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

**Debora Caldarola** [1,*] **Pietro Cagnasso** [1], **Barbara Caputo** [1], **Marco Ciccone** [2,†]
[1] Politecnico di Torino, [2] 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 (**Fed**erated **Glo**bal **S**erver-side **S**harpness), 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 reaches flatter minima and better performance compared to state-of-the-art FL methods across various federated vision benchmarks.

## 1 Introduction

Federated Learning (FL) (McMahan et al., 2017) 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* (Reddi et al., 2021), *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 on the overall dataset. This approach holds immense potential for privacy-sensitive applications, proving its value in areas like healthcare (Liu et al., 2021; Antunes et al., 2022; Rauniyar et al., 2023; Nevrataki et al., 2023), finance (Nevrataki et al., 2023), autonomous driving (Fantauzzo et al., 2022; Shenaj et al., 2023; Miao et al., 2023), IoT (Zhang et al., 2022), and more (Li et al., 2020a; Wen et al., 2023). However, the real-world deployment of FL presents unique challenges stemming from data heterogeneity and communication costs (Li et al., 2020b). Clients gather their data influenced by various factors such as personal habits or geographical locations, leading to inherent differences across devices (Kairouz et al., 2021; Hsu et al., 2020; Shenaj et al., 2023). This results in the global model suffering from degraded performance and slower convergence (Li et al., 2020d; Karimireddy et al., 2020a;b; Caldarola et al., 2022), with instability emerging as client-specific optimization paths diverge from the global one. This phenomenon, known as *client drift* (Karimireddy et al., 2020b), 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 (Li et al., 2020c; Acar et al., 2021; Varno et al., 2022), a recent trend leverages the geometry of the

---

[*]Corresponding authors: debora.caldarola@polito.it, pietro.cagnasso@studenti.polito.it
[†]This work was started while the author was at Politecnico di Torino.

| (a) CIFAR10 $\alpha = 0$ | (b) CIFAR10 $\alpha = 0.05$ | (c) CIFAR100 $\alpha = 0$ | (d) CIFAR100 $\alpha = 0.5$ |

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 ($\alpha \to 0$) 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.

loss landscape to improve generalization (Caldarola et al., 2022; Qu et al., 2022; Sun et al., 2023b; Dai et al., 2023; Sun et al., 2023a). These methods build upon the notion that convergence to sharp minima correlates with poor generalization (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Jiang et al., 2019). FEDSAM (Caldarola et al., 2022; Qu et al., 2022) employs Sharpness-aware Minimization (SAM) (Foret et al., 2021) 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 (Sun et al., 2023b; Dai et al., 2023) demonstrated their effectiveness in various settings, they rely solely on local flatness, assuming that minimizing sharpness locally leads to a globally flat minimum. However, in real-world scenarios with significant data heterogeneity, there can be substantial discrepancies between local and global loss landscapes. As a consequence, **optimizing for local sharpness does not guarantee the global model will reside in a flat region** (Fig. 1). Addressing these limitations, FEDSMOO (Sun et al., 2023a) uses the alternating direction method of multipliers (ADMM) (Boyd et al., 2011) to include global sharpness information in SAM's local training. While this approach reduces the inconsistency between local and global geometries, it increases communication cost by **requiring double the bandwidth** in each round. This hinders its real-world applicability, as FL relies on minimizing communication overhead (*i.e.*, both message size and exchange frequency) to avoid network congestion and account for potential connection failures, that are common in practical deployments.

Given the limitations of existing methods, achieving convergence to global flat minima while maintaining communication efficiency in heterogeneous FL remains a critical challenge. To address this, we propose FEDGLOSS (**Fed**erated **Glo**bal **S**erver-side **S**harpness) 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, FEDGLOSS 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 Varno et al. (2022), 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) (McMahan et al., 2017) 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 (Li et al., 2020d; Acar et al., 2021; Reddi et al., 2021), the latter has explored its applications in real-world settings, spanning diverse domains like autonomous driving (Fantauzzo et al., 2022; Shenaj et al., 2023; Miao et al., 2023) and healthcare (Liu et al., 2021). 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 (Kairouz et al., 2021). 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 (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Petzka et al., 2021). SAM (Sharpness-Aware Minimization) (Foret et al., 2021) tackles this issue by guiding the optimization toward flatter regions, seeking minima that exhibit both low loss and low sharpness. FEDSAM (Caldarola et al., 2022; Qu et al., 2022) 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 (Sun et al., 2023b) uses perturbed gradients as SAM to reduce local overfitting, FEDGAMMA (Dai et al., 2023) combines the stochastic variance reduction of SCAFFOLD with SAM and Shi et al. (2023) 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 (Sun et al., 2023a) applies ADMM (Boyd et al., 2011) to the sharpness measure to enforce global and local consistency. This adds communication overhead, doubling the message size in each round and hindering its real-world practicality. In contrast, our work focuses on **minimizing global sharpness** while maintaining **communication efficiency**. Lastly, building on Stochastic Weight Averaging (Izmailov et al., 2018), other works (Caldarola et al., 2022; 2023) 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 (McMahan et al., 2017), 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 (Karimireddy et al., 2020b). Reddi et al. (2021) shows FEDAVG is equivalent to applying Stochastic Gradient Descent (SGD) (Ruder, 2016) 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** (Foret et al., 2021) **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 communi-

cation 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 (Li et al., 2020c) introduces a term to keep local parameters close to the global model, FEDDYN (Acar et al., 2021) employs ADMM to align local and global convergence points, ADABEST (Varno et al., 2022) adjusts local updates with an adaptive bias estimate, and SCAFFOLD (Karimireddy et al., 2020b) applies stochastic variance reduction. Momentum-based techniques (Sutskever et al., 2013) are also employed to maintain a consistent global trajectory, either on the server side (*e.g.*, FE-DAVGM (Hsu et al., 2019)) or by incorporating global information into local training (Karimireddy et al., 2020a; Kim et al., 2022; Gao et al., 2022; Zaccone et al., 2023). Unlike FEDDYN, where ADMM can lead to parameter explosion (Varno et al., 2022), 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 (Du et al., 2022a), or approximating them (Liu et al., 2022; Du et al., 2022b; Park et al., 2023). DP-SAT (Park et al., 2023) approximates the ascent step with the gradient from the previous iteration, and SAF (Du et al., 2022b) replaces SAM's sharpness approximation with the trajectory of weights learned during training. Aiming to the same goal, **FEDGLOSS approximates 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 (Foret et al., 2021) and FEDSAM (Caldarola et al., 2022; Qu et al., 2022).

## 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 $f(\boldsymbol{w}) : \mathcal{X} \to \mathcal{Y}$ parametrized by $\boldsymbol{w} \in \mathbb{R}^d$, 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:

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

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 (Reddi et al., 2021) solves Eq. (1) as

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

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

where $\Delta_{\boldsymbol{w}}^t$ 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, $\boldsymbol{w}^t$ the global model and $\boldsymbol{w}_k^t$ the local update resulting from training on client $k$'s data with CLIENTOPT for $E$ epochs. FEDAVG (McMahan et al., 2017) computes $\boldsymbol{w}^{t+1}$ as $\sum_{k \in \mathcal{C}^t} N_k/N \, \boldsymbol{w}_k^t$, corresponding to one SGD step on the pseudo-gradient $\Delta_{\boldsymbol{w}}^t$ with $\eta_s = 1$ (Reddi et al., 2021).

(a) Trained on class `sea`

(b) Trained on class `snail`

(c) Trained on class `skyscraper`

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.**

## 3.2 SHARPNESS-AWARE MINIMIZATION

SAM (Foret et al., 2021) jointly minimizes the loss value and the sharpness of the loss landscape by solving the min-max problem

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

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

$$\arg\max_{\|\boldsymbol{\epsilon}\| \leq \rho} f(\boldsymbol{w}) + \boldsymbol{\epsilon}^\top \nabla_{\boldsymbol{w}} f(\boldsymbol{w}) = \rho \frac{\nabla_{\boldsymbol{w}} f(\boldsymbol{w})}{\|\nabla_{\boldsymbol{w}} f(\boldsymbol{w})\|} \triangleq \hat{\boldsymbol{\epsilon}}(\boldsymbol{w}). \tag{5}$$

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

## 3.3 SAM IN FEDERATED LEARNING

FEDSAM (Caldarola et al., 2022; Qu et al., 2022) 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 $\min_{\boldsymbol{w}} \left\{ f^{\text{SAM}}(\boldsymbol{w}) = 1/C \sum_{k \in \mathcal{C}} f_k^{\text{SAM}}(\boldsymbol{w}) \right\}$, with $f_k^{\text{SAM}}(\boldsymbol{w}) \triangleq \max_{\|\boldsymbol{\epsilon}_k\| \leq \rho} f_k(\boldsymbol{w} + \boldsymbol{\epsilon}_k)$ with local perturbation $\boldsymbol{\epsilon}_k$. 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 Section 6.

Fig. 1 compares the loss surfaces of CNNs trained with FEDAVG and FEDSAM. On the easier CIFAR10, FEDSAM exhibits noticeably flatter minima w.r.t. FEDAVG, effectively navigating simpler landscapes. However, their performance difference diminishes with increasing dataset complexity (CIFAR100) and heterogeneity ($\alpha \to 0$). This suggests that **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 randomly selected clients during the last round with FEDSAM, distinguished by the locally seen class (results for all 5 clients in Appendix B). The inconsistency between local

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 $\alpha = 0$ (Caldarola et al., 2022; Varno et al., 2022), hence the lack of results (✗).

