

# Beyond Local Sharpness: Communication-Efficient Global Sharpness-aware Minimization for Federated Learning

Debora Caldarola<sup>1,\*</sup>, Pietro Cagnasso<sup>1</sup>, Barbara Caputo<sup>1</sup>, Marco Ciccone<sup>2,†</sup>

<sup>1</sup>Politecnico di Torino, <sup>2</sup> Vector Institute

## Abstract

Federated learning (FL) enables collaborative model training with privacy preservation. Data heterogeneity across edge devices (*clients*) can cause models to converge to sharp minima, negatively impacting generalization and robustness. Recent approaches use client-side sharpness-aware minimization (SAM) to encourage flatter minima, but the discrepancy between local and global loss landscapes often undermines their effectiveness, as optimizing for local sharpness does not ensure global flatness. This work introduces FEDGLOSS (Federated Global Server-side Sharpness), a novel FL approach that prioritizes the optimization of global sharpness on the server, using SAM. To reduce communication overhead, FEDGLOSS cleverly approximates sharpness using the previous global gradient, eliminating the need for additional client communication. Our extensive evaluations demonstrate that FEDGLOSS consistently achieves flatter minima and better performance compared to state-of-the-art FL methods in various federated vision benchmarks. Code available at [github.com/pietrocagnasso/fedgloss](https://github.com/pietrocagnasso/fedgloss).

## 1. Introduction

Federated Learning (FL) [39] provides a powerful framework to collaboratively train machine learning models on private data distributed across multiple endpoints. Unlike traditional methods, FL enables edge devices (*clients*), like smartphones or IoT (Internet of Things) hardware, to train a shared model without compromising their sensitive information. This is achieved through communication rounds, where clients independently train on their local data and then exchange updated model parameters with a central server, preserving data privacy. The optimization on the server side relies on *pseudo-gradients* [47], i.e., the average of the difference between the global model and the client's update, which serve as an estimate of the true global gradient

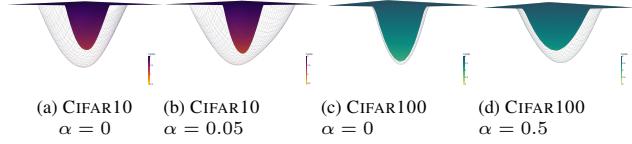


Figure 1. Comparison of **FEDAVG** (solid) and **FEDSAM** (net) loss landscapes with varying degrees of data heterogeneity ( $\alpha$ ) on the CIFAR datasets. **FEDSAM**'s effectiveness in converging to **global flat minima** is highly influenced by the data heterogeneity, where higher heterogeneity ( $\alpha \rightarrow 0$ ) leads to sharper minima, and the complexity of the task, e.g., higher sharpness for the more complex CIFAR100. This highlights the importance of optimizing global sharpness. Model: CNN.

on the overall dataset. This approach holds immense potential for privacy-sensitive applications, proving its value in areas like healthcare [2, 37, 41, 46], finance [41], autonomous driving [11, 40, 51], IoT [62], and more [32, 57]. However, the real-world deployment of FL presents unique challenges stemming from data heterogeneity and communication costs [34]. Clients gather their data influenced by various factors such as personal habits or geographical locations, leading to inherent differences across devices [20, 24, 51]. This results in the global model suffering from degraded performance and slower convergence [4, 25, 26, 36], with instability emerging as client-specific optimization paths diverge from the global one. This phenomenon, known as *client drift* [26], limits the model's ability to generalize to the overall underlying distribution.

While many FL approaches focus on mitigating client drift through client-side regularization [1, 35, 56], a recent trend leverages the geometry of the loss landscape to improve generalization [4, 7, 45, 53, 54]. These methods build upon the notion that convergence to sharp minima correlates with poor generalization [17, 23, 27]. FEDSAM [4, 45] employs Sharpness-aware Minimization (SAM) [12] in local training to guide clients toward flatter loss regions, enhancing global performance. This comes at the cost of increased client-side computation, since SAM requires two forward/backward passes for each local optimization step: a gradient ascent step to compute the maximum sharpness and a descent step for sharpness and loss value minimization. Although FEDSAM and its variants [7, 54] demonstrated their effectiveness in various settings, they rely solely on local flatness, assuming that minimizing sharpness

\*Corresponding author: debora.caldarola@polito.it

†Work mainly carried out while at Politecnico di Torino

locally leads to a globally flat minimum<sup>1</sup>. However, in real-world scenarios with significant data heterogeneity, there can be substantial discrepancies between local and global loss landscapes and optimizing for local sharpness does not guarantee the global model will reside in a flat region (Fig. 1). Addressing these limitations, FEDSMOO [53] uses the alternating direction method of multipliers (ADMM) [3] to include global sharpness information in SAM’s local training. While this approach reduces the inconsistency between local and global geometries, it increases communication costs by requiring double the bandwidth in each round. This hinders its real-world applicability, as FL relies on minimizing communication overhead (*i.e.*, message size and exchange frequency) to avoid network congestion and account for potential connection failures.

Given the limitations of existing methods, achieving global flat minima while ensuring communication efficiency in heterogeneous FL remains a critical challenge. To address this, we propose FEDGLOSS (Federated Global Server-side Sharpness) that directly **optimizes global sharpness by using SAM on the server side**, avoiding additional exchanges over the network. Such adaptation is not straightforward, as SAM would require dual exchanges with each client set per round to solve its optimization problem. Instead, FEDGLOSS approximates the sharpness measure using available previous pseudo-gradients. As a result, it facilitates faster training and keeps communication efficiency. To summarize, our core contributions are the following:

- Empirical proof of local-global discrepancies: we provide the first empirical evidence showing the limitations of approaches that focus solely on local sharpness. Our analysis highlights the inconsistency between local and global loss geometries even when using sharpness-aware approaches like FEDSAM, demonstrating that local flatness does not necessarily ensure a flat global minimum. While reaching flat global solutions in simpler problems, we show that their effectiveness diminishes as data complexity and heterogeneity increase (Fig. 1).
- To bridge this gap and motivated by communication efficiency, our FEDGLOSS algorithm directly optimizes for global sharpness on the server using SAM, reducing the communication overhead and the clients’ computational costs compared to previous works. FEDGLOSS consistently achieves flatter minima and outperforms state-of-the-art methods across various vision benchmarks.
- We show the importance of aligning global and local solutions and illustrate how SAM, especially on the server side, enables effective ADMM use in FL. While typically ADMM-based methods suffer from parameter explosion [56], we show that by targeting flat minima, SAM encourages smaller gradient steps and minimal weight updates,

leading to a significantly more stable algorithm.

## 2. Related works

**Federated framework.** In the last few years, Federated Learning (FL) [39] garnered significant attention from both the machine learning and computer vision communities. While the former has primarily focused on optimizing FL algorithms and guaranteeing their convergence [1, 36, 47], the latter has explored its applications in real-world settings, spanning diverse domains like autonomous driving [11, 40, 51] and healthcare [37]. The key appeal of FL lies in its ability to efficiently learn from privacy-protected, distributed data while complying with regulations and leveraging edge resources. Real-world deployments of FL range across both *cross-silo* and *cross-device* settings [24]. This work focuses on the latter, with up to millions of individual devices at the network edge, with typically limited data and computational power, and potential unavailability due to battery life or network connectivity issues. User-specific factors like geographical location, capturing devices and daily habits introduce inherent *bias* and *statistical heterogeneity* into the local datasets. In this setting, FEDGLOSS aims to learn a global model that generalizes to the overall data distribution under statistical heterogeneity without increasing communication complexity, unlike other algorithms for local-global consistency in heterogeneous FL.

**Flatness search in FL.** Recent research has explored the connection between loss landscape geometry and generalization in heterogeneous FL. Studies suggest that convergence to sharp minima might hinder generalization performance [17, 27, 44]. SAM (Sharpness-Aware Minimization) [12] tackles this issue by guiding the optimization toward flatter regions, seeking minima that exhibit both low loss and low sharpness. FEDSAM [4, 45] deploys SAM in local training, marking the first step toward leveraging loss surface geometry in FL to reduce discrepancies between local and global objectives, ultimately improving the global model’s generalization ability. Following its success, FEDSPEED [54] uses perturbed gradients as SAM to reduce local overfitting, FEDGAMMA [7] combines the stochastic variance reduction of SCAFFOLD with SAM and Shi et al. [52] show FEDSAM’s effectiveness in mitigating the negative effects of differential privacy. However, these approaches rely on *local* sharpness information, assuming its minimization directly translates to a globally flat minimum. This may not always be true, as we hypothesize discrepancies may exist between the geometries of local and global losses. Optimizing local sharpness alone does not guarantee a server model residing in a flat region of the *global* loss landscape (Fig. 1). Addressing these limitations, FEDSMOO [53] applies ADMM [3] to the sharpness measure to enforce global and local consistency. This adds communication overhead, doubling the message size in each round and hindering its

<sup>1</sup>We use the term “global flat minima” to refer to local minima within the global (*i.e.*, server-side) loss landscape.

real-world practicality. In contrast, our work focuses on minimizing global sharpness while maintaining communication efficiency. Lastly, building on Stochastic Weight Averaging [22], other works [4, 5] use a window-based average of global models across rounds to reach wider minima. Being agnostic to the underlying optimization algorithm, they remain orthogonal to our approach.

