
Rethinking the Flat Minima Searching in Federated Learning

Taehwan Lee¹ Sung Whan Yoon^{1 2}

Abstract

Albeit the success of federated learning (FL) in decentralized training, bolstering the generalization of models by overcoming heterogeneity across clients still remains a huge challenge. To aim at improved generalization of FL, a group of recent works pursues flatter minima of models by employing sharpness-aware minimization in the local training at the client side. However, we observe that the global model, i.e., the aggregated model, does not lie on flat minima of the global objective, even with the effort of flatness searching in local training, which we define as *flatness discrepancy*. By rethinking and theoretically analyzing flatness searching in FL through the lens of the discrepancy problem, we propose a method called Federated Learning for Global Flatness (FedGF) that explicitly pursues the flatter minima of the global models, leading to the relieved flatness discrepancy and remarkable performance gains in the heterogeneous FL benchmarks.

1. Introduction

Federated Learning (FL) has drawn great attention as a key framework for enabling decentralized learning across an immense number of distributed clients while preserving data privacy. The essence of FL is keeping local data on the client side and communicating the gradients or model parameters between a server and clients, where the server’s direct access to the local data is prohibited (McMahan et al., 2017). Nonetheless, there still exist daunting challenges that remain unsolved. Most importantly, the diversity or heterogeneity of data distribution across clients is shown to hinder the successful aggregation of global model parameters, leading to deteriorated performance and inhibiting model

convergence (Li et al., 2020b). To overcome the problem, researchers have dedicated to developing FL algorithms that achieve successful aggregation across heterogeneous clients (Li et al., 2020a; Karimireddy et al., 2020; Acar et al., 2021). Although such intensive efforts have been made to break the hurdle of heterogeneity, the agreed rule of thumb methods and principles have not yet been established.

One of the intriguing recent approaches to enhance the generalization is employing particular optimization methods that find flatter minima on loss surface, which is widely observed to enhance the generalization of deep models against the data distribution shifts (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Izmailov et al., 2018; Foret et al., 2021; Cha et al., 2021). The most popular flatness searching method is Sharpness-Aware Minimization (SAM), which incorporates the flatness around the minimum into the cost function (Foret et al., 2021). Building on the promising optimizer, researchers in the FL field have recently confirmed the effectiveness of SAM in strengthening the performance of FL algorithms for heterogeneous settings (Qu et al., 2022; Caldarola et al., 2022). Their approaches basically employ a SAM or SAM-variant optimizer at the local training step at each client to find a flatter local model of the local objective, which indeed yields considerable performance gains for the aggregated global model.

Let us then raise a pivotal question to rethink the flatness searching in FL: “Does the flatness searching in local training truly imply the flatness of the global model for the global objective?” The answer is “Not for the heterogeneous FL cases”. When the heterogeneity across clients becomes severe, we observed that the existing FL methods with the flat minima searching look effective for finding flatter minima in local training, but the global model does not lie on flatter minima for the global objective. One of the recent works initiated a discussion on this issue, but little intention was paid to elaborate on the formal understanding of it (Sun et al., 2023a). The issue is raised particularly in the regime of decentralized training, and it severely degrades the performance of the flatness searching FL methods. We formally define the issue as *flatness discrepancy*.

In this paper, we empirically and theoretically analyze the relationship between heterogeneity across clients and flatness discrepancy: A strong heterogeneity leads to severe

¹Graduate School of Artificial Intelligence, Ulsan National Institute of Science and Technology (UNIST), Ulsan, South Korea ²Department of Electrical Engineering, UNIST, Ulsan, South Korea. Correspondence to: Sung Whan Yoon <shyoon8@unist.ac.kr>.

discrepancy, eventually yielding the degraded performance of the global model. Based on the findings, we propose a method called Federated Learning for Global Flatness (FedGF) that relieves flatness discrepancy, leading to flatter minima of the global model. The key of FedGF is to explicitly consider the sharpness of the global model when running SAM in the local training. Specifically, we utilize the interpolated perturbation for SAM in both views of local and global objectives. We empirically confirm that our method shows remarkable performance gains over prior flatness searching FL methods, ranging up to +5.09% and +10.02% gains in the heterogeneous CIFAR-10 and CIFAR-100 benchmarks, respectively. Also, FedGF shows significantly faster convergence in heterogeneous cases, which is theoretically guaranteed by showing how FedGF suppresses the heterogeneity-related factor in the convergence analysis.