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

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 Hessian eigenvalues presented in Table 1. FEDSAM's local maximum Hessian eigenvalue, denoted by $\lambda_{1,l}$ and computed on each client's individual dataset, is significantly lower than the global eigenvalue $\lambda_{1,g}$, calculated on the overall dataset, on the more complex CIFAR100. This suggests that FEDSAM effectively achieves *local* convergence to flatter regions of the loss landscape on individual devices. However, 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 $\min_{\boldsymbol{w}} \left\{ \mathcal{F}(\boldsymbol{w}) = \frac{1}{C} \sum_{k \in \mathcal{C}} \mathcal{F}_k(\boldsymbol{w}) \right\}$, with $\mathcal{F}_k(\boldsymbol{w}) \triangleq \max_{\|\boldsymbol{\epsilon}\| \leq \rho} f_k(\boldsymbol{w} + \boldsymbol{\epsilon})$, where $\boldsymbol{\epsilon}$ is the global perturbation. Calculating the true $\boldsymbol{\epsilon}$ value requires the global gradient $\nabla_{\boldsymbol{w}} f$ (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 FEDSMOO (Sun et al., 2023a) tackles this issue by using ADMM on the sharpness with the constraint $\boldsymbol{\epsilon} = \boldsymbol{\epsilon}_k$, it necessitates transmitting $\boldsymbol{\epsilon}$ alongside the model parameters $\boldsymbol{w}$ 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 $\boldsymbol{w}^{t+1} \leftarrow \boldsymbol{w}^t - \eta_s \nabla_{\boldsymbol{w}} \mathcal{F}(\boldsymbol{w})|_{\boldsymbol{w}^t + \hat{\boldsymbol{\epsilon}}^t(\boldsymbol{w})}$, where $\hat{\boldsymbol{\epsilon}}^t$ is the global perturbation at each round $t$. However, a key challenge arises: the computation of both $\hat{\boldsymbol{\epsilon}}^t$ 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 $\boldsymbol{w}^t$ to a subset $\mathcal{C}^t$ of clients, which update it using their local data. With the resulting pseudo-gradient $\Delta_{\boldsymbol{w}}^t$, $\hat{\boldsymbol{\epsilon}}^t(\boldsymbol{w}) = \rho(\Delta_{\boldsymbol{w}}^t / \|\Delta_{\boldsymbol{w}}^t\|)$ and the perturbed model $\tilde{\boldsymbol{w}}^t = \boldsymbol{w}^t + \hat{\boldsymbol{\epsilon}}^t(\boldsymbol{w})$.

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

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

| Method | SERVEROPT | CLIENTOPT | Global Flatness | Communication Cost | Local Computation Cost |
|---|---|---|---|---|---|
| FEDSAM (Caldarola et al., 2022; Qu et al., 2022) | SGD | SAM | ✗ | 1× | 2× |
| FEDDYN (Acar et al., 2021) + SAM | SGD | SAM | ✗ | 1× | 2× |
| FEDSPEED (Sun et al., 2023b) | SGD | Similar to SAM | ✗ | 1× | 2× |
| FEDGAMMA (Dai et al., 2023) | SGD | SAM | ✓ | 2× | 2× |
| FEDSMOO (Sun et al., 2023a) | SGD | SAM | ✓ | 2× | 2× |
| **FEDGLOSS** | **SAM** | **Any optimizer** | ✓ | 1× | 1× or 2× |

- **Step 2**: the server transmits $\tilde{\boldsymbol{w}}^t$ to the *same* $\mathcal{C}^t$, which compute their update $\tilde{\boldsymbol{w}}_k^t \; \forall k$. The resulting global pseudo-gradient $\tilde{\Delta}_{\boldsymbol{w}}^t \triangleq \sum_{k \in \mathcal{C}^t} N_k/N (\tilde{\boldsymbol{w}}^t - \tilde{\boldsymbol{w}}_k^t)$ is an estimate of $\nabla_{\boldsymbol{w}} \mathcal{F}(\boldsymbol{w})|_{\boldsymbol{w}^t + \hat{\boldsymbol{\epsilon}}^t(\boldsymbol{w})}$.

This two-step approach, referred to as NAIVEFEDGLOSS, is conceptually simple but suffers from communication *in*efficiency, 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 (Park et al., 2023), FEDGLOSS estimates $\hat{\boldsymbol{\epsilon}}^t$ using the **perturbed global pseudo-gradient from the previous round** $\tilde{\Delta}_{\boldsymbol{w}}^{t-1}$ at each round $t$. This approach leverages available information *without incurring extra communications* and avoiding 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 $\boldsymbol{w}^t$ as

$$① \; \tilde{\boldsymbol{\epsilon}}^t(\boldsymbol{w}) \triangleq \rho \frac{\tilde{\Delta}_{\boldsymbol{w}}^{t-1}}{\|\tilde{\Delta}_{\boldsymbol{w}}^{t-1}\|}$$

$$② \; \tilde{\boldsymbol{w}}^t \leftarrow \boldsymbol{w}^t + \tilde{\boldsymbol{\epsilon}}^t(\boldsymbol{w})$$

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

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

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 Table 2, FEDGLOSS enables SAM on the server side while **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 FEDSMOO, 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 $\tilde{\boldsymbol{\epsilon}}^t$ (FEDGLOSS) and the true $\hat{\boldsymbol{\epsilon}}^t$ (NAIVEFEDGLOSS) is

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

where $\tilde{\Delta}_{\boldsymbol{w}}^{t-1}$ is computed using the updates of the clients in $\mathcal{C}^{t-1}$ and $\tilde{\Delta}_{\boldsymbol{w}}^t$ with $\mathcal{C}^t$. Eq. (6) suggests $\delta_\epsilon^t$ is minimized when $\tilde{\Delta}_{\boldsymbol{w}}^{t-1}$ and $\tilde{\Delta}_{\boldsymbol{w}}^t$ are aligned. However, in real-world heterogeneous FL, *i*) to due 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 (Karimireddy et al., 2020b). As a consequence, $\delta_\epsilon^t \not\to 0$ necessarily.

To align local and global objectives - ensuring client and server gradient alignment and minimizing Eq. (6) - FEDGLOSS leverages the Alternating Direction Method of Multipliers (ADMM) (Boyd et al., 2011) on $\boldsymbol{w}^t$ (Acar et al., 2021; Sun et al., 2023a;b). 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 Karimireddy et al. (2020b); Dai et al. (2023)). In contrast, ADMM has been proved to converge under arbitrary heterogeneity Acar et al. (2021) and can thus be leveraged as a base algorithm for FEDGLOSS, as shown in Algorithm 1 in Appendix A. ADMM makes use of the augmented Lagrangian function $\mathcal{L}(\boldsymbol{w}, \boldsymbol{W}, \sigma) = \sum_{k \in \mathcal{C}} L(\boldsymbol{w}, \boldsymbol{w}_k, \sigma_k)$ where $\boldsymbol{W} = \{\boldsymbol{w}_1, \cdots, \boldsymbol{w}_C\}$ and $\sigma$ is the Lagrangian multiplier. The problem solved by $\mathcal{L}$ is

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

with $\beta > 0$ being an hyperparameter. Eq. (7) is split into $C$ sub-problems of the form $\boldsymbol{w}_{k,E} = \arg\min_{\boldsymbol{w}_k} \{f_k - \sigma_k^\top (\boldsymbol{w}^t - \boldsymbol{w}_k) + \frac{1}{2\beta} \| \boldsymbol{w}^t - \boldsymbol{w}_k^t \|^2 \}$. The local dual variable is updated as $\sigma_k \leftarrow \sigma_k - \frac{1}{\beta} (\boldsymbol{w}_{k,E}^t - \boldsymbol{w}_{k,0}^t)$. The global one $\sigma$ is updated by adding the averaged $\boldsymbol{w}_k - \boldsymbol{w}^t \; \forall k \in \mathcal{C}$.

Figure 4 confirms the effect of ADMM on gradient alignment: the difference between the true and approximated perturbation, $\delta_\epsilon^t$ (Eq. (6)), decreases over training rounds and with the use of Lagrangian multipliers.



Figure 4: 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

Appendix C details implementation and hyperparameter settings.

**Federated datasets.** We leverage established FL benchmarks (Caldas et al., 2019; Hsu et al., 2020; 2019). *Small-scale image classification:* following (Hsu et al., 2020; Caldarola et al., 2022), the federated versions of CIFAR10 (10 classes) and CIFAR100 (100 classes) (Krizhevsky, 2009) split the respective $50k$ training images in 100 clients with 500 images each. The data distribution is controlled by the Dirichlet's parameter $\alpha \in \{0, 0.05\}$ for CIFAR10 and $\{0, 0.5\}$ for CIFAR100 (Hsu et al., 2019). Lower $\alpha$ signifies increased heterogeneity, with $\alpha = 0$ being the most challenging scenario where each client holds samples from one class. *Large-scale image classification:* LANDMARKS-USER-160K (2,028 classes) (Hsu et al., 2020) is the federated Google Landmarks v2 (Weyand et al., 2020) 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 (Hsu et al., 2020; Caldarola et al., 2022), we use a Convolutional Neural Network (CNN) similar to

(a) Local model on class `sea` w/ **FEDGLOSS**

(b) Local model on class `snail` w/ **FEDGLOSS**

(c) Local model on class `sea` w/ **FEDSMOO**

(d) Local model on class `snail` w/ **FEDSMOO**

Figure 5: **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' local model (*net*) as reference. **FEDGLOSS achieves better consistency w.r.t. FEDSMOO.**

LeNet5 (LeCun et al., 1998) on both CIFAR10 ($T = 10k$) and CIFAR100 ($T = 20k$). Experiments with ResNet18 (He et al., 2015) run for $10k$ rounds. For LANDMARKS-USER-160K, we train MobileNetv2 (Sandler et al., 2018; Hsu et al., 2020) ($T = 1.3k$), 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 (Section 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 (McMahan et al., 2017), FedProx (Li et al., 2020c), FedDyn (Acar et al., 2021) and Scaffold (Karimireddy et al., 2020b), while use SAM FedSAM (Caldarola et al., 2022; Qu et al., 2022), FedDyn + SAM, FedSpeed (Sun et al., 2023b), FedGamma (Dai et al., 2023) and FedSmoo (Sun et al., 2023a).