**Heterogeneity in FL.** The de-facto standard algorithm for FL is FEDAVG [39], which updates the global model with a weighted average of the clients' parameters. However, FEDAVG struggles when faced with heterogeneous data distributions, leading to performance degradation and slow convergence due to the local optimization paths diverging from the global one [26]. Reddi et al. [47] shows FEDAVG is equivalent to applying Stochastic Gradient Descent (SGD) [49] with a unitary learning rate on the server side, using the difference between the initial global model parameters and the clients' updates as *pseudo-gradient*, opening the door to alternative optimizers beyond SGD to improve performance and convergence speed. Building on this intuition, this work proposes SAM [12] as a server-side optimizer to enhance generalization by converging toward *global* flat minima. Since SAM requires two optimization steps per iteration, a direct adaptation to the FL setting would double communication exchanges between clients and server; FEDGLOSS overcomes this limitation and maintains communication efficiency through the use of the latest pseudo-gradient as sharpness approximation.

Several approaches address client drift by adding regularization during local training. FedProx [35] introduces a term to keep local parameters close to the global model, FEDDYN [1] employs ADMM to align local and global convergence points, ADABEST [56] adjusts local updates with an adaptive bias estimate, and SCAFFOLD [26] applies stochastic variance reduction. Momentum-based techniques [55] are also employed to maintain a consistent global trajectory, either on the server side (*e.g.*, FEDAVGM [19]) or by incorporating global information into local training [13, 25, 28, 61]. Unlike FEDDYN, where ADMM can lead to parameter explosion [56], our FEDGLOSS successfully leverages ADMM to align global and local solutions, even under extreme heterogeneity, aided by SAM on server.

**Centralized SAM.** To avoid doubling client-server exchanges caused by SAM's two-step process, FEDGLOSS draws on insights from the literature on SAM in centralized settings. Several strategies have been proposed to minimize computational overhead, including reducing the number of parameters needed to compute the sharpness-aware components [9], or approximating them [10, 38, 42]. DP-SAT [42] approximates the ascent step with the gradient from the previous iteration, and SAF [10] replaces SAM's sharpness approximation with the trajectory of weights learned during training. Aiming to the same goal, FEDGLOSS ap-

proximates the sharpness measure with the pseudo-gradient from the previous round on the server side, without incurring in unnecessary exchanges with the clients and effectively guiding the optimization toward globally flat minima.

### 3. Background

This section introduces the FL problem setting and preliminary notations on SAM [12] and FEDSAM [4, 45].

#### 3.1. Problem setting

In FL, a central server communicates with a set of clients  $\mathcal{C}$  for  $T$  rounds. The goal is to learn a global model  $f(\mathbf{w}) : \mathcal{X} \rightarrow \mathcal{Y}$  parametrized by  $\mathbf{w} \in \mathbb{R}^d$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  are the input and the output spaces respectively. In image classification,  $\mathcal{X}$  contains the images and  $\mathcal{Y}$  their corresponding labels. Each client  $k \in \mathcal{C}$  has access to a local dataset  $\mathcal{D}_k$  of  $N_k$  pairs  $\{(x_i, y_i), x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}_{i=1}^{N_k}$ . In realistic heterogeneous settings, clients usually hold different data distributions and quantity, *i.e.*,  $D_i \neq D_j$  and  $N_i \neq N_j \forall i \neq j \in \mathcal{C}$ . The global FL objective is:

$$\min_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{C} \sum_{k \in \mathcal{C}} f_k(\mathbf{w}) \right\}, f_k(\mathbf{w}) \triangleq \mathbb{E} f_k(\mathbf{w}, \xi_k), \quad (1)$$

where  $C \triangleq |\mathcal{C}|$  is the total number of clients,  $f_k$  is the empirical loss on the  $k$ -th client (*e.g.*, cross-entropy loss) and  $\xi_k$  is the data sample randomly drawn from the local data distribution  $D_k$ . The training process is a two-phase optimization approach within each round  $t \in [T]$ . First, due to potential client unavailability, a subset of selected clients  $\mathcal{C}^t \subset \mathcal{C}$  trains the received global model using their local optimizer CLIENTOPT (*e.g.*, SGD, SAM). Then, the server aggregates their updates with a server optimizer, SERVEROPT. FEDOPT [47] solves Eq. (1) as

$$\Delta_{\mathbf{w}}^t \triangleq \sum_{k \in \mathcal{C}^t} \frac{N_k}{N} (\mathbf{w}^t - \mathbf{w}_k^t) \text{ and} \quad (2)$$

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \text{SERVEROPT}(\mathbf{w}^t, \Delta_{\mathbf{w}}^t, \eta_s), \quad (3)$$

where  $\Delta_{\mathbf{w}}^t$  is the **global pseudo-gradient** at round  $t$ ,  $N = \sum_{k \in \mathcal{C}^t} N_k$  the total number of images seen during the current round,  $\eta_s$  the server learning rate,  $\mathbf{w}^t$  the global model and  $\mathbf{w}_k^t$  the local update resulting from training on client  $k$ 's data with CLIENTOPT for  $E$  epochs. FEDAVG [39] computes  $\mathbf{w}^{t+1}$  as  $\sum_{k \in \mathcal{C}^t} N_k/N \mathbf{w}_k^t$ , corresponding to one SGD step on the pseudo-gradient  $\Delta_{\mathbf{w}}^t$  with  $\eta_s = 1$  [47].

#### 3.2. Sharpness-aware Minimization

SAM [12] jointly minimizes the loss value and the sharpness of the loss landscape by solving the min-max problem

$$\min_{\mathbf{w}} \left\{ F(\mathbf{w}) \triangleq \max_{\|\boldsymbol{\epsilon}\| \leq \rho} f(\mathbf{w} + \boldsymbol{\epsilon}) \right\}, \quad (4)$$

where  $\boldsymbol{\epsilon}$  is the perturbation to estimate the sharpness,  $f$  the loss function,  $\rho$  the neighborhood size and  $\|\cdot\|$  the  $\ell_2$  norm.

Using the first-order Taylor expansion of  $f$ , SAM efficiently solves the inner maximization as

$$\operatorname{argmax}_{\|\epsilon\| \leq \rho} f(\mathbf{w}) + \epsilon^\top \nabla_{\mathbf{w}} f(\mathbf{w}) = \rho \frac{\nabla_{\mathbf{w}} f(\mathbf{w})}{\|\nabla_{\mathbf{w}} f(\mathbf{w})\|} \triangleq \hat{\epsilon}(\mathbf{w}). \quad (5)$$

$\hat{\epsilon}$  is the scaled gradient of the loss w.r.t. the current parameters  $\mathbf{w}$ . The *sharpness-aware gradient* is  $\nabla_{\mathbf{w}} f(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$ . Eq. (4) is solved with a first gradient ascent step to compute  $\hat{\epsilon}$  and a descent step with the sharpness-aware gradient, updating the model as  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} f(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$ .

### 3.3. SAM in Federated Learning

FEDSAM [4, 45] aims to improve the clients' models generalization through convergence to flatter regions by using SAM in the local training. From Eqs. (1) and (4), the global objective becomes  $\min_{\mathbf{w}} \{f^{\text{SAM}}(\mathbf{w}) = 1/C \sum_{k \in \mathcal{C}} f_k^{\text{SAM}}(\mathbf{w})\}$ , with  $f_k^{\text{SAM}}(\mathbf{w}) \triangleq \max_{\|\epsilon_k\| \leq \rho} f_k(\mathbf{w} + \epsilon_k)$  with local perturbation  $\epsilon_k$ . The intuition behind this approach is that the improved local models' generalization positively reflects on the global model performance. However, by independently applying Eq. (4) in the local optimization, FEDSAM does not explicitly address global flatness, potentially leading to discrepancies between local and global loss geometries.

## 4. Local-Global Sharpness Inconsistency

This section empirically investigates the hypothesis that discrepancies between local and global loss landscapes impact FEDSAM's performance, using a CNN model on CIFAR10 and CIFAR100 datasets — further details in Sec. 6.