2. Preliminaries and Motivations

We here present the preliminaries for the baseline FL framework called FedAvg (McMahan et al., 2017) and a popular flatness searching FL algorithm, i.e., FedSAM (Qu et al., 2022; Caldarola et al., 2022). Also, we here define the flatness discrepancy, which is the key motivation of our work, and empirically show how it appears to FedSAM.

2.1. Preliminaries of FL: FedAvg

Notations: In the FL setting with N clients and a server, each client contains m_i local data samples, where $i \in [N]$ is the index of the client. $[N]$ indicates the set of integers ranging from 1 to N . A data sample $z_{i,j} = (x_{i,j}, y_{i,j})$ is j -th sample of i -th client with paired input $x_{i,j}$ and its label $y_{i,j}$, and it is drawn from the local data distribution \mathcal{D}_i .

The local objective function of client i is $F_i(w) := \frac{1}{m_i} \sum_{j=1}^{m_i} l(w, z_{i,j})$, where w is the model weight and $l(\cdot, \cdot)$ is the loss function. FL basically aims to find the model weight w^* that minimizes the global objective $F(w)$, i.e.,

$$w^* = \operatorname{argmin}_w \left\{ F(w) := \sum_{i=1}^N \frac{m_i}{m} F_i(w) \right\}, \quad (1)$$

where m is the total number of data samples across clients.

The way to optimize the model weight without accessing the samples on the client side is to adopt a repetitive aggregation process called round, where a round consists of local training of models at each client and aggregation of the locally-trained models at the server.

Local training: At round $r \in [R]$, each client receives the aggregated model w^r from the server and runs local training with K epochs. Specifically, local training is done with empirical risk minimization of the local loss:

$$w_{i,k+1}^r = w_{i,k}^r - \eta_l \nabla F_i(w_{i,k-1}^r), \quad (2)$$

where k is the number of local epochs, η_l is the learning rate and $w_{i,0}^r = w^r$. After K epochs, client i obtains $w_{i,K}^r$.

Aggregation: The updated local models are then uploaded to the server and aggregated to obtain global model w^{r+1} :

$$w^{r+1} = \sum_{i \in \mathcal{S}^r} \frac{m_i}{m} w_{i,K}^r, \quad (3)$$

where $\mathcal{S}^r \subseteq [N]$ is the index of participating clients for round r . After the aggregation, the next local training for round $r+1$ follows by broadcasting the global model to the clients. With a sufficient number of rounds up to R , the global model converges to the optimal weight in Eq. (1).

2.2. Preliminaries of Flatness Searching in FL: FedSAM

FedSAM adopts the SAM optimizer (Foret et al., 2021) for flatness searching in the local training of FedAvg (Qu et al., 2022; Caldarola et al., 2022).

SAM optimizer: The SAM optimizer transforms a loss function $f(w)$ into a min-max cost function as follows:

$$\min_w \max_{\|\delta\| \leq \rho} F(w + \delta), \quad (4)$$

where ρ is a positive real number and $\|\delta\|$ is L2-norm of δ . As a key factor, δ works as the perturbation that maximally raises the loss value so that the SAM optimizer can find flat minima. The perturbation can be simply approximated as the gradient direction, which points to the steepest direction of the loss surface.

FedSAM: By adopting the min-max problem of Eq. (4) in local training, FedSAM perturbs local model $w_{i,k}^r$:

$$\tilde{w}_{i,k}^r = w_{i,k}^r + \delta = w_{i,k}^r + \rho g_{i,k}^r / \|g_{i,k}^r\| \quad (5)$$

$$w_{i,k+1}^r = w_{i,k}^r - \eta_l \tilde{g}_{i,k}^r, \quad (6)$$

where $g_{i,k}^r = \nabla F_i(w_{i,k}^r)$ is the gradient computed at $w_{i,k}^r$, $\tilde{w}_{i,k}^r$ is the perturbed model weight, and $\tilde{g}_{i,k}^r = \nabla F_i(\tilde{w}_{i,k}^r)$ is the gradient computed at the perturbed model. Thus, FedSAM finds the local model with flatter minima, leading to the improved performance of the aggregated global model.

2.3. Motivations: Flatness Discrepancy

Now, we are ready to discuss the flatness discrepancy issue: *Flatness searching in local training does not imply the flatness of the global model*. First, we formally define the discrepancy, i.e., $\Delta_{\mathcal{F}}$, as the difference gap of the flatness between the global and local models. For simplicity, we here drop the notations of round r and local epoch k .

