0.6	50.9±0.5

Table 7. CIFAR-10 and ImageNet-200 results for all methods.

Method	CIFAR-10				ImageNet-200			
	$C = 16$	$C = 64 \times 0.25 (100)$	MFLOPs ↓	Param (M) ↓	$C = 16$	$C = 32 \times 0.125 (50)$	GFLOPs ↓	Param (M) ↓
FedAvg	81.9±0.6	78.9±0.3	87.3	0.61	60.7±0.4	52.7±0.2	18.1	11.22
FedProx	81.9±0.2	78.9±0.7	87.3	1.21	61.0±0.4	52.5±0.3	18.1	22.42
MOON	82.9±0.4	79.4±0.5	262.2	2.21	61.1±0.2	54.3±0.2	54.4	19.96
Mixup	80.3±0.4	80.5±0.5	87.3	0.61	61.0±0.3	52.3±0.3	18.1	11.22
StochDepth	82.2±0.2	80.8±0.7	82.4	0.61	60.5±0.2	52.9±0.2	17.3	11.22
GradAug ( $n = 2$ )	84.6±0.6	83.8±0.3	170.7	0.61	63.5±0.4	55.6±0.1	34.4	11.22
GradAug ( $n = 1$ )	84.0±0.2	82.3±0.5	133.9	0.61	62.8±0.3	54.4±0.4	25.3	11.22
<b>FedAlign</b>	82.3±0.3	82.3±0.3	89.1	0.61	62.0±0.1	55.1±0.5	19.3	11.22

1.02x difference in FLOPs, FedAlign realizes a significant  $\sim 4.0\%$  accuracy improvement over the FedAvg baseline. For the FL algorithms FedProx and MOON, they not only have much lower accuracy than FedAlign, but also require substantially more compute and/or memory. Particularly, FedAlign achieves a  $\sim 1.9\%$  accuracy improvement over MOON, while reducing the local compute overhead by over 65% and the memory requirements by over 70%. Furthermore, FedAlign realizes a critical  $\sim 47\%$  and  $\sim 33\%$  reduction in compute needs compared to GradAug with ( $n = 2$ ) and ( $n = 1$ ), without sacrificing accuracy.

#### 4.1. FedAlign Experiments

We further verify the effectiveness of our method across various settings and datasets. In Table 6, we examine the performance of FedAlign with the same ablations as in Section 3.5, where FedAlign exhibits strong performance in many settings. We also investigate FedAlign and all other methods across two additional datasets: CIFAR-10 and ImageNet-200. For ImageNet-200, we randomly sample 200 classes from the classic ImageNet-1k [22] dataset. We employ ResNet56 and ResNet18 [5] as our models on CIFAR-10 and ImageNet-200, respectively. For FedAlign,  $\omega_S = 0.25$  and  $\mu = 0.45$  in all results. Hyperparameters for all other methods are those described in Section 3.2 (with  $\mu = 2.0$  for GradAug ( $n = 1$ ) as in Table 5). For additional analyses, please refer to the supplementary material.

For CIFAR-10, we ran a 16 client synchronous and 64 client case with sampling in Table 7. We note similar trends to CIFAR-100; regularization methods perform well, particularly in the more realistic client sampling case. On ImageNet-200, we also ran synchronous and sampling settings. Here, both GradAug and FedAlign maintain higher performance than other methods. FedAlign provides competitive accuracy with GradAug ( $n = 1$ ) and even ( $n = 2$ ) in the sampling case, while reducing computational needs by a significant margin. Interestingly, StochDepth does not perform as well in the ImageNet-200 cases. As mentioned in the original paper [6], Stochastic Depth performs better

with deeper networks. However, with ResNet18, the overall depth of the network is reduced compared to that in the CIFAR cases. Therefore, as most deployable networks favor width over depth, regularizing with respect to the width of a network is more applicable to the FL setting. *This highlights an additional benefit of FedAlign, which operates using width reduction in the final block and maintains relatively high accuracy despite low resource needs.*

## 5. Conclusion and Discussion

In this work, we study the data heterogeneity challenge of FL from a simple yet unique perspective of local learning generality. To this end, we present a thorough study of various methods in FL settings, and further propose FedAlign, which achieves competitive SOTA accuracy with excellent resource efficiency. One limitation of our study is that we only focused on image tasks and models for the experiments. Natural language processing applications of FL are also a common setting, and therefore could be explored in future work. Nonetheless, we note that FedAlign can easily be applied to language applications, as it operates in the feature space and does not have a fundamental reliance on the input type. On the other hand, GradAug is primarily designed for vision data, employing a random transformation and applying it to the input of sub-networks.