Fig. 1 compares the loss surfaces of CNNs trained with FEDAVG and FEDSAM. On the easier CIFAR10, FEDSAM exhibits flatter minima w.r.t. FEDAVG, effectively navigating simpler landscapes. However, their difference diminishes with increasing dataset complexity (CIFAR100) and heterogeneity ( $\alpha \rightarrow 0$ ). This suggests larger discrepancies between local and global geometries arise as tasks become more complex and data distributions more diverse.

To highlight the existing difference between local and global behavior, Fig. 2 investigates the behavior of client models at the end of local training when tested on their own data  $\mathcal{D}_k$  (*bottom landscape*), prior to server-side aggregation, w.r.t. the overall dataset  $\mathcal{D}$  (*top landscape*). Each plot shows the behavior of one of the randomly selected clients during the last round with FEDSAM, distinguished by the locally seen class (results for all 5 clients in App. B). The inconsistency between local and global behavior can be easily appreciated: locally, each model lands in a flat region; differently, the same model is close to saddle points (Fig. 2a) or sharp minima (Figs. 2b and 2c) in the global landscape. These findings are further corroborated by the maximum Hessian eigenvalues presented in Tab. 1, computed using each client's local dataset ( $\lambda_{1,l}$ ) and the overall one ( $\lambda_{1,g}$ ). FEDSAM performs well on simple datasets like CIFAR10

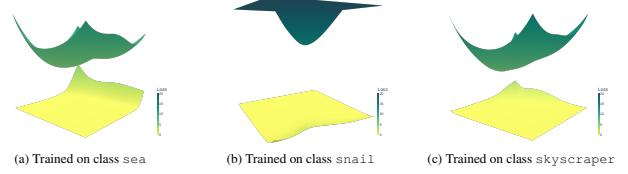


Figure 2. **Global vs. local perspective on FEDSAM.** CIFAR100  $\alpha = 0$  @ 20k rounds on CNN. Local models trained on one class, tested on the local (*bottom landscape*) or global dataset (*top landscape*). Models trained with FEDSAM present significant differences between local and global behaviors.

Table 1. **Maximum Hessian eigenvalues of local models**, computed on global ( $\lambda_{1,g}$ ) and local datasets ( $\lambda_{1,l}$ ). CIFAR10 and CIFAR100,  $\alpha = 0$ . Each client is identified via its local class. The lowest  $\lambda_{1,g}$  in **bold**. FEDDYN does not converge on CIFAR100 with  $\alpha = 0$  [4, 56], hence the lack of results (**X**).

Local Class	FEDAVG		FEDSAM		FEDDYN		FEDDYN + SAM		FEDSMO		FEDGLOSS		
	$\lambda_{1,l}$	$\lambda_{1,g}$											
CIFAR10	airplane	9.1	239.1	100.6	36.4	752.5	347.8	199.6	12.0	122.1	26.5	190.1	<b>4.3</b>
	cat	424.2	273.6	28.8	16.5	59.9	242.3	122.0	11.1	82.4	26.9	106.9	<b>3.9</b>
	bird	18.4	237.0	106.4	35.7	894.0	371.2	200.2	12.0	134.2	25.7	200.1	<b>4.1</b>
	airplane	483.5	269.5	103.2	30.6	761.6	348.9	206.9	12.3	122.8	25.2	207.8	<b>4.0</b>
	frog	263.2	259.0	68.1	32.9	528.9	286.0	155.6	11.7	79.3	33.5	84.8	<b>4.1</b>
	sea	251.0	224.5	0.1	238.5				33.2	31.4	28.3	<b>19.6</b>	
CIFAR100	snail	91.2	267.0	0.2	149.1				331.2	102.2	260.8	<b>40.7</b>	
	bear	108.4	215.2	6.7	129.3	<b>X</b>	<b>X</b>	<b>X</b>	428.6	121.0	220.3	<b>49.6</b>	
	skyscraper	613.3	300.1	1.3	194.6				143.5	40.2	269.2	<b>22.2</b>	
	possum	37.9	259.6	15.3	142.6				455.5	90.9	392.4	<b>39.0</b>	

with global eigenvalues smaller than local ones. However,  $\lambda_{1,l} \ll \lambda_{1,g}$  on the more complex CIFAR100. This suggests that FEDSAM effectively achieves *local* convergence to flatter regions of the loss landscape on individual devices, but the higher global eigenvalue indicates limitations in reaching a *globally* flat minimum. The challenge of achieving flat regions under high heterogeneity and the gap between local and global flatness support the introduction of FEDGLOSS.

## 5. FL with Global Server-side Sharpness

FEDGLOSS (Federated Global Server-side Sharpness) overcomes FEDSAM's limitations by efficiently optimizing both global flatness and consistency.

### 5.1. Rethinking SAM in Federated Learning

Aiming to optimize SAM's objective (Eq. (4)) on the global function, FEDGLOSS solves  $\min_{\mathbf{w}} \{\mathcal{F}(\mathbf{w}) = \frac{1}{C} \sum_{k \in \mathcal{C}} \mathcal{F}_k(\mathbf{w})\}$ , with  $\mathcal{F}_k(\mathbf{w}) \triangleq \max_{\|\epsilon\| \leq \rho} f_k(\mathbf{w} + \epsilon)$ , where  $\epsilon$  is the global perturbation. Calculating the true  $\epsilon$  value requires the global gradient  $\nabla_{\mathbf{w}} f$  (Eq. (5)) computed on the entire dataset  $\mathcal{D} \triangleq \cup_{k \in \mathcal{C}} \mathcal{D}_k$ , which is not available in FL due to data privacy and communication constraints. While FEDSMO [53] tackles this issue by using ADMM on the sharpness with the constraint  $\epsilon = \epsilon_k$ , it necessitates transmitting  $\epsilon$  alongside the model parameters  $\mathbf{w}$  to both clients and server in each round, hindering its practicality in real-world scenarios with limited communication budgets. This observation motivates the question: *how to minimize global sharpness while maintaining communication efficiency?*

#### 5.1.1. Challenges of Server-side SAM

We address this question by applying SAM on the server side, directly optimizing for global sharpness and eliminating the need to align local sharpness on the clients. The global model has to be updated as  $\mathbf{w}^{t+1} \leftarrow$

Table 2. Overview of FL methods using SAM. Differently from previous works, FEDGLOSS uses SAM as server optimizer and allows any local optimizer.

Method	SERVOPT	CLIENTOPT	Global Flatness	Communication Cost	Local Computation Cost
FEDSAM [4, 45]	SGD	SAM	x	1x	2x
FEDDYN [1] + SAM	SGD	SAM	x	1x	2x
FEDSPEED [54]	SGD	Similar to SAM	x	1x	2x
FEDGAMMA [7]	SGD	SAM	✓	2x	2x
FEDSMO [53]	SGD	SAM	✓	2x	2x
<b>FEDGLOSS</b>	SAM	Any optimizer	✓	1x	1x or 2x

$\mathbf{w}^t - \eta_s \nabla_{\mathbf{w}} \mathcal{F}(\mathbf{w})|_{\mathbf{w}^t + \hat{\epsilon}^t(\mathbf{w})}$ , where  $\hat{\epsilon}^t$  is the global perturbation at each round  $t$ . However, a key challenge arises: the computation of both  $\hat{\epsilon}^t$  and the sharpness-aware gradient necessitates two transmissions with the clients, making its direct application in server-side FL non-trivial. A straightforward solution is to emulate SAM’s double computation step through two communication exchanges  $\forall t \in [T]$ .

- **Step 1:** the server sends the global model  $\mathbf{w}^t$  to a subset  $\mathcal{C}^t$  of clients, which update it using their local data. With the resulting pseudo-gradient  $\Delta_{\mathbf{w}}^t$ ,  $\hat{\epsilon}^t(\mathbf{w}) = \rho(\Delta_{\mathbf{w}}^t / \|\Delta_{\mathbf{w}}^t\|)$  and the perturbed model  $\tilde{\mathbf{w}}^t = \mathbf{w}^t + \hat{\epsilon}^t(\mathbf{w})$ .
- **Step 2:** the server transmits  $\tilde{\mathbf{w}}^t$  to the *same*  $\mathcal{C}^t$ , which compute their update  $\tilde{\mathbf{w}}_k^t \forall k$ . The resulting global pseudo-gradient  $\tilde{\Delta}_{\mathbf{w}}^t \triangleq \sum_{k \in \mathcal{C}^t} N_k / N (\tilde{\mathbf{w}}^t - \tilde{\mathbf{w}}_k^t)$  is an estimate of  $\nabla_{\mathbf{w}} \mathcal{F}(\mathbf{w})|_{\mathbf{w}^t + \hat{\epsilon}^t(\mathbf{w})}$ .