Definition 1. Flatness discrepancy $\Delta_{\mathcal{F}}$ of the global model

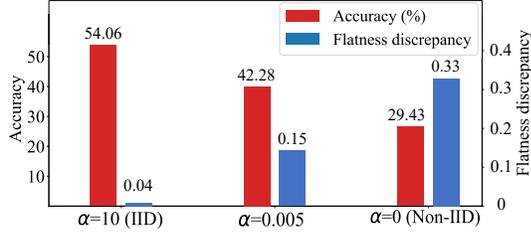


Figure 1: The performance and the flatness discrepancy ($\Delta_{\mathcal{F}}$) of FedSAM for the CIFAR-100 experiment

w and the local models $\{w_i\}_{i=1}^N$ is defined as:

$$\Delta_{\mathcal{F}} := \left| \max_{\|\delta\| \leq \rho} F(w + \delta) - F(w) - \left[\sum_{i=1}^N \frac{m_i}{m} \max_{\|\delta_i\| \leq \rho} F_i(w_i + \delta_i) - F_i(w_i) \right] \right|. \quad (7)$$

When $\Delta_{\mathcal{F}}$ is small, it means that the increasing amount of losses of the global and local objectives are similar, i.e., the degree of the flatness is similar. When $\Delta_{\mathcal{F}}$ is large, the increasing amount of losses of the global and local objectives are not the same, i.e., the flatness is discrepant. Herein, we want to provide the motivating empirical results to show how the discrepancy issue arises in flatness searching FL.

As a preliminary experiment, we compute the flatness discrepancy values of FedSAM for the CIFAR-100 FL benchmark. As shown in Fig. 1, FedSAM shows significant performance degradation as the heterogeneity increases (when α decreases to 0, the data distribution across clients becomes heterogeneous, i.e., non-IID (non-independently and identically distributed)). Interestingly, along with the performance degradation, we observe that the flatness discrepancy increases when the heterogeneity gets worse. It empirically reveals that a naive application of SAM to local training of FL does not guarantee the flatness of the global model for the global objective. Also, we visualize how the discrepancy appears on the loss surface. As presented in Fig. 2, FedSAM shows flatter minima around the local model, but the global model does not lie on flatter minima.

We provide an intuition of how the heterogeneity causes the larger flatness discrepancy. With the strong heterogeneity, the local and global objectives, i.e., F_i and F , respectively, trivially deviate from each other due to the data distribution gap, so the flatness searching in the local training does not imply the flatness of the global model. On the other hand, if the heterogeneity is not severe, i.e., close to the IID case, the local objective is ideally the same as the global objective, so the flatness searching for local training directly links to the global model with flatter minima of global objective.

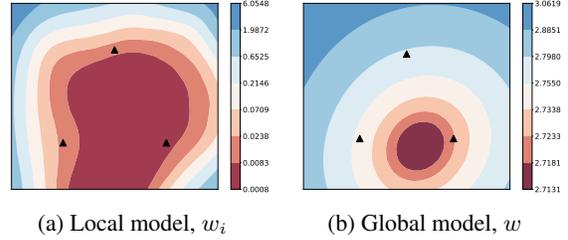


Figure 2: Visualization of the loss surface of FedSAM for the CIFAR-100 case ($\alpha = 0$).

3. Related Work

3.1. Heterogeneity Issues in FL

In past years, various strategies have been proposed to solve the heterogeneity issues in FL. A main branch of prior approaches focuses on regularizing local training to alleviate the divergence of the local models. As early works, FedAvgM (Hsu et al., 2019) and SCAFFOLD (Karimireddy et al., 2020) utilize the update of the global model as momentum in global and local training, respectively, to regularize the divergence of local gradients. FedProx (Li et al., 2020a) adopts a regularization term, which is called the proximal term, to prevent local models from largely deviating from the global model. FedDyn (Acar et al., 2021) utilizes the dynamic regularization term, which is tailored to each client, so suppressing the discrepancy between the global and local models. FLIX (Gasnov et al., 2021) utilizes the interpolation between global and local model parameters. Our method, FedGF, is based on flat minima searching, which is clearly different from the regularization-based methods.