While no one presented regularization method is perfect in all respects, we emphasize that local learning is extremely important in federated settings. Furthermore, methods that particularly focus on promoting learning generality inherently improve global FL aggregation and optimization to a surprising degree. By introducing methods like GradAug in FL, we propose a rethinking of federated optimization and how to tackle its challenges. As a step further in this direction, FedAlign provides strong improvement over classic baselines and state-of-the-art FL methods while addressing the local computational restraints of an FL system.

**Acknowledgement:** This work is supported by the NSF/Intel Partnership on MLWiNS under Grant No. 2003198.



## References

- [1] Durmus Alp Emre Acar, Yue Zhao, Ramon Matas, Matthew Mattina, Paul Whatmough, and Venkatesh Saligrama. Federated learning based on dynamic regularization. In *International Conference on Learning Representations*, 2021. 1, 3, 4, 11
- [2] Xiangning Chen and Cho-Jui Hsieh. Stabilizing differentiable architecture search via perturbation-based regularization. In *International Conference on Machine Learning*, pages 1554–1565. PMLR, 2020. 4
- [3] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. *Advances in Neural Information Processing Systems*, 31:10727–10737, 2018. 3
- [4] Chaoyang He, Songze Li, Jinhyun So, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Praneeth Vepakomma, Abhishek Singh, Hang Qiu, Li Shen, Peilin Zhao, Yan Kang, Yang Liu, Ramesh Raskar, Qiang Yang, Murali Annamalai, and Salman Avestimehr. Fedml: A research library and benchmark for federated machine learning. *arXiv preprint arXiv:2007.13518*, 2020. 4, 11
- [5] 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. 4, 8
- [6] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *European conference on computer vision*, pages 646–661. Springer, 2016. 3, 8
- [7] M.F. Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communication in Statistics- Simulation and Computation*, 18:1059–1076, 01 1989. 4
- [8] Yiding Jiang\*, Behnam Neyshabur\*, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2020. 4
- [9] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, and et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021. 1
- [10] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020. 1, 3
- [11] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. 2016. 4
- [12] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research). 4
- [13] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>. 1
- [14] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. *arXiv preprint arXiv:2102.02079*, 2021. 1
- [15] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10713–10722, 2021. 1, 2, 3, 4, 11
- [16] 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. *arXiv e-prints*, pages arXiv–2106, 2021. 7
- [17] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 1, 2, 11
- [18] Seungeun Oh, Jihong Park, Eunjeong Jeong, Hyesung Kim, Mehdi Bennis, and Seong-Lyun Kim. Mix2fld: Downlink federated learning after uplink federated distillation with two-way mixup. *IEEE Communications Letters*, 24(10):2211–2215, 2020. 3
- [19] Giambattista Parascandolo, Alexander Neitz, Antonio Orvieto, Luigi Gresele, and Bernhard Schölkopf. Learning explanations that are hard to vary. *arXiv preprint arXiv:2009.00329*, 2020. 2, 5
- [20] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 4
- [21] Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant gradient variances for out-of-distribution generalization. *arXiv preprint arXiv:2109.02934*, 2021. 2, 5
- [22] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 8
- [23] Anit Kumar Sahu, Tian Li, Maziar Sanjabi, Manzil Zaheer, Ameet Talwalkar, and Virginia Smith. On the convergence of federated optimization in heterogeneous networks. *CoRR*, abs/1812.06127, 2018. 1, 2
- [24] Yuzhang Shang, Bin Duan, Ziliang Zong, Liqiang Nie, and Yan Yan. Lipschitz continuity guided knowledge distillation, 2021. 7, 12
- [25] MyungJae Shin, Chihoon Hwang, Joongheon Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Xor mixup: Privacy-preserving data augmentation for one-shot federated learning. *arXiv preprint arXiv:2006.05148*, 2020. 3
- [26] Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun, and Christoph Bregler. Efficient object localization using convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 648–656, 2015. 3

- [27] Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In *International Conference on Learning Representations*, 2019. 3
- [28] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 2020. 3
- [29] Taojiannan Yang, Sijie Zhu, and Chen Chen. Gradaug: A new regularization method for deep neural networks. *Advances in Neural Information Processing Systems*, 33, 2020. 3
- [30] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. Pyhessian: Neural networks through the lens of the hessian. In *2020 IEEE International Conference on Big Data (Big Data)*, pages 581–590. IEEE, 2020. 4
- [31] Zhewei Yao, Amir Gholami, Qi Lei, Kurt Keutzer, and Michael W Mahoney. Hessian-based analysis of large batch training and robustness to adversaries. *Advances in Neural Information Processing Systems*, 31:4949–4959, 2018. 4, 7
- [32] Tehrim Yoon, Sumin Shin, Sung Ju Hwang, and Eunho Yang. Fedmix: Approximation of mixup under mean augmented federated learning. In *International Conference on Learning Representations*, 2021. 3
- [33] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032, 2019. 3
- [34] Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang, and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In *International Conference on Machine Learning*, pages 7252–7261. PMLR, 2019. 3
- [35] Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox, and Frank Hutter. Understanding and robustifying differentiable architecture search. In *International Conference on Learning Representations*, 2020. 4
- [36] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. 3

## Supplementary Material

The supplementary material is organized into the following sections:

- Section **A**: Analysis of both communication and compute efficiency for all explored methods.
- Section **B**: Second-order analysis of FedAlign.
- Section **C**: Hyperparameter ablations for FedAlign.
- Section **C**: Details and visualization for the non-IID data partitioning scheme.
- Section **E**: Implementation details for transmitting matrices and FedAlign training.