This two-step approach, referred to as NAIVEFEDGLOSS, is conceptually simple but suffers from communication inefficiency, doubling the communication cost w.r.t. FedAvg, while requiring the same set of clients  $\mathcal{C}^t$  to remain active for two consecutive exchanges. This may be unrealistic in real-world settings often characterized by network failures. These limitations highlight the need for an efficient alternative that accounts for practical real-world FL factors.

## 5.2. FedGloss

To overcome the challenges posed by NAIVEFEDGLOSS, following [42], FEDGLOSS estimates  $\hat{\epsilon}^t$  using the perturbed global pseudo-gradient from the previous round  $\tilde{\Delta}_{\mathbf{w}}^{t-1}$  at each round  $t$ . This approach leverages available information *without incurring extra communications* and avoids unnecessary computations. Intuitively, the use of the previous pseudo-gradient to minimize the sharpness allows FEDGLOSS to access information on the *global* loss landscape geometry, thus guiding the *global* optimization towards flatter minima. From Eqs. (3) and (5), FEDGLOSS updates the global model  $\mathbf{w}^t$  as

$$\textcircled{1} \quad \tilde{\epsilon}^t(\mathbf{w}) \triangleq \rho \frac{\tilde{\Delta}_{\mathbf{w}}^{t-1}}{\|\tilde{\Delta}_{\mathbf{w}}^{t-1}\|} \quad \textcircled{2} \quad \tilde{\mathbf{w}}^t \leftarrow \mathbf{w}^t + \tilde{\epsilon}^t(\mathbf{w})$$

$$\textcircled{3} \quad \text{Obtain } \tilde{\mathbf{w}}_k^t \text{ from clients and } \tilde{\Delta}_{\mathbf{w}}^t = \sum_{k \in \mathcal{C}^t} \frac{N_k}{N} (\tilde{\mathbf{w}}^t - \tilde{\mathbf{w}}_k^t)$$

$$\textcircled{4} \quad \mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \text{FEDGLOSS}(\mathbf{w}^t, \tilde{\Delta}_{\mathbf{w}}^t, \eta_s) = \mathbf{w}^t - \eta_s \tilde{\Delta}_{\mathbf{w}}^t,$$

where with a slight abuse of notation SERVEROPT from Eq. (3) is substituted with the server-side strategy proposed by FEDGLOSS. The notation follows the colors of Fig. 3, which depicts our approach. Notably, as summarized in Tab. 2, FEDGLOSS enables SAM on the server side while

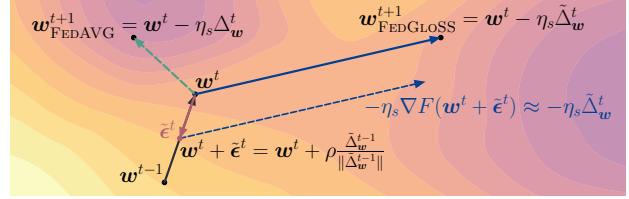


Figure 3. Illustration of FEDGLOSS. The model  $\mathbf{w}^t$  is perturbed using  $\tilde{\Delta}_{\mathbf{w}}^{t-1}$ . The sharpness-aware direction (dashed) is used to compute  $\mathbf{w}^{t+1}$  (solid), which lands in a flat region. Compared to FEDAVG.

allowing any CLIENTOPT for local training, with computational costs varying based on the chosen optimizer. This differs from previous methods constrained to the more computationally expensive SAM. In addition, differently from FEDSMO, FEDGLOSS maintains FEDAVG’s communication complexity while optimizing for global flatness.

### 5.2.1. Promoting Global Consistency with ADMM

The difference in using the approximation  $\tilde{\epsilon}^t$  (FEDGLOSS) and the true  $\hat{\epsilon}^t$  (NAIVEFEDGLOSS) is

$$\delta_{\epsilon}^t \triangleq \|\epsilon^t(\mathbf{w}) - \hat{\epsilon}^t(\mathbf{w})\| = \rho \left\| \frac{\tilde{\Delta}_{\mathbf{w}}^{t-1}}{\|\tilde{\Delta}_{\mathbf{w}}^{t-1}\|} - \frac{\Delta_{\mathbf{w}}^t}{\|\Delta_{\mathbf{w}}^t\|} \right\|, \quad (6)$$

where  $\tilde{\Delta}_{\mathbf{w}}^{t-1}$  is computed using the updates of the clients in  $\mathcal{C}^{t-1}$  and  $\tilde{\Delta}_{\mathbf{w}}^t$  with  $\mathcal{C}^t$ . Eq. (6) suggests  $\delta_{\epsilon}^t$  is minimized when  $\tilde{\Delta}_{\mathbf{w}}^{t-1}$  and  $\tilde{\Delta}_{\mathbf{w}}^t$  are aligned, which occurs when clients’ updates are directionally consistent. However, in real-world heterogeneous FL, *i*) due to clients’ unavailability, only a subset of them participates in training at each round, with  $\mathcal{C}^t$  likely differing from  $\mathcal{C}^{t-1}$ , and *ii*) clients hold different data distributions, *i.e.*, local optimization paths likely converge towards different local minima, leading to unstable global updates [26]. As a consequence,  $\delta_{\epsilon}^t \not\rightarrow 0$  necessarily.

To align local and global objectives - guiding client and server updates in the same direction and minimizing Eq. (6) - FEDGLOSS leverages the Alternating Direction Method of Multipliers (ADMM) [3] on  $\mathbf{w}^t$  [1, 53, 54]. While alternative approaches could be used, they either lack full immunity to data heterogeneity or have shown poor performance on realistic scenarios (*e.g.*, variance reduction [7, 26]). In contrast, ADMM has been proved to converge under arbitrary heterogeneity [1] and can thus be leveraged as a base algorithm for FEDGLOSS, as shown in Alg. 1 in App. A. ADMM makes use of the augmented Lagrangian function  $\mathcal{L}(\mathbf{w}, \mathbf{W}, \sigma) = \sum_{k \in \mathcal{C}} L(\mathbf{w}, \mathbf{w}_k, \sigma_k)$  where  $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_C\}$  and  $\sigma$  is the Lagrangian multiplier. The problem solved by  $\mathcal{L}$  is

$$\frac{1}{C} \sum_{k \in \mathcal{C}} (f_k + \sigma_k^\top (\mathbf{w}^t - \mathbf{w}_k^t) + \frac{1}{2\beta} \|\mathbf{w}^t - \mathbf{w}_k^t\|^2) \text{ s.t. } \mathbf{w} = \mathbf{w}_k \quad (7)$$

with  $\beta > 0$  being an hyperparameter. Eq. (7) is split into  $C$  sub-problems of the form  $\mathbf{w}_{k,E} = \operatorname{argmin}_{\mathbf{w}_k} \{f_k - \sigma_k^\top (\mathbf{w}^t - \mathbf{w}_k) + \frac{1}{2\beta} \|\mathbf{w}^t - \mathbf{w}_k\|^2\}$ . The local dual variable is updated as  $\sigma_k \leftarrow \sigma_k - \frac{1}{\beta} (\mathbf{w}_{k,E}^t - \mathbf{w}_{k,0}^t)$ . The global one  $\sigma$  is updated by adding the averaged  $\mathbf{w}_k - \mathbf{w}^t \forall k \in \mathcal{C}$ . We note that ADMM introduces a key limitation of our approach, namely reliance on stateful clients. Fig. 5 confirms

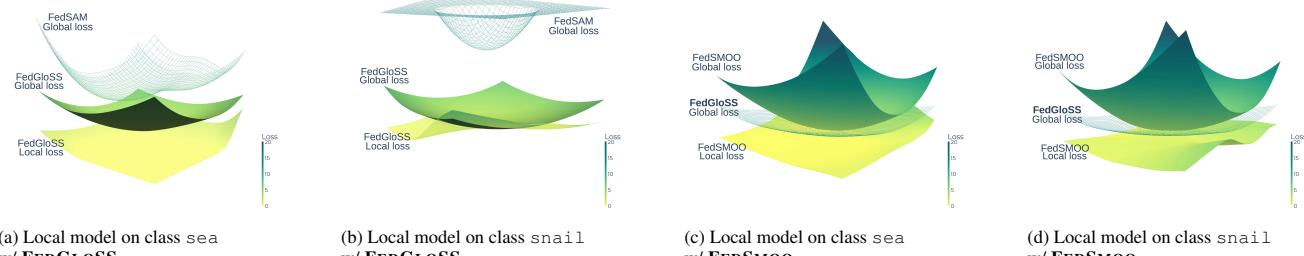


Figure 4. **Global vs. local perspective of FEDGLOSS and FEDSMOO.** Loss landscapes of clients models trained on one class, tested on the local (“*Local loss*”) or global dataset (“*Global loss*”). CIFAR100  $\alpha = 0$  with SAM as local optimizer @  $t = 20k$ , CNN. (a)-(b): Models trained with FEDGLOSS. Global loss of FEDSAM’s local model (*net*) as reference. (c)-(d): Models trained with FEDSMOO. Global loss of FEDGLOSS’s local model (*net*) as reference. **FEDGLOSS achieves better consistency w.r.t. FEDSMOO.**

the directional consistency of local and global updates enforced by ADMM reduces  $\delta_\epsilon^t$  (Eq. (6)), *i.e.*, the difference between the true and approximated perturbation. The gap between the directions of  $\tilde{\Delta}_{\mathbf{w}}^{t-1}$  and  $\tilde{\Delta}_{\mathbf{w}}^t$  remains constant after an initial phase, suggesting the most challenging loss landscape direction is largely stable over time.