## 6.2  ACHIEVING LOCAL-GLOBAL SHARPNESS CONSISTENCY

To assess the effectiveness of FEDGLOSS in promoting consistency between local and global loss landscapes, Fig. 5 replicates the analysis previously conducted on FEDSAM (Fig. 2) for direct comparison. The behavior of local models is shown from both local and global perspectives ("*Local loss*" and "*Global loss*", respectively). Appendix 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. 5a and 5b). 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. 5c and 5d 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.

Table 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 (Section 6.1) on vision tasks in heterogeneous federated settings. Appendix B.4 discusses results in homogeneous FL.

(a) C10 $\alpha = 0$ CNN FEDGLOSS *vs.* FEDSAM    (b) C10 $\alpha = 0$ CNN FEDGLOSS *vs.* FEDSMOO    (c) C100 $\alpha = 0$ CNN FEDGLOSS *vs.* FEDSAM    (d) C100 $\alpha = 0$ CNN FEDGLOSS *vs.* FEDSMOO    (e) C10 $\alpha = 0.05$ ResNet18 FEDGLOSS *vs.* FEDSAM    (f) C10 $\alpha = 0.05$ ResNet18 FEDGLOSS *vs.* FEDSMOO

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.

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**.

| | Method | Comm. Cost | CNN | | | | ResNet18 | |
| | | | CIFAR10 | | CIFAR100 | | CIFAR10 | CIFAR100 |
| | | | $\alpha = 0$ | $\alpha = 0.05$ | $\alpha = 0$ | $\alpha = 0.5$ | $\alpha = 0.05$ | $\alpha = 0.5$ |
|---|---|---|---|---|---|---|---|---|
| **Client SGD** | FEDAVG | $1\times$ | $59.9_{\pm0.4}$ | $65.7_{\pm1.0}$ | $28.6_{\pm0.7}$ | $38.5_{\pm0.5}$ | $72.6_{\pm0.1}$ | $37.4_{\pm0.2}$ |
| | FEDPROX | $1\times$ | $59.8_{\pm0.5}$ | $65.6_{\pm1.0}$ | $28.8_{\pm0.7}$ | $38.7_{\pm0.4}$ | $72.2_{\pm0.2}$ | $37.6_{\pm0.1}$ |
| | FEDDYN | $1\times$ | $65.5_{\pm0.3}$ | $70.1_{\pm1.2}$ | ✗ | ✗ | $70.2_{\pm0.6}$ | $38.8_{\pm0.6}$ |
| | SCAFFOLD | $2\times$ | $25.1_{\pm3.7}$ | $54.0_{\pm2.6}$ | ✗ | $30.0_{\pm1.1}$ | $70.8_{\pm0.6}$ | $38.6_{\pm0.1}$ |
| | **FEDGLOSS** | $1\times$ | $\mathbf{69.5_{\pm0.4}}$ | $\mathbf{75.5_{\pm0.3}}$ | $\mathbf{42.5_{\pm0.6}}$ | $\mathbf{47.9_{\pm0.5}}$ | $\mathbf{79.1_{\pm0.5}}$ | $\mathbf{46.7_{\pm0.6}}$ |
| **Client SAM** | FEDSAM | $1\times$ | $70.2_{\pm0.9}$ | $71.5_{\pm1.08}$ | $28.7_{\pm0.5}$ | $39.6_{\pm0.5}$ | $72.8_{\pm0.1}$ | $38.5_{\pm0.1}$ |
| | FEDDYN | $1\times$ | $79.3_{\pm3.1}$ | $81.5_{\pm0.6}$ | ✗ | ✗ | $72.6_{\pm0.2}$ | $39.6_{\pm0.8}$ |
| | FEDSPEED | $1\times$ | $70.9_{\pm0.4}$ | $72.3_{\pm1.1}$ | $28.9_{\pm0.5}$ | $39.7_{\pm0.5}$ | $72.6_{\pm0.1}$ | $38.8_{\pm0.6}$ |
| | FEDGAMMA | $2\times$ | $58.9_{\pm1.8}$ | $61.9_{\pm1.8}$ | ✗ | $29.4_{\pm1.4}$ | $72.2_{\pm0.1}$ | $38.8_{\pm0.3}$ |
| | FEDSMOO | $2\times$ | $81.3_{\pm0.5}$ | $82.8_{\pm0.6}$ | $47.8_{\pm0.5}$ | $51.7_{\pm0.46}$ | $75.3_{\pm0.6}$ | $44.8_{\pm0.5}$ |
| | **FEDGLOSS** | $1\times$ | $\mathbf{83.9_{\pm0.4}}$ | $\mathbf{84.4_{\pm0.5}}$ | $\mathbf{50.6_{\pm0.6}}$ | $\mathbf{53.4_{\pm0.5}}$ | $\mathbf{80.0_{\pm0.3}}$ | $\mathbf{47.2_{\pm0.2}}$ |

### 6.3.1 FEDGLOSS ON STANDARD FEDERATED BENCHMARKS

Table 3 presents results on CIFAR100 and CIFAR10 with varying levels of heterogeneity on the CNN model. Several observations highlight the advantages of FEDGLOSS. It is straightforward to notice how FEDGLOSS achieves the **best results** among both SGD and SAM-based approaches while maintaining **communication efficiency**. FEDGLOSS with local SAM consistently outperforms the best-performing SOTA FEDSMOO by $\approx 2.5$ percentage points in accuracy across all dataset configurations *with half the communication cost*. FEDGLOSS also reaches the **flattest global minima** (*e.g.*, $\lambda_1^{\text{FEDGLOSS}} = 2.03$ *vs.* $\lambda_1^{\text{FEDSMOO}} = 15.37$ on CIFAR10 with $\alpha = 0$), 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. As studied in (Varno et al., 2022), FEDDYN suffers from parameter explosion in highly heterogeneous settings, failing to converge on CIFAR100. Differently, FEDGLOSS successfully employs ADMM to align global and local solutions, with the best results under extreme heterogeneity. Similarly, studies showed SCAFFOLD performs poorly in complex heterogeneous environments (Li et al., 2022; Caldarola et al., 2022), resulting in its inability to converge on CIFAR100 alongside FEDGAMMA. Lastly, the last two columns of Table 3 further confirm FEDGLOSS' 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).



(a) CIFAR10           (b) CIFAR100

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

Figure 8: Trend of model parameters norm, $\| \boldsymbol{w}^t \|_2$, 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**.

### 6.3.2 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. Table 4 compares FEDGLOSS with local SAM against the best-performing baselines. FEDGLOSS is among the few methods, alongside FEDSAM and FEDSMOO, outperforming FEDAVG. Similarly to the CIFAR100 results (Section 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.

Table 4: **MobileNetv2** on LANDMARKS-USER-160K.

| Method | Comm. cost | Accuracy |
|---|---|---|
| FEDAVG | 1× | 56.3±0.2 |
| FEDPROX | 1× | 55.0±0.2 |
| FEDDYN | 1× | 55.2±0.6 |
| SCAFFOLD | 2× | ✗ |
| FEDSAM | 1× | 56.7±0.1 |
| FEDDYN + SAM | 1× | 56.0±1.3 |
| FEDGAMMA | 2× | ✗ |
| FEDSMOO | 2× | 59.5±0.1 |
| FEDGLOSS | 1× | **59.7±1.2** |

### 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 Varno et al. (2022). This occurs as multiple gradients accumulate in the global dual variable $\sigma$ (Section 5), causing the parameter norms to grow uncontrollably. However, empirical results indicate that this issue is mitigated with SAM (*e.g.*, see FEDDYN vs. FEDGLOSS in Table 3). We hypothesize this is due 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.

### 6.3.4 COMMUNICATION EFFICIENCY WITH FEDGLOSS

Table 5: **Communication costs** comparison w.r.t. FEDAVG. "-" for not reached accuracy, "✗" for non-convergence. GLDV2 is LANDMARKS-USER-160K.

| | Method | CNN | | | | ResNet18 | | | | MobileNetv2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | CIFAR10 $\alpha = 0$ | | CIFAR100 $\alpha = 0$ | | CIFAR10 $\alpha = 0.05$ | | CIFAR100 $\alpha = 0.5$ | | GLDV2 | |
| | | *Rounds* | *Cost* | *Rounds* | *Cost* | *Rounds* | *Cost* | *Rounds* | *Cost* | *Rounds* | *Cost* |
| Client SGD | 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 | 2.0k | **0.2B** | | ✗ | - | - | 3.5k | 0.4B | - | - |
| | SCAFFOLD | - | - | - | - | - | - | 8.9k | 1.8B | | ✗ |
| | FEDGLOSS | 3.4k | 0.3B | 5k | **0.3B** | 2.4k | **0.2B** | 1.9k | **0.2B** | 1.3k | 1B |
| Client SAM | FEDSAM | 6.3k | 0.6B | 18.3k | 0.9B | 9.2k | 0.9B | 7.8k | 0.8B | 1.3k | 1B |
| | FEDDYN | 3.0k | 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 |
| | FEDGLOSS | 2.2k | **0.2B** | 6.3k | **0.3B** | 2.4k | **0.2B** | 1.9k | **0.2B** | 200 | **0.2B** |

Communication cost is the main bottleneck in FL (Li et al., 2020a), making its optimization a relevant challenge. As already previously highlighted, FEDGLOSS considers communication efficiency

Figure 9: Loss barriers resulting from interpolating NAÏVEFEDGLOSS and FEDGLOSS' models, which are found in the same basin. CIFAR datasets, CNN.