### A. Communication and Compute Efficiency

Communication cost is another critical factor in FL systems, as participating client devices are often on slow or congested networks [17]. Therefore, total efficiency in FL systems includes both the ability to reduce the local computational burden, as well as the communication overhead. We evaluate all methods with such measures in Table 8. We maintain the CIFAR-100 setting described in Section 3.2 of the main paper, except that we do not limit the number of rounds, but rather allow all methods to reach a common accuracy of 60%. This allows us to analyze the total costs of each method consistently. FedAlign proves to be the most efficient in all respects, achieving the target accuracy in less number of rounds, less communication cost, and less local computation.

Table 8. Number of rounds (Rounds), local compute (MFLOPs), and communication cost (Comm Cost) required by each method to achieve 60% accuracy on CIFAR-100. Local computation is computed as the sum total over all nodes and samples for all completed rounds. Communication costs are calculated by the number of parameters of the model transferred as 32 bit weights with all completed rounds.

Method	Rounds	MFLOPs	Comm Cost (Gb)
FedAvg	84	7332	26.23
FedProx	78	6809	24.36
MOON	55	14421	17.18
Mixup	71	6198	22.17
StochDepth	46	3790	14.37
GradAug ( $n = 2$ )	41	6999	12.81
GradAug ( $n = 1$ )	44	5892	13.74
<b>FedAlign</b>	<b>37</b>	<b>3297</b>	<b>11.56</b>

### B. Second-order Analysis

In Table 9, we show the second-order analysis results for FedAlign along with the other methods. The Liptzshitz-focused distillation loss of FedAlign effectively reduces the Lipshitz constant  $\lambda_{max}$  considerably as intended (FedAlign results in the lowest  $\lambda_{max}$  across all

methods), and therefore helps provide stronger generalization and performance. However, one limitation to FedAlign is that it does not directly translate to a strong reduction in  $H_N$ . Therefore, a promising direction for future work could extend FedAlign to also consider this aspect.

### C. Hyperparameter Ablations of FedAlign

The default hyperparameter setting used throughout the paper is  $\omega_S = 0.25$  and  $\mu = 0.45$ . The performance of FedAlign with various hyperparameters on the CIFAR-100 basic setting (described in Section 3.2 of the main paper) is shown in Table 10. We vary  $\mu$  and  $\omega_S$  independently, meaning  $\omega_S = 0.25$  for the  $\mu$  ablations, and  $\mu = 0.45$  when varying  $\omega_S$ . Table 10 shows that FedAlign is more sensitive to  $\omega_S$  than  $\mu$ ; nonetheless, we found  $\omega = 0.25$  to be a versatile choice in practice. Furthermore, hyperparameters only need to be decided once, as they transfer well across a variety of other datasets and FL settings as shown in Section 4.1 of the main paper.

Table 9. Results for accuracy (%) on CIFAR-100 and second-order metrics indicating the smoothness of the loss space ( $\lambda_{max}$ ,  $H_T$ ) and cross-client consistency ( $H_N$ ,  $H_D$ ) for each method.

Method	Acc. $\uparrow$	$\lambda_{max}\downarrow$	$H_T\downarrow$	$H_N\downarrow$	$H_D\uparrow$
FedAvg	52.9	297	6240	11360	0.98
FedProx	53.0	270	6132	6522	0.98
MOON	55.3	252	5520	5712	0.97
Mixup	54.0	216	5468	15434	<b>0.99</b>
StochDepth	55.5	215	3970	8267	0.97
GradAug ( $n = 2$ )	<b>57.1</b>	167	<b>2597</b>	2924	0.96
GradAug ( $n = 1$ )	56.8	179	3620	<b>2607</b>	0.97
<b>FedAlign</b>	56.7	<b>143</b>	4409	9655	<b>0.99</b>

Table 10. FedAlign hyperparameter ablations on with CIFAR-100

Method	$\mu = 0.35$	$\mu = 0.45$	$\mu = 0.55$	$\omega_S = 0.1$	$\omega_S = 0.25$	$\omega_S = 0.4$
FedAlign	56.0	<b>56.7</b>	56.1	54.9	<b>56.7</b>	55.2