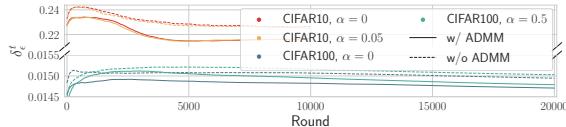


Figure 5. Trend of the difference  $\delta_\epsilon^t$  (Eq. (6)), which decreases as ADMM is used and over training rounds. CIFAR datasets, CNN.

## 6. Experiments

### 6.1. Experimental Setting

App. C details implementation and hyperparameter settings.

**Federated datasets.** We leverage established FL benchmarks [6, 19, 20]. *Small-scale image classification:* following [4, 20], the federated versions of CIFAR10 (10 classes) and CIFAR100 (100 classes) [29] split the respective 50k training images in 100 clients with 500 images each. The data distribution is controlled by Dirichlet’s parameter  $\alpha \in \{0, 0.05, 1, 5, 10\}$  for CIFAR10 and  $\{0, 0.5\}$  for CIFAR100 [19]. Lower  $\alpha$  signifies increased heterogeneity, with  $\alpha = 0$  being the most challenging scenario (each client holds samples from one class). *Large-scale image classification:* LANDMARKS-USER-160K (2,028 classes) [20] is the federated Google Landmarks v2 [58] with 164,172 pictures of worldwide locations, split among 1,262 realistic clients.

**Models.** The effectiveness of FEDGLOSS is shown using multiple model architectures. As in [4, 20], we use a Convolutional Neural Network (CNN) similar to LeNet5 [30] on both CIFAR10 ( $T = 10k$ ) and CIFAR100 ( $T = 20k$ ). Experiments with ResNet18 [16] run for 10k rounds. For LANDMARKS-USER-160K, we train MobileNetv2 [20, 50] ( $T = 1.3k$ ), considering the limited resources at the edge.

**Baselines.** To study real-world settings with varying participation, a small fraction of clients is sampled at each round, with participation rate set to 5% with the CNN and 10% with ResNet18 on both CIFARS, and to 50 clients per round in LANDMARKS-USER-160K ( $\approx 4\%$ ). FEDGLOSS is compatible with any local optimizer (Sec. 5). We choose SGD and SAM to comply with previous works and compare it with state-of-the-art (SOTA) methods for statistical

heterogeneity in FL, distinguishing the results by optimizer type to highlight performance differences. SGD-based approaches are FedAvg [39], FedProx [35], FedDyn [1] and Scaffold [26], while use SAM FedSAM [4, 45], FedDyn + SAM, FedSpeed [54], FedGamma [7] and FedSmoo [53].

### 6.2. Achieving Local-Global Sharpness Consistency

To assess the effectiveness of FEDGLOSS in promoting consistency between local and global loss landscapes, Fig. 4 replicates the analysis previously conducted on FEDSAM (Fig. 2). The behavior of local models is shown from both local (“*Local loss*”) and global perspectives (“*Global loss*”). App. B.1 offers visualizations for the remaining clients. Compared to FEDSAM, the gap between local and global loss landscapes in FEDGLOSS is significantly smaller, and both global and local loss surfaces are found in *flat and low-loss regions* (Figs. 4a and 4b). This suggests our method effectively promotes convergence toward aligned low-loss flat regions, minimizing the discrepancy between local and global geometries. This results in a global model residing in a flat minimum in the global landscape (Figs. 6a and 6c). Figs. 4c and 4d instead compare FEDGLOSS with the best-performing SOTA FEDSMOO, where the position in the global landscape of FEDGLOSS’ local models is added for reference. While FEDSMOO improves consistency between local and global sharpness compared to FEDSAM, it falls short of FEDGLOSS in reaching a flatter global minimum.

Tab. 1 confirms these claims. By combining ADMM for consistency and server-side SAM for global flatness, FEDGLOSS prioritizes achieving a flatter *global* region during training, as proven by the lowest global maximum eigenvalue  $\lambda_{1,g}$  and larger  $\lambda_{1,l}$ , across all clients and methods.

### 6.3. Benchmarking FedGloSS against SOTA

This section compares FEDGLOSS with SOTA methods on vision tasks in heterogeneous FL. App. B.4 presents results for  $\alpha \in \{1, 5, 10\}$  on CIFAR10 and homogeneous settings.

#### 6.3.1. FedGloSS on Standard Federated Benchmarks

Tab. 3 presents results on CIFAR100 and CIFAR10 with varying levels of heterogeneity on the CNN model. Several observations highlight the advantages of FEDGLOSS. FEDGLOSS achieves the best results among both SGD and SAM-based approaches while maintaining communication

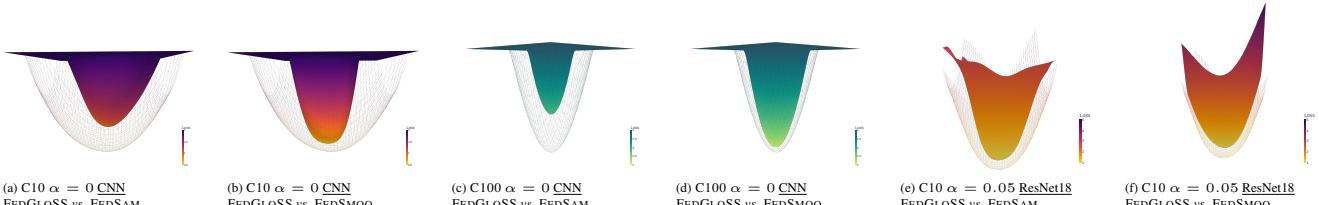
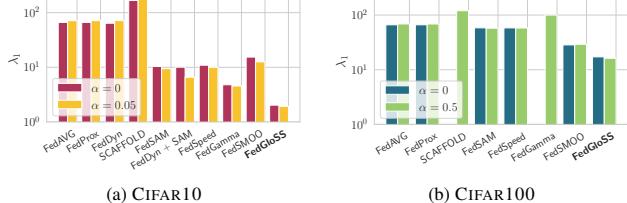
(a) C10  $\alpha = 0$  CNN(b) C10  $\alpha = 0$  CNN(c) C100  $\alpha = 0$  CNN(d) C100  $\alpha = 0$  CNN(e) C10  $\alpha = 0.05$  ResNet18(f) C10  $\alpha = 0.05$  ResNet18

Figure 6. Loss landscapes of models trained with FEDGLOSS (net) vs. FEDSAM and FEDSMOO (solid) on CIFAR10/100. (a) - (c) - (e): The flatter regions reached by FEDGLOSS w.r.t. FEDSAM prove the effectiveness of optimizing for global flatness. (b) - (d) - (f): FEDGLOSS achieves flatter minima and lower loss values w.r.t. FEDSMOO.



(a) CIFAR10

(b) CIFAR100

Figure 7. Maximum Hessian eigenvalues ( $\lambda_1$ ), CNN. Values shown only if algorithm converged. FEDGLOSS reaches the flattest global minima.

Table 3. FEDGLOSS vs. the state of the art on CIFAR datasets, distinguished by local optimizer, SGD (top) and SAM (bottom), in terms of communication cost and accuracy (%). Best results in bold. Model: CNN.