Table 6: **FEDGLOSS** *vs.* **NAÏVEFEDGLOSS** 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% | $\lambda_1$ ($\downarrow$) |
| NAÏVEFEDGLOSS | 2× | $77.6_{\pm0.3}$ | $\mathbf{83.9_{\pm0.2}}$ | $2.78_{\pm0.13}$ | $\mathbf{42.6_{\pm0.8}}$ | $\mathbf{50.8_{\pm0.1}}$ | $16.93_{\pm0.27}$ |
| FEDGLOSS | 1× | $\mathbf{78.9_{\pm0.5}}$ | $\mathbf{83.9_{\pm0.4}}$ | $\mathbf{2.03_{\pm0.05}}$ | $39.5_{\pm0.9}$ | $50.6_{\pm0.6}$ | $17.18_{\pm0.97}$ |
| NAÏVEFEDGLOSS | 2× | $78.7_{\pm0.1}$ | $\mathbf{84.4_{\pm0.2}}$ | $2.75_{\pm0.09}$ | $\mathbf{49.4_{\pm0.5}}$ | $\mathbf{53.7_{\pm0.3}}$ | $15.84_{\pm0.52}$ |
| FEDGLOSS | 1× | $\mathbf{79.7_{\pm0.4}}$ | $\mathbf{84.4_{\pm0.5}}$ | $\mathbf{1.93_{\pm0.03}}$ | $47.2_{\pm1.1}$ | $53.4_{\pm0.5}$ | $16.22_{\pm0.35}$ |
| **Centralized** | | 87.1 SAM | 86.3 $\widetilde{\text{SAM}}$ | | 58.4 SAM | 57.6 $\widetilde{\text{SAM}}$ | |

its primacy concern. Defined $B$ the number of bits exchanged by FEDAVG in $T$ training rounds, Table 5 studies FEDGLOSS's communication cost against the SOTA baselines in terms of rounds necessary to reach FEDAVG's performance and quantity of exchanged bits. 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 Appendix B.5.

## 6.4 ABLATION STUDIES

### 6.4.1 COMMUNICATION-EFFICIENT SHARPNESS

This section studies the efficacy of using the pseudo-gradient from the previous round $\tilde{\Delta}_{\boldsymbol{w}}^{t-1}$ (Eq. (6)) as an estimate of the sharpness measure. FEDGLOSS uses past gradients as a reliable indication on the global loss landscape and, by aligning global and local optimization paths through ADMM, enables consistent trajectories across rounds.

Table 6 compares FEDGLOSS with its baseline, NAÏVEFEDGLOSS (Section 5), which computes the true perturbation $\hat{\epsilon}^t$ at the expense of doubled communication costs. FEDGLOSS achieves accuracy comparable to NAÏVEFEDGLOSS while maintaining communication efficiency, with minimal or negligible gap in performance. This aligns with the observed sharpness of the achieved minima ($\lambda_1$). To ensure a fair comparison in communication cost, Table 6 also compares FEDGLOSS' final accuracy with NAÏVEFEDGLOSS' performance at 50% training progress, showing that FEDGLOSS achieves higher accuracy with the same number of exchanges, benefiting from the additional global optimization steps deriving from its communication-efficient strategy. This shows our approximation does not slow convergence. In addition, our strategy in centralized settings lowers the performance w.r.t. SAM, thus reducing our centralized upper bound w.r.t. NAÏVEFEDGLOSS. At equal performance, FEDGLOSS narrows the gap to the upper bound: $-2.4\%$ on CIFAR10 with $\alpha = 0$ and $-1.9\%$ with $\alpha = 0.05$ *vs.* respectively $-3.2\%$ and $-2.7\%$ of NAÏVEFEDGLOSS 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. Lastly, our sharpness approximation does not steer the optimization path: models trained with FEDGLOSS and NAÏVEFEDGLOSS end up in the same basin (no loss barrier), with

Table 7: 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 | ✗ | ✗ | $70.2_{\pm 0.9}$ | $28.7_{\pm 0.5}$ | $72.8_{\pm 0.1}$ | $38.5_{\pm 0.1}$ |
| | FEDDYN | ✓ | ✗ | $79.3_{\pm 3.1}$ | ✗ | $72.6_{\pm 0.2}$ | $39.6_{\pm 0.8}$ |
| | FEDGLOSS | ✓ | ✓ | $\mathbf{83.9_{\pm 0.4}}$ | $\mathbf{50.6_{\pm 0.6}}$ | $\mathbf{80.0_{\pm 0.3}}$ | $\mathbf{47.2_{\pm 0.2}}$ |
| SGD | FEDAVG | ✗ | ✗ | $59.9_{\pm 0.4}$ | $28.6_{\pm 0.7}$ | $72.6_{\pm 0.1}$ | $37.4_{\pm 0.2}$ |
| | FEDDYN | ✓ | ✗ | $65.5_{\pm 0.3}$ | ✗ | $70.2_{\pm 0.6}$ | $38.8_{\pm 0.6}$ |
| | FEDGLOSS | ✓ | ✓ | $\mathbf{69.5_{\pm 0.4}}$ | $\mathbf{42.5_{\pm 0.6}}$ | $\mathbf{79.1_{\pm 0.5}}$ | $\mathbf{46.7_{\pm 0.6}}$ |

similar flatness (Fig. 9). Aiming to reduce the communication bottleneck while achieving superior performance, **these results confirm our choice of FEDGLOSS over NAIVEFEDGLOSS**.

### 6.4.2 THE ROLE OF GLOBAL CONSISTENCY AND FLATNESS

Table 7 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 (✗). The flatness of FEDGLOSS' solutions w.r.t. FEDSAM in Fig. 6 confirms the efficacy of its strategy.

## 7 DISCUSSION

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 issue, we introduced Federated Global Server-side Sharpness (FEDGLOSS), which finds flat minima in the *global* loss landscape using server-side Sharpness-Aware Minimization (SAM). FEDGLOSS achieves communication efficiency by approximating SAM's sharpness through past global pseudo-gradients, distinguishing it from prior approaches. In addition, by not constraining the clients to use SAM as local optimizer, FEDGLOSS' required computational cost can be adapted depending on the local available resources.

Importantly, this work revealed SAM prevents ADMM-related parameter explosion by guiding optimization along flat directions, which reflects in reduced model parameters' norm, enabling stable updates in heterogeneous FL. A promising future direction could involve exploring alternative approaches to ADMM to align local and global solutions, with the goal of avoiding stateful clients.

Extensive evaluations showed FEDGLOSS outperforms state-of-the-art methods in accuracy, flatness of the solution and communication efficiency, making it a strong candidate for real-world FL applications.

### REFERENCES

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, 3, 4, 7, 8, 9, 25

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

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, 3, 8

Debora Caldarola, Barbara Caputo, and Marco Ciccone. Improving generalization in federated learning by seeking flat minima. In *European Conference on Computer Vision*, pp. 654–672. Springer, 2022. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 25, 26

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*, pp. 2263–2271, 2023. 3, 25

Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečnỳ, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *Workshop on Data Privacy and Confidentiality*, 2019. 8

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. 2, 3, 7, 8, 9

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*, pp. 248–255. Ieee, 2009. 25

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*, 2022a. 4

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, 2022b. 4

Lidia Fantauzzo, Eros Fanì, 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)*, pp. 11504–11511. IEEE, 2022. 1, 3

Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *Int. Conf. Machine Learn.*, 2021. 2, 3, 4, 5, 26

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*, pp. 10112–10121, 2022. 4

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. 26

Noah Golmant, Zhewei Yao, Amir Gholami, Michael Mahoney, and Joseph Gonzalez. pytorch-hessian-eigenthings: efficient pytorch hessian eigendecomposition, October 2018. URL https://github.com/noahgolmant/pytorch-hessian-eigenthings. 26

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. corr abs/1512.03385 (2015), 2015. 9, 24

Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997. 2, 3

Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-iid data quagmire of decentralized machine learning. In *Int. Conf. Machine Learn.*, pp. 4387–4398. PMLR, 2020. 24

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. 4, 8, 24

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*, pp. 76–92. Springer, 2020. 1, 8, 9, 24, 25

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015. 24

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

Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *ICLR*, 2019. URL https://openreview.net/forum?id=SJgIPJBFvH. 2

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, 3

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*, 2020a. 1, 4

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.*, pp. 5132–5143. PMLR, 2020b. 1, 3, 4, 8, 9, 22

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. 2, 3

Geeho Kim, Jinkyu Kim, and Bohyung Han. Communication-efficient federated learning with acceleration of global momentum. *arXiv preprint arXiv:2201.03172*, 2022. 4

Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. 8

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. 9, 24

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. 26

Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning. *Computers & Industrial Engineering*, 149:106854, 2020a. 1, 11

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)*, pp. 965–978. IEEE, 2022. 10

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, 2020b. 1

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, 2020c. 1, 4, 9

Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *ICLR*, 2020d. 1, 3

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*, pp. 1013–1023, 2021. 1, 3

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*, pp. 12360–12370, 2022. 4

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*, pp. 1273–1282. PMLR, 2017. 1, 3, 4, 9

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*, pp. 8042–8052, 2023. 1, 3

Theodora Nevrataki, Anastasia Iliadou, George Ntolkeras, Ioannis Sfakianakis, Lazaros Lazaridis, George Maraslidis, Nikolaos Asimopoulos, and George F Fragulis. A survey on federated learning applications in healthcare, finance, and data privacy/data security. In *AIP Conference Proceedings*, volume 2909. AIP Publishing, 2023. 1

Jinseong Park, Hoki Kim, Yujin Choi, and Jaewook Lee. Differentially private sharpness-aware training. *Int. Conf. Machine Learn.*, 2023. 4, 7

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*, pp. 8024–8035. Curran Associates, Inc., 2019. 24

Henning Petzka, Michael Kamp, Linara Adilova, Cristian Sminchisescu, and Mario Boley. Relative flatness and generalization, 2021. URL https://arxiv.org/abs/2001.00939. 3

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.*, pp. 18250–18280. PMLR, 2022. 2, 3, 4, 5, 7, 9