### D. Data Partitioning

As is common in the literature [1, 4, 15], we partition the employed datasets into  $K$  unbalanced subsets using a Dirichlet distribution  $Dir(\alpha)$ . The distribution for all three datasets at  $\alpha = 0.5$  is visualized in Fig. 3 (a), (c), and (d). Additionally, (b) shows the distribution for CIFAR-100 with  $\alpha = 0.1$  as studied in Section 3.5 of the main paper. Overall, we see that the number of samples for each class varies considerably across clients, and often times a client will not have any samples from multiple classes. This enhances the FL setting by making it more realistic and challenging. For implementation, we utilize the same data partitioning script as that in [4].

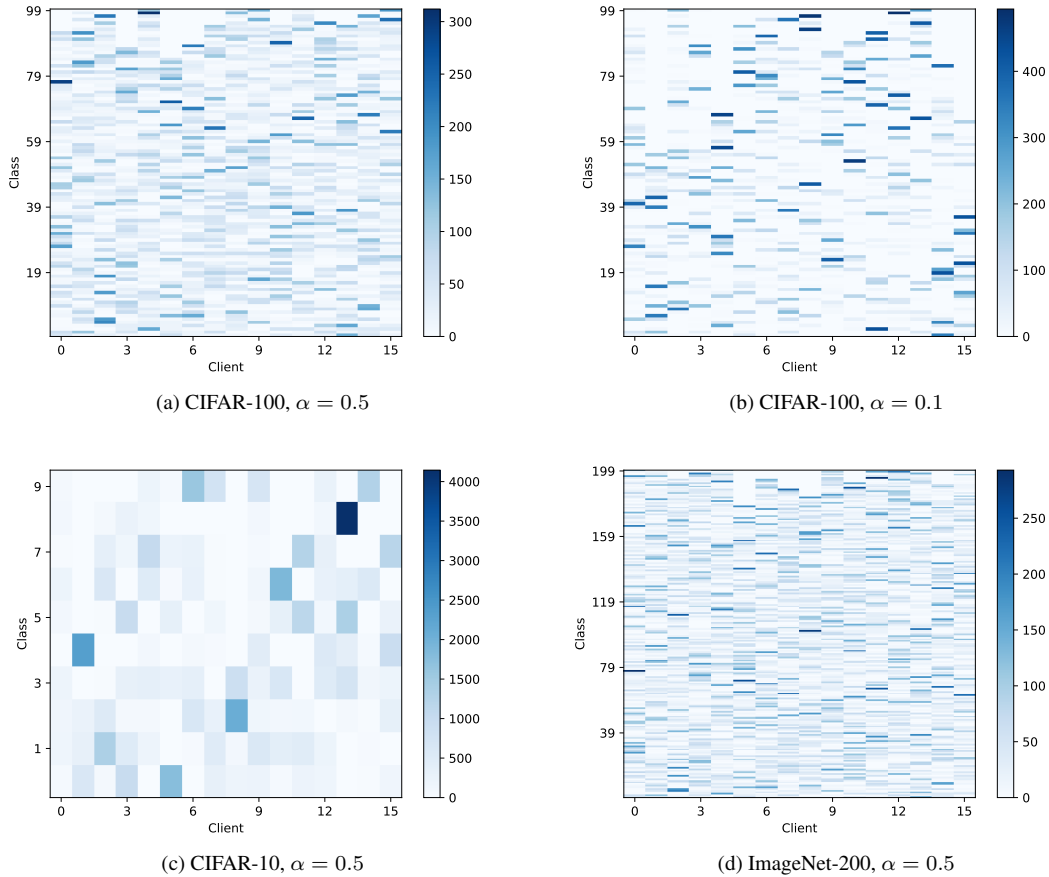


Figure 3. Data distribution visualization for  $Dir(\alpha)$  and  $C = 16$  across multiple datasets. Each column shows the number of samples per class allocated to a client.

### E. Additional Implementation Details

When calculating  $\mathbf{X}_F$  and  $\mathbf{X}_S$ , the input and output features involved will typically be of different spatial sizes in practice. Therefore, [24] utilizes an adaptive average pool operation in PyTorch to reduce the spatial size of the larger feature map to that of the smaller one. We likewise employ this operation.

Prior to performing backpropagation, we apply a relative scale to  $\mathcal{L}_{Lip}$  along with the  $\mu$  scaling parameter. In PyTorch-style pseudocode:  $loss_{lip} = \mu * (loss_{ce.item()} / loss_{lip.item()}) * loss_{lip}$ . This is to ensure that  $\mathcal{L}_{Lip}$  is on relatively the same scale with  $\mathcal{L}_{CE}$ . A gradient clip is also applied.