3.2. Flat Minima Searching

In centralized learning: From an early finding of the benefits of flat minima over sharp minima of the model parameters on loss surface (Keskar et al., 2017), the potential of flat minima for enhancing the generalization ability of deep models is largely investigated in both empirical and theoretical ways. A group of works with Stochastic Weight Averaging (SWA) has been suggested as a simple heuristic method for finding flatter minima (Izmailov et al., 2018; Cha et al., 2021). As another branch of tools, Sharpness-Aware Minimization (SAM) embeds the flatness term into the cost function of the optimizer for seeking flatter minima (Foret et al., 2021; Kwon et al., 2021).

Correlation to the generalization: Some researchers have raised doubts about the correlation between the generalization and the flatness. Sharp minima are shown to be able to generalize (Dinh et al., 2017), and large models, e.g., transformers, empirically seem not to well correlate its flatness to the generalization capability (Andriushchenko et al., 2023;

Kim et al., 2023). Also, a flatness metric should be chosen carefully to show the positive correlation to the generalization (Bisla et al., 2022). The recent active debate argues that flatter minima do not directly imply better generalization, but it does not contradict the shown advantages of flat minima searching in standard training (Izmailov et al., 2018; Foret et al., 2021; Kwon et al., 2021), FL (Qu et al., 2022; Caldarola et al., 2022; Dai et al., 2023; Sun et al., 2023a;b), and out-of-distribution generalization (Cha et al., 2021).

In the FL setting: The existing FL methods that pursue flatter minima are based on applying SAM or SAM-variants to FL, e.g., FedSAM, FedASAM, MoFedSAM, FedGAMMA, FedSMOO, and FedSpeed (Caldarola et al., 2022; Qu et al., 2022; Dai et al., 2023; Sun et al., 2023a;b). FedSAM simply applies the SAM optimizer to local training. FedASAM controls the size of the perturbation along with the magnitude of model parameters. On the other hand, MoFedSAM utilizes the momentum, which is the update of the global model, when running SAM on the client side. Its strategy is analogous to that of SCAFFOLD (Karimireddy et al., 2020), which utilizes the update of the global model as momentum in the local updates of FedAvg. Moreover, both two algorithms, FedSpeed (Sun et al., 2023b) and FedGAMMA (Dai et al., 2023), utilize the gradient computed at the SAM-based perturbed weight. Specifically, FedSpeed tunes the local perturbed gradient with the proximal term, and FedGAMMA tunes it with the local perturbation of other clients. FedSMOO (Sun et al., 2023a) further tunes the local perturbation by leveraging the global perturbation approximated by Taylor expansion. Our FedGF is closely related to the methods that tune the local perturbations to pursue better performance. However, FedGF is unique in interpolating the local and global perturbations and managing the interpolation by observing the divergence between local and global models, yielding remarkable gains over others.

4. Proposed Method: FedGF

4.1. Training Process of FedGF

Based on the preliminaries in Section 2, we here focus on the local training and the aggregation steps of FedGF.

Local training: At the beginning of round r , client i receives the aggregated global model w^r , and runs K local epochs. For local epoch k , client i computes the two perturbed models, $\tilde{w}_{i,k}^r$ and \tilde{w}^r as follows:

$$g_{i,k}^r = \nabla F_i(w_{i,k}^r, \zeta_{i,k}) \quad (8)$$

$$\tilde{w}_{i,k}^r = w_{i,k}^r + \rho g_{i,k}^r / \|g_{i,k}^r\| \quad (\text{perturbed local model}) \quad (9)$$

$$\Delta^r = w^{r-1} - w^r \quad (10)$$

$$\tilde{w}^r = w^r + \rho \Delta^r / \|\Delta^r\| \quad (\text{perturbed global model}) \quad (11)$$

As FedSAM works, the perturbation of the local model

is computed (referring to Eq. (9)). When only using the perturbed local model, we already found that the aggregated global model is not located on flatter loss surface. To pursue the flatness of the global model for the global objective, we should consider the perturbation in the view of the global model and objective. However, the perturbation of the global model for the global objective cannot be tractable because each client cannot access the globally aggregated data samples across clients. To detour the hardship, we utilize the difference between the previous and the current global model, as formulated by Eq. (10), which is an approximated direction of the gradients for the global objective, i.e., $\nabla F(w^r)$. Based on the global perturbation, FedGF computes the perturbed global model in Eq. (11), which can be understood as the perturbed model with the maximally raised global objective loss. FedGF then interpolates the perturbed global and local models to compute $\tilde{w}_{i,k,c}^r$:

$$\tilde{w}_{i,k,c}^r = c\tilde{w}^r + (1-c)\tilde{w}_{i,k}^r, \quad (12)$$

where $0 \leq c \leq 1$ is the interpolation coefficient to control the position of $\tilde{w}_{i,k,c}^r$ between the global and local models. When c is close to 0, it indicates that FedGF becomes FedSAM. When c is close to 1, FedGF leans toward the global model to find flat minima around the global model. FedGF then computes the gradient based on the local epoch $\zeta_{i,k} \sim \mathcal{D}_i$ at position $\tilde{w}_{i,k,c}^r$ to update the local model:

$$w_{i,k+1}^r = w_{i,k}^r - \eta_l \nabla F_i(\tilde{w}_{i,k,c}^r, \zeta_{i,k}). \quad (13)$$

After K epochs, local training for round r is terminated.

Aggregation: The updated local model is then aggregated at the server to newly update the global model for the next round, i.e., w^{r+1} . Here, we adopt the global learning rate, η_g , suggested in (Reddi et al., 2021) (referring to Eq. (14)). When we set a global learning rate $\eta_g = 1$, it becomes the basis form of aggregation in Eq. (3).

$$w^{r+1} = w^r - \eta_g \sum_{i \in \mathcal{S}^r} \frac{m_i}{m} \Delta_i^r, \quad (14)$$

where $\Delta_i^r = w^r - w_{i,K}^r$.

4.2. Interpolation Coefficient c

Interpolation coefficient c controls the perturbed model in-between the local and global perturbations. Our key strategy to determine c is based on the following wisdom: When the local model largely deviates from the global model, i.e., a non-IID case, a larger c is preferred for focusing on the global model flatness; otherwise, i.e., an IID case, a smaller c is preferred. However, it is non-trivial to control c based on the non-IIDness because the server cannot access the local data distribution, which makes it difficult to measure the heterogeneity of the given setting. Thus, FedGF determines

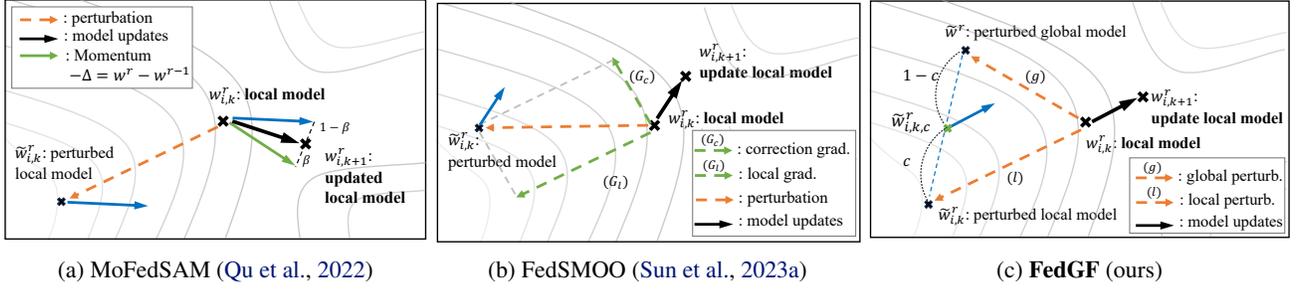


Figure 3: Schematic of MoFedSAM, FedSMOO, and FedGF. The gray line illustrates the loss landscape of local distribution.

c based on the divergence metric D^r , which represents how local models deviate from global model:

$$D^r = \frac{1}{|S^r|} \sum_{i \in S^r} \|w^r - w_{i,K}^r\|_2. \quad (15)$$

For mapping D^r to a value between 0 and 1, we adopt the thresholding: $I^r = \mathbf{I}[D^r > T_D]$, where $\mathbf{I}[\cdot]$ is an indicator function and $T_D > 0$ is a hyperparameter. For stability, we use the averaged I^r across recent W rounds in computing c :

$$c = \frac{1}{W} \sum_{i=r-W+1}^r I^i. \quad (16)$$

In the non-IID cases, as the round goes on, we empirically observe that c starts to increase at a relatively earlier round, which means that non-zero c is preferred to pursue the global flatness. Otherwise, c stays near zero for the IID case. Also, we further theoretically and empirically analyze that c is strongly related to the faster convergence of FedGF.

A pseudocode of FedGF is provided in Appendix C.