Ashish Rauniyar, Desta Haileselassie Hagos, Debesh Jha, Jan Erik Håkegård, 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

Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *ICLR*, 2021. 1, 3, 4

Jae Hun Ro, Ananda Theertha Suresh, and Ke Wu. FedJAX: Federated learning simulation with JAX. *arXiv preprint arXiv:2108.02117*, 2021. 24

Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016. 3

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*, pp. 4510–4520, 2018. 9, 25

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*, pp. 444–454, 2023. 1, 3

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*, pp. 24552–24562, 2023. 3

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.*, 2023a. 2, 3, 6, 7, 8, 9

Yan Sun, Li Shen, Tiansheng Huang, Liang Ding, and Dacheng Tao. Fedspeed: Larger local interval, less communication round, and higher generalization accuracy. *ICLR*, 2023b. 2, 3, 7, 8, 9

Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *Int. Conf. Machine Learn.*, pp. 1139–1147. PMLR, 2013. 4

Farshid Varno, Marzie Saghayi, Laya Rafiee Sevyeri, Sharut Gupta, Stan Matwin, and Mohammad Havaei. Adabest: Minimizing client drift in federated learning via adaptive bias estimation. In *European Conference on Computer Vision*, pp. 710–726. Springer, 2022. 1, 3, 4, 6, 10, 11

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

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*, pp. 2575–2584, 2020. 8

Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018. 24

Peng Xu, Bryan He, Christopher De Sa, Ioannis Mitliagkas, and Chris Re. Accelerated stochastic power iteration. In *International Conference on Artificial Intelligence and Statistics*, pp. 58–67. PMLR, 2018. 26

Riccardo Zaccone, Carlo Masone, and Marco Ciccone. Communication-efficient heterogeneous federated learning with generalized heavy-ball momentum. *arXiv preprint arXiv:2311.18578*, 2023. 4

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

APPENDIX

This appendix is organized as follows:

## A  ALGORITHM

Algorithm 1 summarizes FEDGLOSS, using as an example SGD or SAM as local optimizers, differently highlighted. The comment colors replicate those of Fig. 3.

---

**Algorithm 1** FEDGLOSS with SAM or SGD as local optimizers

---

1: **Input:** Global model $\boldsymbol{w}$, clients $\mathcal{C}$, rounds $T$, local iterations $E$, clients' learning rate $\eta$, clients' SAM neighborhood size $\rho_l$, FEDGLOSS neighborhood size $\rho$, Lagrangian hyperparameter $\beta$.
2: **Initialize:** $\boldsymbol{w}^0$, $\sigma^0 = \sigma_k = 0$, $\Delta_{\boldsymbol{w}}^0 = 0$.
3: **for** each round $t \in [1, T]$ **do**
4:     $\tilde{\boldsymbol{\epsilon}}^t(\boldsymbol{w}) = \rho \frac{\Delta_{\boldsymbol{w}}^{t-1}}{\|\Delta_{\boldsymbol{w}}^{t-1}\|_2}$                ▷ Global perturbation with previous pseudo-grad
5:     $\tilde{\boldsymbol{w}}^t = \boldsymbol{w}^t + \tilde{\boldsymbol{\epsilon}}^t(\boldsymbol{w})$             ▷ Server-side approximated **FEDGLOSS ascent step**
6:     Randomly select a subset of clients $\mathcal{C}^t \subset \mathcal{C}$
7:     **for** each client $k \in \mathcal{C}^t$ in parallel **do**
8:         $\boldsymbol{w}_{k,0} = \tilde{\boldsymbol{w}}^t$           ▷ Initialize local model with *perturbed* global model $\tilde{\boldsymbol{w}}^t$
9:         Set iteration counter $i = 1$
10:        **for** each epoch $e \in [1, E]$ **do**
11:           **for** each batch $\mathcal{B} \in \mathcal{D}_k$ **do**
12:             $\boldsymbol{g}_{k,i} = \nabla f_{\mathcal{B}}(\boldsymbol{w}_{k,i-1})$                   ▷ Local SGD gradient
13:             $\hat{\boldsymbol{\epsilon}}_{k,i} = \rho_l \frac{\boldsymbol{g}_{k,i}}{\|\boldsymbol{g}_{k,i}\|_2}$                   ▷ SAM local perturbation
14:             $\boldsymbol{g}_{k,i} = \nabla f_{\mathcal{B}}(\boldsymbol{w}_{k,i-1} + \hat{\boldsymbol{\epsilon}}_{k,i})$           ▷ Local sharpness-aware gradient
15:             $\boldsymbol{w}_{k,i} \leftarrow \boldsymbol{w}_{k,i-1} - \eta[\boldsymbol{g}_{k,i} - \sigma_k + (\boldsymbol{w}_{k,i-1} - \boldsymbol{w}_{k,0})/\beta]$    ▷ Local update with ADMM
16:             $i \leftarrow i + 1$
17:           **end for**
18:        **end for**
19:         $\sigma_k \leftarrow \sigma_k - (\boldsymbol{w}_{k,E} - \tilde{\boldsymbol{w}}^t)/\beta$              ▷ Update local dual variable
20:         Send back to the server the local updated model $\boldsymbol{w}_k^t = \boldsymbol{w}_{k,E}$
21:     **end for**
22:     $\sigma^{t+1} = \sigma^t - \frac{1}{\beta|\mathcal{C}|} \sum_{k \in \mathcal{C}} (\boldsymbol{w}_k^t - \boldsymbol{w}^t)$           ▷ Update global dual variable
23:     $\tilde{\Delta}_{\boldsymbol{w}}^t = \sum_{k \in \mathcal{C}^t} \frac{N_k}{N}(\tilde{\boldsymbol{w}}^t - \boldsymbol{w}_k^t)$              ▷ **Global pseudo-gradient**
24:     $\boldsymbol{w}^{t+1} = \boldsymbol{w}^t - \tilde{\Delta}_{\boldsymbol{w}}^t - \beta\sigma^{t+1}$    ▷ FEDGLOSS descent step w/ pseudo-grad using ADMM
25: **end for**

---

## B  THE BENEFITS OF FEDGLOSS

### B.1  ACHIEVING CONSISTENCY ON LOCAL AND GLOBAL SHARPNESS WITH FEDGLOSS

This section completes the analyses presented in Sections 4 and 6.2, showing how FEDGLOSS guides towards consistent local and global flat loss landscapes. Fig. 10 extends Figs. 5a and 5b from the main paper with the 3 remaining clients (out of 5) selected in the last round, comparing the behavior of clients' models trained with FEDGLOSS from the local and global perspectives. These results highlight the effectiveness of FEDGLOSS in achieving local models with consistent local-global behavior. Indeed, these models end up in flat regions both in local and global loss landscapes. The comparison with FEDSAM (*net* surface) further demonstrates the effectiveness of using our global flatness-aware approach.

(a) Local model trained
on class bear with FEDGLOSS

(b) Local model trained on
class skyscraper with FEDGLOSS

(c) Local model trained
on class possum with FEDGLOSS

Figure 10: **Global *vs.* local perspective on FEDGLOSS**. CIFAR100 $\alpha = 0$ with SAM as local optimizer @ $20k$ rounds on CNN. **(a) - (c):** Local models trained on one class, tested on the local ("*Local loss*") or global dataset ("*Global loss*"). Corresponding global perspective of local model trained with FEDSAM (*net*) added as reference.



(a) Local model trained
on class bear with FEDGLOSS

(b) Local model trained on
class skyscraper with FEDGLOSS

(c) Local model trained
on class possum with FEDGLOSS

Figure 11: **Global *vs.* local perspective on FEDSMOO**. CIFAR100 $\alpha = 0$ with SAM as local optimizer @ $20k$ rounds on CNN. **(a) - (c):** Local models trained on one class, tested on the local ("*Local loss*") or global dataset ("*Global loss*"). Corresponding global perspective of local model trained with FEDGLOSS (*net*) added as reference.

Fig. 11 instead extends Figs. 5c and 5d, offering a comparative analysis of FEDSMOO's clients' models w.r.t. FEDGLOSS (*net* landscape). The local solutions found by FEDGLOSS achieve flatter, better (*i.e.*, lower loss) and more consistent convergence points in the *global* loss landscape w.r.t. our main competitor FEDSMOO.

## B.2  ACHIEVING FLATTER GLOBAL MINIMA



(a) CIFAR10 $\alpha = 0$

(b) CIFAR10 $\alpha = 0.05$

(c) CIFAR100 $\alpha = 0$

(d) CIFAR100 $\alpha = 0.5$

Figure 12: Visualization of the loss landscapes of the CNN trained with **FEDGLOSS** (*net*) and the de-facto standard optimization algorithm in FL **FEDAVG** (*solid*). Comparison with varying degrees of heterogeneity on CIFAR10 (left) and CIFAR100 (right). **FEDGLOSS consistently achieves flatter minima and lower loss values in the *global* loss landscape.**

**FEDGLOSS *vs.* FEDAVG, FEDSAM, FEDSMOO.**   Fig. 12 compares the loss landscapes of global models trained with FEDGLOSS and FEDAVG, showing how the former consistently achieves flatter minima and lower loss values in the *global* loss landscapes. This confirms the behaviors already appreciated in Fig. 6. In addition, Fig. 13 shows the global loss surfaces of FEDGLOSS' solutions

(a) CIFAR10 $\alpha = 0.05$
FEDGLOSS *vs.* FEDSAM

(b) CIFAR10 $\alpha = 0.05$
FEDGLOSS *vs.* FEDSMOO

(c) CIFAR100 $\alpha = 0.5$
FEDGLOSS *vs.* FEDSAM

(d) CIFAR100 $\alpha = 0.5$
FEDGLOSS *vs.* FEDSMOO

Figure 13: Visualization of the loss landscapes of the CNN trained with **FEDGLOSS** (*net*) **and FEDSAM or FEDSMOO** (*solid*). Comparison with $\alpha = 0.05$ CIFAR10 (left) and $\alpha = 0.5$ CIFAR100 (right) extending Fig. 6. **FEDGLOSS consistently achieves flatter minima and lower loss values in the *global* loss landscape.**