Method	Comm. Cost	CIFAR10		CIFAR100	
		$\alpha = 0$	$\alpha = 0.05$	$\alpha = 0$	$\alpha = 0.05$
<b>Client SGD</b>					
FEDAVG	1 X	59.9 ± 0.4	65.7 ± 1.0	28.6 ± 0.7	38.5 ± 0.5
FEDPROX	1 X	59.8 ± 0.5	65.6 ± 1.0	28.8 ± 0.7	38.7 ± 0.4
FEDDYN	1 X	65.5 ± 0.3	70.1 ± 1.2	X	X
SCAFFOLD	2 X	25.1 ± 3.7	54.0 ± 2.6	X	30.0 ± 1.1
<b>FEDGLOSS</b>	1 X	<b>69.5 ± 0.4</b>	<b>75.5 ± 0.3</b>	<b>42.5 ± 0.6</b>	<b>47.9 ± 0.5</b>
<b>Client SAM</b>					
FEDSAM	1 X	70.2 ± 0.9	71.5 ± 1.08	28.7 ± 0.5	39.6 ± 0.5
FEDDYN	1 X	79.3 ± 3.1	81.5 ± 0.6	X	X
FEDSPEED	1 X	70.9 ± 0.4	72.3 ± 1.1	28.9 ± 0.5	39.7 ± 0.5
FEDGAMMA	2 X	58.9 ± 1.8	61.9 ± 1.8	X	29.4 ± 1.4
FEDSMOO	2 X	81.3 ± 0.5	82.8 ± 0.6	47.8 ± 0.5	51.7 ± 0.46
<b>FEDGLOSS</b>	1 X	83.9 ± 0.4	84.4 ± 0.5	50.8 ± 0.8	53.4 ± 0.5

efficiency. FEDGLOSS with local SAM consistently outperforms the best-performing SOTA FEDSMOO by  $\approx 2.5$  percentage points in accuracy across all configurations *with half the communication cost*. Our method reaches the flattest global minima (e.g.,  $\lambda_1^{\text{FEDGLOSS}} = 2.03$  vs.  $\lambda_1^{\text{FEDSMOO}} = 15.37$  on CIFAR10  $\alpha = 0$ ), as shown in Figs. 6 and 7, achieving the best overall performance. FEDGLOSS with local SGD overcomes by  $\approx 4$  percentage points *all* SGD-based approaches. FEDDYN suffers from parameter explosion in highly heterogeneous settings [56], failing to converge on CIFAR100. Differently, FEDGLOSS successfully employs ADMM to align global and local solutions, with the best results under extreme heterogeneity. Studies showed SCAFFOLD performs poorly in complex heterogeneous environments [4, 33], resulting in its inability to converge on CIFAR100 alongside FEDGAMMA. Tab. 4 confirms FEDGLOSS’s effectiveness, consistently outperforming SOTA methods with the more complex ResNet18 architecture, with  $\approx 8$  points higher accuracy w.r.t. FEDAVG with both SGD and SAM, and +5 w.r.t. FEDSMOO, with the flattest solutions (Figs. 6e and 6f).

### 6.3.2. FedGloss on Real-World Large-Scale Datasets

To further highlight FEDGLOSS’s effectiveness, we evaluate it on *large-scale image classification* using the challenging LANDMARKS-USER-160K dataset. Tab. 5 compares FEDGLOSS with local SAM against the best-performing

Table 4. ResNet18 on CIFAR100  
 $\alpha = 0.5$  and CIFAR10  $\alpha = 0.05$ .

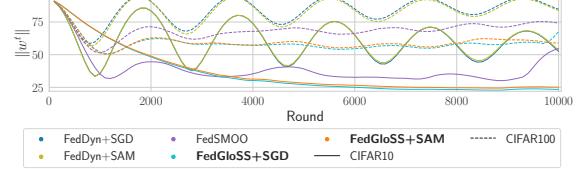
Table 5. MobileNetv2 on LANDMARKS-USER-160K.

Method	Comm. cost	C100 Acc.	C10 Acc.
FEDAVG	1 X	37.4 ± 0.2	72.6 ± 0.1
FEDPROX	1 X	37.6 ± 0.1	72.2 ± 0.2
FEDDYN	1 X	38.8 ± 0.6	70.2 ± 0.6
SCAFFOLD	2 X	38.6 ± 0.1	70.8 ± 0.6
<b>FEDGLOSS</b>	1 X	<b>46.7 ± 0.6</b>	<b>79.1 ± 0.5</b>
FEDSAM	1 X	38.5 ± 0.1	72.8 ± 0.1
FEDDYN	1 X	39.6 ± 0.8	72.6 ± 0.2
FEDSPEED	1 X	38.7 ± 0.6	72.6 ± 0.1
FEDGAMMA	2 X	38.6 ± 0.3	72.2 ± 0.1
FEDSMOO	2 X	44.8 ± 0.5	75.3 ± 0.6
<b>FEDGLOSS</b>	1 X	<b>47.2 ± 0.2</b>	<b>80.0 ± 0.3</b>

baselines. FEDGLOSS is among the few methods, alongside FEDSAM and FEDSMOO, outperforming FEDAVG. Similarly to the CIFAR100 results (Sec. 6.3.1), both SCAFFOLD and FEDGAMMA fail to converge. Importantly, FEDGLOSS achieves the best overall performance (+3.4% w.r.t. FEDAVG) with reduced communication overhead.

### 6.3.3. ADMM and SAM Interaction in FedGloss

ADMM-based methods are often prone to parameter explosion in highly heterogeneous FL settings with many clients [56]. This occurs as multiple gradients accumulate in the global dual variable  $\sigma$  (Sec. 5), causing the parameter norms to grow uncontrollably. However, empirical results indicate this issue is mitigated with SAM (e.g., see FEDDYN vs. FEDGLOSS in Tab. 3). We attribute this to SAM’s nature: by targeting flat minima, it promotes smaller gradient steps and minimal weight updates, resulting in a more stable algorithm. Fig. 8 confirms our hypothesis by showing SAM’s stability effectively lowers parameter norms and the consequent risk of explosion, particularly when SAM is applied directly to the global model, as in FEDGLOSS.

Figure 8. Trend of model parameters norm,  $\|w^t\|_2$ , on SAM-based methods with ResNet18 on CIFAR datasets. SAM reduces the norm and the risk of parameters explosion, successfully enabling ADMM in heterogeneous FL.

### 6.3.4. Communication Efficiency with FedGloss

Communication cost is the main bottleneck in FL [32], making its optimization a relevant challenge. As already previously highlighted, FEDGLOSS considers communication efficiency its primacy concern. Defined  $B$  the number of bits exchanged by FEDAVG in  $T$  training rounds, Tab. 6 studies FEDGLOSS’s communication cost against the SOTA baselines in terms of rounds necessary to reach FEDAVG’s performance and quantity of exchanged bits.

Table 6. **Communication costs** comparison w.r.t. FEDAVG. “-” for not reached accuracy, “ $\times$ ” for non-convergence. GLDV2 is LANDMARKS-USER-160K.

Method	CNN				ResNet18				MobileNetv2			
	Rounds	Cost	Rounds	Cost	Rounds	Cost	Rounds	Cost	Rounds	Cost	Rounds	Cost
FEDAVG	10k	1B	20k	1B	10k	1B	10k	1B	1.3k	1B	-	-
FEDPROX	7.6k	0.8B	18.7k	0.9B	8.8k	0.9B	8.3k	0.8B	-	-	-	-
FEDDYN	2k	0.2B	-	-	-	-	3.5k	0.4B	-	-	-	-
SCAFFOLD	-	-	-	-	-	-	8.9k	1.8B	-	-	-	-
<b>FEDGLOSS</b>	3.4k	0.3B	5k	0.3B	2.4k	0.2B	1.9k	0.2B	1.3k	1B	-	-
FEDSAM	6.3k	0.6B	18.3k	0.9B	9.2k	0.9B	7.8k	0.8B	1.3k	1B	-	-
FEDDYN	3k	0.3B	-	-	4.1k	0.4B	3.5k	0.4B	-	-	-	-
FEDSPEED	6.3k	0.6B	18.3k	0.9B	8.3k	0.8B	8.3k	0.8B	1.3k	1B	-	-
FEDGAMMA	-	-	-	-	9.3k	1.9B	8.1k	1.6B	-	-	-	-
FEDSMOO	1.9k	0.4B	4.5k	0.5B	2.4k	0.5B	2.3k	0.5B	200	0.4B	-	-
<b>FEDGLOSS</b>	2.2k	0.2B	6.3k	0.3B	2.4k	0.2B	1.9k	0.2B	200	0.2B	-	-

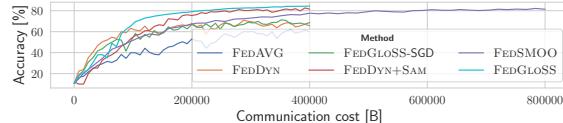


Figure 9. Accuracy gain vs. communication cost. CNN, CIFAR10  $\alpha = 0$ .

The ADMM-based methods are usually faster, with FEDGLOSS being the fastest with ResNet18 and MobileNetv2. While FEDSMOO is faster when using the CNN model, the transmitted bits double due to its increased communication cost, making FEDGLOSS the most efficient method in all cases. Analyses on left-out settings in App. B.5. Fig. 9 further shows the *total* communication cost of the FL algorithms (#model exchanges per round  $\times |\mathcal{C}^t| \times T$ ). The doubled cost of FEDSMOO is clear, while FEDGLOSS achieves nearly double FEDAVG’s accuracy at the same cost.