(a) CIFAR10 $\alpha = 0.05$
FEDGLOSS (*net*) *vs.*
FEDAVG (*solid*)

(b) CIFAR10 $\alpha = 0.05$
FEDGLOSS (*net*) *vs.*
FEDSAM (*solid*)

(c) CIFAR10 $\alpha = 0.05$
FEDGLOSS (*net*) *vs.*
FEDSMOO (*solid*)

(d) CIFAR100 $\alpha = 0.5$
FEDGLOSS (*net*) *vs.*
FEDAVG (*solid*)

(e) CIFAR100 $\alpha = 0.5$
FEDGLOSS (*net*) *vs.*
FEDSAM (*solid*)

(f) CIFAR100 $\alpha = 0.5$
FEDGLOSS (*net*) *vs.*
FEDSMOO (*solid*)

Figure 14: Visualization of loss landscapes of the ResNet18 trained with **FEDGLOSS** (*net*) **and FEDAVG, or FEDSAM or FEDSMOO** (*solid*). Comparison with CIFAR10 $\alpha = 0.05$ (*a-c*) and CIFAR100 $\alpha = 0.5$ (*d-f*). **FEDGLOSS achieves flatter and lower-loss regions in the global landscape.**

against models trained with FEDSAM and FEDSMOO. These plots extend Fig. 6 with the less heterogeneous scenarios $\alpha = 0.05$ for CIFAR10 and $\alpha = 0.5$ on CIFAR100. They confirm FEDGLOSS' effectiveness in reaching flatter and lower-loss solutions with respect to its main direct competitors.

**FEDGLOSS on ResNet18.** Fig. 14 extends Fig. 6 from the main paper and shows the flatter loss landscapes reached by FEDGLOSS when using ResNet18 on CIFAR10 and CIFAR100.

**Hessian Eigenvalues.** Table 8 reports the values of the maximum Hessian eigenvalues, as already shown in Fig. 7 in the main paper. First, as expected, we note that SAM-based methods achieve flatter minima w.r.t. the counterpart. Notably, our main competitor FEDSMOO presents higher sharpness than FEDSAM in the simpler CIFAR10, regardless of the data distribution. In addition, FEDGLOSS with local SAM achieves the lowest sharpness (*i.e.*, lowest $\lambda_1$) on *all* configurations, outperforming the state of the art and, specifically, *all* sharpness-aware methods.

Table 8: **FEDGLOSS *vs.* the state of the art**, distinguished by local optimizer - SGD (*top*) and SAM (*bottom*). Comparison in terms of communication cost and maximum Hessian eigenvalue $\lambda_1$. Best results in **bold**. Model: CNN.

| | Method | Comm. Cost | CIFAR10 $\alpha = 0$ $\lambda_1(\downarrow)$ | CIFAR10 $\alpha = 0.05$ $\lambda_1(\downarrow)$ | CIFAR100 $\alpha = 0$ $\lambda_1(\downarrow)$ | CIFAR100 $\alpha = 0.5$ $\lambda_1(\downarrow)$ |
|---|---|---|---|---|---|---|
| **Client SGD** | FEDAVG | 1× | $66.23 \pm 0.50$ | $71.14 \pm 4.07$ | $\mathbf{66.30 \pm 3.08}$ | $68.77 \pm 0.96$ |
| | FEDPROX | 1× | $66.19 \pm 0.52$ | $71.41 \pm 4.40$ | $66.34 \pm 3.75$ | $\mathbf{68.63 \pm 1.37}$ |
| | FEDDYN | 1× | $63.94 \pm 4.41$ | $71.44 \pm 8.73$ | - | - |
| | SCAFFOLD | 2× | $166.54 \pm 6.93$ | $180.51 \pm 30.08$ | - | $120.01 \pm 0.76$ |
| | **FEDGLOSS** | 1× | $\mathbf{58.26 \pm 3.49}$ | $\mathbf{56.28 \pm 4.19}$ | $96.01 \pm 9.00$ | $107.35 \pm 7.5$ |
| **Client SAM** | FEDSAM | 1× | $10.35 \pm 0.07$ | $9.43 \pm 0.28$ | $58.38 \pm 2.93$ | $57.54 \pm 1.21$ |
| | FEDDYN | 1× | $10.04 \pm 5.38$ | $6.58 \pm 0.20$ | - | - |
| | FEDSPEED | 1× | $10.92 \pm 0.17$ | $9.97 \pm 0.12$ | $58.23 \pm 3.18$ | $58.00 \pm 1.86$ |
| | FEDGAMMA | 2× | $4.79 \pm 0.20$ | $4.55 \pm 0.20$ | - | $99.86 \pm 6.74$ |
| | FEDSMOO | 2× | $15.37 \pm 1.67$ | $12.57 \pm 0.56$ | $28.43 \pm 1.97$ | $29.23 \pm 0.17$ |
| | **FEDGLOSS** | 1× | $\mathbf{2.03 \pm 0.05}$ | $\mathbf{1.93 \pm 0.03}$ | $\mathbf{17.18 \pm 0.97}$ | $\mathbf{16.22 \pm 0.35}$ |

(a) $\alpha = 0$ SAM     (b) $\alpha = 0$ SGD     (c) $\alpha = 0.05$ SAM     (d) $\alpha = 0.05$ SGD

Figure 15: CIFAR10 with varying degrees of heterogeneity ($\alpha \in \{0, 0.05\}$). Results of centralized runs (*dashed lines*) added as reference. Comparison of FEDGLOSS with state-of-the-art approaches, distinguished in SAM-based methods (**a**, **c**) and SGD-based ones (**b**, **d**). FEDGLOSS consistently achieves the best performance. Model: CNN.



(a) $\alpha = 0$ SAM     (b) $\alpha = 0$ SGD     (c) $\alpha = 0.5$ SAM     (d) $\alpha = 0.5$ SGD

Figure 16: CIFAR100 with varying degrees of heterogeneity ($\alpha \in \{0, 0.5\}$) with **CNN**. Results of centralized runs (*dashed lines*) added as reference. Comparison of FEDGLOSS with state-of-the-art approaches, distinguished in SAM-based methods (**a**, **c**) and SGD-based ones (**b**, **d**). FEDGLOSS consistently achieves the best performance.



(a) CIFAR100 $\alpha = 0.5$ SAM    (b) CIFAR100 $\alpha = 0.5$ SGD    (c) CIFAR10 $\alpha = 0.05$ SAM    (d) CIFAR10 $\alpha = 0.05$ SGD

Figure 17: Accuracy trends with **ResNet18** on CIFAR100 (*left*) and CIFAR10 (*right*). Comparison of FEDGLOSS with state-of-the-art approaches, distinguished in SAM-based methods (**a**, **c**) and SGD-based ones (**b**, **d**). FEDGLOSS consistently achieves the best performance, both in terms of final accuracy and convergence speed.

## B.3   INCREASING CONVERGENCE SPEED

Figs. 15 and 16 show the accuracy trends of FEDGLOSS compared to state-of-the-art methods for heterogeneous FL on CIFAR10 and CIFAR100 respectively. For a clearer understanding, we distinguish between SAM-based and SGD-based methods depending on the used local optimizer. For a fair comparison, we report FEDSMOO with and without the scheduling of $\rho$ (+*wp* in the figure). For additional details on the scheduling, refer to Appendix C. FEDGLOSS consistently achieves the best performances and convergence speedup in each group. In addition, we remind that FEDGLOSS communicates *half* the number of bits at each round w.r.t. FEDSMOO. Fig. 17 reports the results obtained with ResNet18 on both datasets: FEDGLOSS consistently achieves the best speed of convergence and final accuracy, with both SAM and SGD as CLIENTOPT.

## B.4   EFFECTIVENESS IN MULTIPLE SCENARIOS AND WITH SEVERAL MODEL ARCHITECTURES

As already shown in Section 6, FEDGLOSS outperforms the state of the art in terms of speed of convergence, final performance in accuracy, flatness of the solution and communication efficiency. This remains true across multiple datasets (CIFAR100, CIFAR10, LANDMARKS-USER-160K) and

Table 9: FEDGLOSS against SOTA FL methods on homogeneous CIFAR settings, compared in terms of communication costs, accuracy (%) and maximum Hessian eigenvalue $\lambda_1$. Best result in **bold** and second best underlined. Model: CNN.

| Method | Comm. Cost | ADMM | CIFAR10 $\alpha = 100$ Accuracy | $\lambda_1$ | CIFAR100 $\alpha = 1000$ Accuracy | $\lambda_1$ |
|---|---|---|---|---|---|---|
| FEDAVG | 1× | ✗ | 84.0 | 68.4 | 50.1 | 49.4 |
| FEDSAM | 1× | ✗ | 84.7 | 36.2 | 53.4 | 32.6 |
| FEDDYN (SGD) | 1× | ✓ | 83.8 | 47.8 | 51.9 | 91.7 |
| FEDDYN (SAM) | 1× | ✓ | 84.5 | 27.9 | 52.5 | 46.0 |
| FEDSMOO | 2× | ✓ | **85.1** | <u>6.4</u> | 53.9 | 24.6 |
| FEDGLOSS (SGD) | 1× | ✗ | 84.0 | 67.7 | 50.5 | 50.8 |
|  | 1× | ✓ | 83.1 | 7.1 | 51.7 | 47.9 |
| FEDGLOSS (SAM) | 1× | ✗ | <u>84.8</u> | 36.2 | **55.8** | <u>13.9</u> |
|  | 1× | ✓ | <u>84.8</u> | **2.8** | <u>55.7</u> | **11.8** |

model architectures (CNN, ResNet18, MobileNetv2). The main paper focuses on results in more challenging heterogeneous FL scenarios. The next paragraph discusses the behavior of FEDGLOSS in homogeneous settings.

**FEDGLOSS in Homogeneous Settings.** Table 9 evaluates FEDGLOSS against the main methods FEDAVG, FEDSAM, FEDDYN and FEDSMOO in homogeneous FL settings. Here, client gradients are naturally more aligned due to reduced client drift (Karimireddy et al., 2020b). We thus test FEDGLOSS with and without ADMM for global consistency. As expected, FEDGLOSS achieves similar accuracy with or without ADMM, particularly when using SAM as the local optimizer. However, ADMM facilitates convergence to flatter minima (evidenced by lower $\lambda_1$ values) by aligning local and global convergence points. Notably, FEDGLOSS achieves the flattest minima (lowest $\lambda_1$) across both datasets, and the best accuracy on the more complex CIFAR100. While FEDSMOO achieves slightly higher accuracy on CIFAR10, FEDGLOSS finds a flatter minimum and achieves competitive accuracy with significantly lower communication costs (halved).

## B.5 REDUCING COMMUNICATION COSTS

**Analysis on communication cost.** Table 10 extends the analysis presented in Section 6.3.4 by evaluating the communication costs across all scenarios considered in this work. Notably, since FEDGLOSS maintains the per-round communication cost of FEDAVG, it remains advantageous even when matching the convergence speed of the best-performing algorithm (FEDSMOO), due to its *halved communication costs*.

Table 10: Communication costs comparison w.r.t. FEDAVG. "-" for not reached accuracy, "✗" for non-convergence.

| Method | CNN Cifar-10 $\alpha = 0$ Rounds | Cost | Cifar-10 $\alpha = 0.05$ Rounds | Cost | Cifar-100 $\alpha = 0$ Rounds | Cost | Cifar-100 $\alpha = 0.5$ Rounds | Cost | ResNet18 Cifar-10 $\alpha = 0.05$ Rounds | Cost | ResNet18 Cifar-100 $\alpha = 0.5$ Rounds | Cost | MobileNetv2 LANDMARKS-USER-160K Rounds | Cost |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **SGD-based** | | | | | | | | | | | | | | |
| FEDAVG | 10k (1×) | 1B | 10k (1×) | 1B | 20k (1×) | 1B | 20k (1×) | 1B | 10k (1×) | 1B | 10k (1×) | 1B | 1.3k (1×) | 1B |
| FEDPROX | 7.6k (1.3×) | 0.8B | 7.9k (1.3×) | 0.8B | 18.7k (1.1×) | 0.9B | 18.6k (1.1×) | 0.9B | 8.8k (1.1×) | 0.9B | 8.3k (1.2×) | 0.8B | - | - |
| FEDDYN | 2k (5×) | **0.2B** | 1.9k (5×) | **0.2B** | ✗ | | ✗ | | - | - | 3.5k (2.9×) | 0.4B | - | - |
| SCAFFOLD | - | | - | | - | | - | | - | | 8.9k (1.1×) | 1.8B | ✗ | |
| FEDGLOSS | 3.4k (2.9×) | 0.3B | 3.8k (2.6×) | 0.4B | 5k (4×) | **0.3B** | 4.7k (4.3×) | **0.2B** | 2.4k (4.2×) | **0.2B** | 1.9k (5.3×) | **0.2B** | - | - |
| **SAM-based** | | | | | | | | | | | | | | |
| FEDSAM | 6.3k (1.6×) | 0.6B | 7.8k (1.3×) | 0.8B | 18.3k (1.1×) | 0.9B | 16.3k (1.2×) | 0.8B | 9.2k (1.1×) | 0.9B | 7.8k (1.3×) | 0.8B | 1.3k (1×) | 1B |
| FEDDYN | 3k (3.3×) | 0.3B | 4.2k (2.4×) | 0.4B | ✗ | | ✗ | | 4.1k (2.4×) | 0.4B | 3.5k (2.9×) | 0.4B | - | - |
| FEDSPEED | 6.3k (1.6×) | 0.6B | 6.9k (1.4×) | 0.7B | 18.3k (1.1×) | 0.9B | 15.7k (1.3×) | 0.8B | 9.3k (1.1×) | 1.9B | 8.1k (1.2×) | 1.6B | 1.3k (1×) | 1B |
| FEDGAMMA | - | | - | | - | | - | | 9.3k (1.1×) | 1.9B | 8.1k (1.2×) | 1.6B | ✗ | |
| FEDSMOO | 1.9k (5.3×) | 0.4B | 2.2k (4.5×) | 0.4B | 4.5k (4.4×) | 0.5B | 6.5k (3.1×) | 0.7B | 2.4k (4.2×) | 0.5B | 2.3k (4.3×) | 0.5B | 200 (6.5×) | 0.3B |
| FEDGLOSS | 2.2k (4.5×) | **0.2B** | 2.2k (4.5×) | **0.2B** | 6.3k (3.2×) | **0.3B** | 5.2k (3.8×) | **0.3B** | 2.4k (4.2×) | **0.2B** | 1.9k (5.3×) | **0.2B** | 200 (6.5×) | **0.2B** |

**FEDGLOSS vs. NAIVEFEDGLOSS.** Fig. 18 deepens the comparison between FEDGLOSS and its baseline NAIVEFEDGLOSS, discussed in Section 6.4. In particular, this plot reports the accuracy trends of the two methods, showing that NAIVEFEDGLOSS is slightly faster ($\approx 1.1\times$) than FEDGLOSS in CIFAR100, while FEDGLOSS surpasses the speed of the baseline after $\approx 25\%$ of training in CIFAR10. However, both methods reach the same accuracy at the of training. In addition, it is to be noted that FEDGLOSS *halves* the communication cost w.r.t. NAIVEFEDGLOSS by transmitting half the number of bits at each round. Additionally, it also reduces the communication cost by half,

Figure 18: **Accuracy trends of FEDGLOSS *vs.* NAIVEFEDGLOSS.** The comparison includes the centralized upper bounds of SAM and $\widetilde{\text{SAM}}$ (the adaptation of FEDGLOSS' strategy to the centralized scenario). CNN on CIFAR10 and CIFAR100 with varying heterogeneity degree ($\alpha$). NAIVEFEDGLOSS is $\approx 1.1\times$ faster than its efficient alternative FEDGLOSS in CIFAR100, while FEDGLOSS shows increased convergence speed after $\approx 25\%$ of training rounds in CIFAR10. However, both methods reach the same accuracy at the of training. These results motivate the choice of FEDGLOSS over NAIVEFEDGLOSS.

as it eliminates the need to invoke the clients twice to compute the updates. These insights further support our choice of the efficient strategy of FEDGLOSS over NAIVEFEDGLOSS.

**Convergence speed.** Since all methods are compared over a fixed amount of communication rounds, **higher final accuracy implies faster convergence**. Since FEDGLOSS consistently outperforms the other state-of-the-art algorithms taken into account, it is guaranteed to converge faster, as also shown in the accuracy trends (Figs. 15 to 17).

## B.6 THE IMPACT OF SERVER-SIDE $\rho_s$

Fig. 19 analyzes the impact of $\rho_s$ on the performance of the global model, both in terms of accuracy (Fig. 19a) and flatness of the solution (Fig. 19b). In details, Fig. 19a compares the accuracy of the global model trained on CIFAR100 when varying the model architecture (CNN *vs.* ResNet18) and the data heterogeneity ($\alpha = 0$ *vs.* $\alpha = 0.5$). In all the configurations, we note that a smaller value of $\rho_s$ usually leads to the best results. Fig. 19b instead focuses on the CNN in the most heterogeneous setting ($\alpha = 0$) and compares the reached accuracy with the corresponding maximum Hessian eigenvalue $\lambda_1$ when varying $\rho_s$. A larger server-side $\rho_s$ corresponds to a smaller $\lambda_1$, *i.e.*, a flatter region in the global loss landscape.



(a) FEDGLOSS $\rho_s$ with different architectures



(b) FEDGLOSS $\rho_s$ *vs.* $\lambda_1$

Figure 19: CIFAR100. **(a):** Accuracy (%) of FEDGLOSS when varying the server-side SAM $\rho_s$, with different heterogeneity and architecture (CNN with $\alpha = 0$ and ResNet18 with $\alpha = 0.5$). Smaller values of $\rho_s$ lead to better performances. **(b):** FEDGLOSS $\rho_s$ *vs.* maximum Hessian eigenvalue $\lambda_1$ on CNN with $\alpha = 0$. Larger values of $\rho_s$ lead to lower eigenvalues with minimum loss in accuracy.

Figure 20: CIFAR10 (*left*) and CIFAR100 (*right*) data distribution across clients with the heterogeneity levels tested in the experiments. On top of each chart we report the average number of classes seen by each client.

## C IMPLEMENTATION DETAILS

This section delves into a comprehensive description of the datasets and models utilized throughout this paper, specifying the Deep Learning framework and the hardware employed (Appendix C.1). Additionally, we present the area of the hyper-parameters' space explored during the fine-tuning process in order to yield optimal results (Appendix C.2).

### C.1 DATASETS AND MODELS

Table 11 provides a comprehensive overview of each dataset's general statistics. This includes the number of training clients participating in the process and the total number of samples used to construct both the training and test sets.

#### C.1.1 CIFAR10 AND CIFAR100

We adapted these two well-known image classification datasets to the FL scenario by replicating the splits among clients proposed by (Hsu et al., 2019). Both datasets are split evenly among 100 clients, thus each of them has access to 500 data samples. This partitioning is performed according to a Latent Dirichlet Allocation (LDA) on the labels. In practice, each local dataset follows a multinomial distribution drawn according to a symmetric Dirichlet distribution with concentration parameter $\alpha$. The higher the value of this parameter is, the more the local datasets resemble a homogeneous scenario, in the limit case $\alpha = 0$ each client has access to one only class of images. In our experiments we tested $\alpha \in \{0, 0.05, 100\}$ and $\alpha \in \{0, 0.5, 1000\}$ for CIFAR10 and CIFAR100, respectively. Figs. 20a and 20b show how data is distributed across clients in all the experimental settings for these two datasets. Both datasets are pre-processed by applying random crops and random horizontal flips.