## 6.4. Ablation Studies

### 6.4.1. Communication-efficient Sharpness

This section evaluates using the pseudo-gradient from the previous round  $\tilde{\Delta}_{\mathbf{w}}^{t-1}$  (Eq. (6)) to estimate sharpness. FEDGLOSS leverages past gradients as a reliable indication on the global loss landscape and aligns global and local paths via ADMM for consistent trajectories across rounds.

Tab. 7 compares FEDGLOSS with its baseline, NAIVEFEDGLOSS (Sec. 5), which computes the exact perturbation  $\hat{\epsilon}^t$  at the expense of doubled communication costs. FEDGLOSS matches NAIVEFEDGLOSS in accuracy while maintaining communication efficiency, with minimal or negligible gap in performance, consistent with similar sharpness ( $\lambda_1$ ). To fairly assess communication costs, Tab. 7 also reports FEDGLOSS’ final accuracy vs. NAIVEFEDGLOSS’ performance at 50% training progress, showing FEDGLOSS’ higher accuracy given the same number of exchanges. Overall, given  $B$  bits transmitted by FEDGLOSS, the cost of NAIVEFEDGLOSS is  $1.96B$  on C10  $\alpha = 0$ ,  $2B$  for  $\alpha = 0.05$ ,  $1.94B$  on C100  $\alpha = 0$  and  $1.97B$  for  $\alpha = 0.5$ . Thus, FEDGLOSS reduces communication costs by up to 50%, with no trade-off in model performance. If NAIVEFEDGLOSS were preferred to maintain exact gradients, it would still achieve the *best results* with costs *lower or comparable* to FEDSMOO. In addition, our strategy in centralized settings lowers the performance w.r.t. NAIVEFEDGLOSS. At equal performance, FEDGLOSS narrows the gap to the upper bound:  $-2.4\%$  on CIFAR10 with  $\alpha = 0$

Table 7. **FEDGLOSS vs. NAIVEFEDGLOSS** in terms of communication cost, accuracy (50% and 100% of training) and maximum Hessian eigenvalue  $\lambda_1$ . SAM as CLIENTOPT. CIFAR datasets with  $\alpha = 0$  (top) and  $\alpha = 0.05/0.5$  (bottom).  $\widetilde{\text{SAM}}$  is SAM with the sharpness approximation of FEDGLOSS, using the previous gradient.

Method	Comm. Cost	CIFAR10			CIFAR100		
		Acc@50%	Acc@100%	$\lambda_1 (\downarrow)$	Acc@50%	Acc@100%	$\lambda_1 (\downarrow)$
NAIVEFEDGLOSS	2x	77.6 $\pm$ 0.3	83.9 $\pm$ 0.2	2.78 $\pm$ 0.13	42.6 $\pm$ 0.8	50.8 $\pm$ 0.1	16.93 $\pm$ 0.27
FEDGLOSS	1x	78.9 $\pm$ 0.5	83.9 $\pm$ 0.4	2.03 $\pm$ 0.05	39.5 $\pm$ 0.9	50.6 $\pm$ 0.6	17.18 $\pm$ 0.97
NAIVEFEDGLOSS	2x	78.7 $\pm$ 0.1	84.4 $\pm$ 0.2	2.75 $\pm$ 0.09	49.4 $\pm$ 0.6	53.7 $\pm$ 0.3	15.84 $\pm$ 0.62
FEDGLOSS	1x	79.7 $\pm$ 0.4	84.4 $\pm$ 0.5	1.93 $\pm$ 0.03	47.2 $\pm$ 1.1	53.4 $\pm$ 0.5	16.22 $\pm$ 0.35
Centralized	-	87.1 SAM	86.3 $\widetilde{\text{SAM}}$	-	58.4 SAM	57.6 $\widetilde{\text{SAM}}$	-

and  $-1.9\%$  with  $\alpha = 0.05$  vs. respectively  $-3.2\%$  and  $-2.7\%$  of NAIVEFEDGLOSS w.r.t. SAM. In CIFAR100 instead,  $-7\%$  on  $\alpha = 0$  and  $-4.2\%$  on  $\alpha = 0.5$  of FEDGLOSS vs.  $-8.1\%$  and  $-5.2\%$  of its baseline. Aiming to reduce communication bottlenecks and improve performance, these results validate choosing FEDGLOSS over NAIVEFEDGLOSS. Further analyses in App. B.5.

### 6.4.2. The Role of Global Consistency and Flatness

Tab. 8 isolates the impact of global consistency and global sharpness minimization in FEDGLOSS. We recall FEDAVG with client-side SAM is FEDSAM and using ADMM only for aligning local and global convergence points is FEDDYN. Both components significantly impact performance, with their combination leading FEDGLOSS to the best overall results. FEDGLOSS is not prone to parameter explosion, achieving the best results even where FEDDYN fails to converge ( $\times$ ). The flatness of FEDGLOSS’ solutions w.r.t. FEDSAM in Fig. 6 confirms the efficacy of its strategy.

Table 8. Efficacy of global sharpness minimization in FEDGLOSS: ADMM for global consistency and server-side SAM for global sharpness minimization lead to the best performance. CIFARS, CNN with  $\alpha = 0$  and ResNet18 with  $\alpha \in \{0.05, 0.5\}$ .

CLIENT OPT	Method	Global Consistency	Global Flatness	CNN		ResNet18	
				CIFAR10	CIFAR100	CIFAR10	CIFAR100
SAM	FEDSAM	x	x	70.2 $\pm$ 0.9	28.7 $\pm$ 0.5	72.8 $\pm$ 0.1	38.5 $\pm$ 0.1
	FEDDYN	✓	-	79.3 $\pm$ 3.1	x	72.6 $\pm$ 0.2	39.6 $\pm$ 0.8
	FEDGLOSS	✓	✓	83.9 $\pm$ 0.4	50.6 $\pm$ 0.6	80.0 $\pm$ 0.3	47.2 $\pm$ 0.2
SGD	FEDAVG	x	x	59.9 $\pm$ 0.4	28.6 $\pm$ 0.7	72.6 $\pm$ 0.1	37.4 $\pm$ 0.2
	FEDDYN	✓	x	65.5 $\pm$ 0.3	x	70.2 $\pm$ 0.6	38.8 $\pm$ 0.6
	FEDGLOSS	✓	✓	69.5 $\pm$ 0.4	42.5 $\pm$ 0.6	79.1 $\pm$ 0.5	46.7 $\pm$ 0.6

## 7. Conclusion

This work tackled the challenge of limited generalization in heterogeneous Federated Learning (FL), prioritizing communication efficiency for real-world use. Building on research linking poor generalization to sharp minima in the loss landscape, we showed data heterogeneity worsens discrepancies between local and global loss surfaces, a problem not resolved by methods focusing only on local sharpness. To address this, our Federated Global Server-side Sharpness (FedGloss) finds flat minima in the *global* loss landscape with server-side Sharpness-Aware Minimization and achieves communication efficiency by approximating SAM’s sharpness through past global pseudo-gradients, distinguishing it from prior approaches. This work revealed SAM prevents ADMM-related parameter explosion by guiding optimization along flat directions, enabling stable updates in heterogeneous FL. Extensive evaluations showed FedGloss outperforms SOTA methods in accuracy, flatness and communication efficiency.

**Acknowledgments.** This study was carried out within the FAIR - Future Artificial Intelligence Research and received funding from the European Union Next-GenerationEU (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR) – MISSIONE 4 COMPONENTE 2, INVESTIMENTO 1.3 – D.D. 1555 11/10/2022, PE00000013). This manuscript reflects only the authors’ views and opinions, neither the European Union nor the European Commission can be considered responsible for them. We acknowledge the CINECA award under the ISCRA initiative for the availability of high-performance computing resources and support.