**Models.** We trained a Convolutional Neural Network (CNN) inspired by the LeNet-5 architecture (LeCun et al., 1998), as proposed by (Hsu et al., 2020). The network comprises two 64-channels convolutional layers, both using $5 \times 5$ kernels and followed by $2 \times 2$ max-pooling layers. This is succeeded by two fully-connected layers with 384 and 192 units, respectively. The final output layer is adapted to the specific number of classes in the dataset.

To explore deeper and more expressive architectures, we also trained a ResNet18 (He et al., 2015) on CIFAR100 with $\alpha = 0.5$ and CIFAR10 with $\alpha = 0.05$. We replaced the standard Batch Normalization layers (Ioffe & Szegedy, 2015) with Group Normalization layers (Wu & He, 2018) due to their demonstrated effectiveness in handling skewed data distributions in FL (Hsieh et al., 2020). We carried out the experiments with the CNN model using PyTorch (Paszke et al., 2019) and the ResNet18 ones using FedJAX (Ro et al., 2021).

Table 11: Datasets' description with their general statistics on the size and number of clients.

| Dataset | Train clients | Train samples | Test samples |
|---|---|---|---|
| CIFAR10 | 100 | 50,000 | 10,000 |
| CIFAR100 | 100 | 50,000 | 10,000 |
| LANDMARKS-USER-160K | 1262 | 164,172 | 19,526 |

Table 12: General training hyper-parameters common to all methods, distinguished by dataset and model architecture. Symbols: local epochs $E$, local learning rate $\eta$, weight decay $wd$, client-side momentum $\beta_l$, batch size $B$.

| Dataset | Model | Rounds | Clients per round | E | $\eta$ | $wd$ | $\beta_l$ | $B$ |
|---|---|---|---|---|---|---|---|---|
| CIFAR10 | CNN | 10000 | 5 | 1 | $10^{-2}$ | $4 \cdot 10^{-4}$ | 0 | 64 |
| CIFAR100 | CNN | 20000 | 5 | 1 | $10^{-2}$ | $4 \cdot 10^{-4}$ | 0 | 64 |
| | ResNet18-GN | 10000 | 10 | 1 | $10^{-2}$ | $10^{-5}$ | 0.7 | 64 |
| LANDMARKS-USER-160K | MobileNetv2 | 1300 | 50 | 5 | 0.1 | $4 \cdot 10^{-5}$ | 0 | 64 |

Table 13: Search grid used to find optimal hyper-parameters for each combination of method, dataset and model. We highlight the best performing values in **bold**.

| Method | HParam | CIFAR10 CNN | CIFAR10 ResNet18 | CIFAR100 CNN | CIFAR100 ResNet18 | LANDMARKS-USER-160K |
|---|---|---|---|---|---|---|
| FEDSAM | $\rho$ | [0.05, 0.1, **0.15**, 0.2] | [**0.01**, 0.02, 0.05] | [0.005, **0.01**, 0.02, 0.05] | [**0.01**, 0.02, 0.05] | [**0.05**] |
| FEDPROX | $\mu$ | [0.001, 0.01, **0.1**] | [0.001, **0.01**, 0.1] | [0.001, 0.01, **0.1**] | [0.001, 0.01, **0.1**] | [0.001, 0.01, **0.1**] |
| FEDDYN | $\alpha$ | [0.001, **0.01**, 0.1] | [0.001, **0.01**, 0.1] | [0.001, **0.01**, 0.1] | [**0.01**] | [**0.001**, 0.01] |
| | $\rho$ (SAM-based only) | [**0.15**] | [**0.01**] | [0.01, **0.02**] | [**0.01**] | [**0.05**] |
| FEDSPEED | $\rho$ | [0.05, 0.1, **0.15**, 0.2] | [**0.01**] | [0.005, **0.01**, 0.02, 0.05] | [**0.01**] | [**0.05**] |
| | $\alpha$ | [0.9, **0.95**, 0.99] | [0.9, **0.95**] | [0.9, **0.95**, 0.99] | [0.9, **0.95**] | [**0.95**] |
| | $\lambda$ | [10, **100**, 1000] | [10, 100, **1000**] | [10, 100, **1000**] | [10, 100, **1000**] | [100, **1000**] |
| FEDGAMMA | $\rho$ | [**0.15**] | [**0.01**] | [**0.01**] | [**0.01**] | [**0.05**] |
| FEDSMOO | $\rho$ | [0.05, 0.1, **0.15**, 0.2] | [**0.01**] | [0.005, 0.01, 0.05, **0.1**, 0.2] | [**0.01**] | [0.05, **0.1**, 0.2, 0.3] |
| | $\beta$ | [5, **10**, 100] | [1, 2, **5**, 10] | [**10**, 100] | [5, **10**, 100] | [10, **50**, 100, 1000] |
| FEDGLOSS (ours) | $\rho_s$ | [0.01, 0.1, **0.15**] | [**0.01**, 0.05, 0.1, 0.5] | [**0.01**, 0.05, 0.1, 0.2] | [**0.01**, 0.05, 0.1, 0.5] | [**0.005**, 0.01, 0.02] |
| | $\rho$ | [0.05, 0.1, **0.15**, 0.2] | [**0.01**] | [0.05, 0.1, **0.2**] | [**0.01**] | [0.05, 0.1, 0.2, **0.3**] |
| | $\beta$ | [5, **10**, 100] | [**1**, 2, 5, 10] | [**10**, 100] | [5, **10**, 100] | [10, **50**, 100] |
| | $T_s$ | [1000, **2000**, 4000] | [**0**] | [1000, 2000, 5000, 10000, **15000**] | [**0**] | [**0**] |

### C.1.2 LANDMARKS-USER-160K

To achieve a comprehensive understanding of the efficacy of the proposed method, a thorough analysis was undertaken utilizing large-scale real-world datasets. Specifically, we use the LANDMARKS-USER-160K dataset (Hsu et al., 2020), a repository encompassing 164,172 images depicting 2028 distinct landmarks, distributed among 1262 clients.

**Models.** The model employed for training is a MobileNetV2 (Sandler et al., 2018; Hsu et al., 2020), replacing batch normalization with group normalization layers and pre-trained on ImageNet (Deng et al., 2009). The tests on this dataset were carried out on a cluster of NVIDIA A100 40GB, using our FedJAX codebase.

### C.2 HYPERPARAMETERS

In Table 12 we report the training hyperparameters associated to each dataset and model pairing. While Table 13 summarizes the hyper-parameters search grid tested for each method (in bold the chosen ones). All runs are averaged over 3 seeds. In addition, following previous works (Acar et al., 2021; Caldarola et al., 2022; 2023), we report the averaged accuracy of the last 100 rounds to reduce the noise typical of heterogeneous FL settings.

While running our experiments on FEDGLOSS, we noticed that a larger value of local $\rho$ allowed to reach the best final accuracy, while a smaller $\rho$ achieved faster convergence in the initial training stages. Following this insight, we schedule the value of $\rho$ for $T_s$ rounds as

$$\rho(t) = \begin{cases} \rho_0 + \frac{(\rho - \rho_0)}{T_s} \cdot t & \text{if } t \leq T_s \\ \rho & \text{otherwise,} \end{cases}$$

starting from the value $\rho_0 = 0.001$.

## D FLATNESS ANALYSIS

This section describes the procedure to compute the visualization of the loss landscapes and the Hessian eigenvalues.

## D.1   VISUALIZING 1D LOSS LANDSCAPES

Fig. 9 in the main paper shows the interpolation of FEDGLOSS and NAIVEFEDGLOSS's models. Given their respective weights $\boldsymbol{w}_{\text{FEDGLOSS}}$ and $\boldsymbol{w}_{\text{NAIVEFEDGLOSS}}$, the interpolation is computed using a coefficient $\gamma$ as

$$\boldsymbol{w} = \gamma \cdot \boldsymbol{w}_{\text{FEDGLOSS}} + (1 - \gamma) \cdot \boldsymbol{w}_{\text{NAIVEFEDGLOSS}} . \tag{8}$$

For each $\gamma \in [-1, 2]$, $\boldsymbol{w}$ is tested on the training or test sets, and the plot reports the computed loss. The resulting interpolation indicates that there is no loss barrier between the two models, suggesting they lie within the same basin. Additionally, the 1D geometry of the emerging loss landscape reveals that both models converge to a flat minimum when evaluated on both CIFAR100 and CIFAR10.

## D.2   VISUALIZING 3D LOSS LANDSCAPES

We leverage techniques from (Li et al., 2018) to visualize the loss landscapes of our models. We adapt their code to work with our specific datasets and network architectures. The process involves calculating the loss function along random directions in the parameter space. This is achieved by perturbing the model's parameters within a defined range. In our visualizations, we constrain these perturbations to occur within the range of $[-1, 1]$ for both the $x$ and $y$ axes. To ensure consistent comparisons across models (*e.g.*, as seen in Fig. 6), we utilize the same set of random directions for all models.

## D.3   HESSIAN EIGENVALUES FOR FLATNESS MEASURE

Following prior work (Foret et al., 2021; Garipov et al., 2018; Li et al., 2018; Caldarola et al., 2022), we use the spectrum of the Hessian matrix to quantify the *flatness* of the loss landscape. Here, lower maximum eigenvalues correspond to flatter landscapes, implying less sharpness. To compute these eigenvalues (denoted by $\lambda_1$ in the main paper), we employ the stochastic power iteration method (Xu et al., 2018) with a maximum of 20 iterations, referring to the code of Golmant et al. (2018).