## 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. *ICLR*, 2021. [1](#), [2](#), [3](#), [5](#), [6](#), [16](#), [18](#), [19](#)
- [2] Rodolfo Stoffel Antunes, Cristiano André da Costa, Arne Küderle, Imrana Abdullahi Yari, and Björn Eskofier. Federated learning for healthcare: Systematic review and architecture proposal. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 13(4):1–23, 2022. [1](#)
- [3] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3 (1):1–122, 2011. [2](#), [5](#)
- [4] Debora Caldarola, Barbara Caputo, and Marco Ciccone. Improving generalization in federated learning by seeking flat minima. In *European Conference on Computer Vision*, pages 654–672. Springer, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [16](#), [18](#), [19](#), [20](#)
- [5] Debora Caldarola, Barbara Caputo, and Marco Ciccone. Window-based model averaging improves generalization in heterogeneous federated learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2263–2271, 2023. [3](#), [19](#)
- [6] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *Workshop on Data Privacy and Confidentiality*, 2019. [6](#)
- [7] Rong Dai, Xun Yang, Yan Sun, Li Shen, Xinmei Tian, Meng Wang, and Yongdong Zhang. Fedgamma: Federated learning with global sharpness-aware minimization. *IEEE Transactions on Neural Networks and Learning Systems*, 2023. [1](#), [2](#), [5](#), [6](#)
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [19](#)
- [9] Jiawei Du, Hanshu Yan, Jiashi Feng, Joey Tianyi Zhou, Liangli Zhen, Rick Siow Mong Goh, and Vincent YF Tan. Efficient sharpness-aware minimization for improved training of neural networks. *ICLR*, 2022. [3](#)
- [10] Jiawei Du, Daquan Zhou, Jiashi Feng, Vincent Tan, and Joey Tianyi Zhou. Sharpness-aware training for free. *Advances in Neural Information Processing Systems*, 35: 23439–23451, 2022. [3](#)
- [11] Lidia Fantauzzo, Eros Fani, Debora Caldarola, Antonio Tavera, Fabio Cermelli, Marco Ciccone, and Barbara Caputo. Feddrive: Generalizing federated learning to semantic segmentation in autonomous driving. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 11504–11511. IEEE, 2022. [1](#), [2](#)
- [12] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *Int. Conf. Machine Learn.*, 2021. [1](#), [2](#), [3](#), [20](#)
- [13] Liang Gao, Huazhu Fu, Li Li, Yingwen Chen, Ming Xu, and Cheng-Zhong Xu. Feddc: Federated learning with non-iid data via local drift decoupling and correction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10112–10121, 2022. [3](#)
- [14] Timur Garipov, Pavel Izmailov, Dmitrii Podoprikhin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31, 2018. [20](#)
- [15] Noah Golmant, Zhewei Yao, Amir Gholami, Michael Mahoney, and Joseph Gonzalez. pytorch-hessian-eigenthings: efficient pytorch hessian eigendecomposition, 2018. [20](#)
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *corr abs/1512.03385* (2015), 2015. [6](#), [19](#)
- [17] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997. [1](#), [2](#)
- [18] Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-iid data quagmire of decentralized machine learning. In *Int. Conf. Machine Learn.*, pages 4387–4398. PMLR, 2020. [19](#)
- [19] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *Neurips Workshop on Federated Learning*, 2019. [3](#), [6](#), [18](#)
- [20] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Federated visual classification with real-world data distribution. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*, pages 76–92. Springer, 2020. [1](#), [6](#), [16](#), [18](#), [19](#)
- [21] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. [19](#)
- [22] Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *Conference on Uncertainty in Artificial Intelligence*, 2018. [3](#)
- [23] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *ICLR*, 2019. [1](#)
- [24] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista

- Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021. [1](#), [2](#)
- [25] Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. Mime: Mimicking centralized stochastic algorithms in federated learning. *Advances in Neural Information Processing Systems*, 2020. [1](#), [3](#)
- [26] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *Int. Conf. Machine Learn.*, pages 5132–5143. PMLR, 2020. [1](#), [3](#), [5](#), [6](#), [16](#)
- [27] 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. *ICLR*, 2017. [1](#), [2](#)
- [28] Geeho Kim, Jinkyu Kim, and Bohyung Han. Communication-efficient federated learning with acceleration of global momentum. *arXiv preprint arXiv:2201.03172*, 2022. [3](#)
- [29] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. [6](#)
- [30] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [6](#), [19](#)
- [31] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018. [20](#)
- [32] Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning. *Computers & Industrial Engineering*, 149:106854, 2020. [1](#), [7](#)
- [33] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pages 965–978. IEEE, 2022. [7](#)
- [34] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60, 2020. [1](#)
- [35] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020. [1](#), [3](#), [6](#)
- [36] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *ICLR*, 2020. [1](#), [2](#)
- [37] Quande Liu, Cheng Chen, Jing Qin, Qi Dou, and Pheng-Ann Heng. Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1013–1023, 2021. [1](#), [2](#)
- [38] Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. Towards efficient and scalable sharpness-aware minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12360–12370, 2022. [3](#)
- [39] 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](#), [3](#), [6](#)
- [40] Jiaxu Miao, Zongxin Yang, Leilei Fan, and Yi Yang. Fedseg: Class-heterogeneous federated learning for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8042–8052, 2023. [1](#), [2](#)
- [41] Theodora Nevrataki, Anastasia Iliadou, George Ntolkeras, Ioannis Sfakianakis, Lazaros Lazaridis, George Maraslis, Nikolaos Asimopoulos, and George F Fragulis. A survey on federated learning applications in healthcare, finance, and data privacy/data security. In *AIP Conference Proceedings*. AIP Publishing, 2023. [1](#)
- [42] Jinseong Park, Hoki Kim, Yujin Choi, and Jaewook Lee. Differentially private sharpness-aware training. *Int. Conf. Machine Learn.*, 2023. [3](#), [5](#)
- [43] 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 *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. [19](#)
- [44] Henning Petzka, Michael Kamp, Linara Adilova, Cristian Sminchisescu, and Mario Boley. Relative flatness and generalization, 2021. [2](#)
- [45] Zhe Qu, Xingyu Li, Rui Duan, Yao Liu, Bo Tang, and Zhuo Lu. Generalized federated learning via sharpness aware minimization. In *Int. Conf. Machine Learn.*, pages 18250–18280. PMLR, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#)
- [46] Ashish Rauniyar, Desta Haileselassie Hagos, Debesh Jha, Jan Erik Häkegård, Ulas Bagci, Danda B Rawat, and Vladimir Vlassov. Federated learning for medical applications: A taxonomy, current trends, challenges, and future research directions. *IEEE Internet of Things Journal*, 2023. [1](#)
- [47] Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *ICLR*, 2021. [1](#), [2](#), [3](#)
- [48] Jae Hun Ro, Ananda Theertha Suresh, and Ke Wu. FedJAX: Federated learning simulation with JAX. *arXiv preprint arXiv:2108.02117*, 2021. [19](#)
- [49] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016. [3](#)
- [50] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted

- residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 6, 19
- [51] Donald Shenaj, Eros Fanì, Marco Toldo, Debora Caldarola, Antonio Tavera, Umberto Michieli, Marco Ciccone, Pietro Zanuttigh, and Barbara Caputo. Learning across domains and devices: Style-driven source-free domain adaptation in clustered federated learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 444–454, 2023. 1, 2
- [52] Yifan Shi, Yingqi Liu, Kang Wei, Li Shen, Xueqian Wang, and Dacheng Tao. Make landscape flatter in differentially private federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24552–24562, 2023. 2
- [53] Yan Sun, Li Shen, Shixiang Chen, Liang Ding, and Dacheng Tao. Dynamic regularized sharpness aware minimization in federated learning: Approaching global consistency and smooth landscape. *Int. Conf. Machine Learn.*, 2023. 1, 2, 4, 5, 6, 16, 18
- [54] Yan Sun, Li Shen, Tiansheng Huang, Liang Ding, and Dacheng Tao. Fedspeed: Larger local interval, less communication round, and higher generalization accuracy. *ICLR*, 2023. 1, 2, 5, 6, 16, 18
- [55] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *Int. Conf. Machine Learn.*, pages 1139–1147. PMLR, 2013. 3
- [56] Farshid Varno, Marzie Saghati, Laya Rafiee Seyyeri, Sharut Gupta, Stan Matwin, and Mohammad Havaei. Adabest: Minimizing client drift in federated learning via adaptive bias estimation. In *European Conference on Computer Vision*, pages 710–726. Springer, 2022. 1, 2, 3, 4, 7
- [57] Jie Wen, Zhixia Zhang, Yang Lan, Zhihua Cui, Jianghui Cai, and Wensheng Zhang. A survey on federated learning: challenges and applications. *International Journal of Machine Learning and Cybernetics*, 14(2):513–535, 2023. 1
- [58] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2575–2584, 2020. 6
- [59] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 19
- [60] Peng Xu, Bryan He, Christopher De Sa, Ioannis Mitliagkas, and Chris Re. Accelerated stochastic power iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 58–67. PMLR, 2018. 20
- [61] Riccardo Zaccone, Carlo Masone, and Marco Ciccone. Communication-efficient heterogeneous federated learning with generalized heavy-ball momentum. *arXiv preprint arXiv:2311.18578*, 2023. 3, 16, 18
- [62] Tuo Zhang, Lei Gao, Chaoyang He, Mi Zhang, Bhaskar Krishnamachari, and A Salman Avestimehr. Federated learning for the internet of things: Applications, challenges, and opportunities. *IEEE Internet of Things Magazine*, 5(1):24–29, 2022. 1