4.3. In-depth Comparison to the Existing Methods

Fig. 3 shows how FedGF is differentiated from the related methods. We here illustrate the schematics of our FedGF and two state-of-the-art methods, i.e., MoFedSAM and FedSMOO, while omitting the figure of FedSAM which is quite straightforward. Black-colored arrows indicate the model update at each local epoch. FedSAM and MoFedSAM only utilize the perturbation computed in view of the local objective. However, FedGF additionally considers the global model perturbation (see the orange-colored dotted arrows). Blue-colored arrows represent the gradient computed at the perturbed model. MoFedSAM uses the momentum for finding the model with a lower global objective (see the green-colored solid line), but FedGF utilizes the reverse direction of the momentum, as the global perturbation (the orange-colored arrow with a marker (g)), aiming to find flatter minima. FedSMOO adjusts the local perturbations by introducing a correction gradient (G_c), and it is closely

related to FedGF in adjusting perturbation vectors. The key differences are that FedGF explicitly utilizes the local and global perturbations and adaptively controls the interpolation between two perturbations by observing the divergence between the local and global models. Also, we found that FedSMOO strongly relies on auxiliary regularizations for restricting the local model to allocate near to the global model. FedGF solely works without the additional regularizations.

5. Theoretical Analysis

We provide the theoretical analysis of FedGF, including the convergence behavior and the flatness discrepancy. Before that, we introduce the following assumptions:

Assumption 1. (Smoothness of loss function) F_i is Lipschitz-smooth for all $i \in [N]$, i.e.,

$$\|\nabla F_i(w) - \nabla F_i(v)\| \leq L\|w - v\|$$

for all w, v in its domain and $i \in [N]$.

Assumption 2. (Bounds of gradients) The global variability of the local gradient is bounded by σ_g^2 , i.e.,

$$\|\nabla F_i(w^r) - \nabla F(w^r)\|^2 \leq \sigma_g^2,$$

for all $i \in [N]$ and r .

Assumption 3. (Bounds of the stochastic gradients) The stochastic gradient $\nabla F_i(w, \zeta_i)$, computed by client i with model parameter w using mini-batch ζ_i is an unbiased estimator of $\nabla F_i(w)$ with variance bounded by σ_l^2 , i.e.,

$$\mathbb{E}_{\zeta_i} \left\| \frac{\nabla F_i(w, \zeta_i)}{\|\nabla F_i(w, \zeta_i)\|} - \frac{\nabla F_i(w)}{\|\nabla F_i(w)\|} \right\|^2 \leq \sigma_l^2,$$

for all $i \in [N]$.

Assumption 1 and 2 are largely accepted by the prior non-convex FL convergence studies to assume the smoothness of loss function and the bounded heterogeneity (McMahan et al., 2017; Karimireddy et al., 2020; Reddi et al., 2021). Assumption 3, which bounds the variance of stochastic gradients, is from the work of FedSAM (Qu et al., 2022).

5.1. Convergence Analysis of FedGF

We present the convergence analysis of FedGF. Here, ϵ is the error in estimating the direction of the global perturbation of FedGF: $\epsilon := \|\Delta^r / \|\Delta^r\| - \nabla F(w^r) / \|\nabla F(w^r)\|\|$.

Theorem 1. *Let the learning rate be $\eta_l = \mathcal{O}(\frac{1}{\sqrt{RKL}})$, $\eta_g = \sqrt{KN}$, and the amplitude of perturbation is proportional to the learning rate, e.g., $\rho = \mathcal{O}(\frac{1}{\sqrt{R}})$. Under Assumptions 1 - 3 and full client participation, the average of the norm of the gradient generated by the iterative rounds of FedGF satisfies:*

$$\mathcal{O}\left(\frac{FL}{\sqrt{RKN}} + \frac{(1-c)^2}{R}\sigma_g^2 + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KN}}\sigma_l^2 + \frac{L^2c^2\epsilon^2}{R}\right), \quad (17)$$

where $F = F(\tilde{w}^0) - F(\tilde{w}^*)$ and $F(\tilde{w}^*) = \min_{\tilde{w}} F(\tilde{w})$. For the partial client participation strategy, if we choose the learning rates $\eta_l = \mathcal{O}(\frac{1}{\sqrt{RKL}})$, $\eta_g = \sqrt{KS}$ and $\rho = \mathcal{O}(\frac{1}{\sqrt{R}})$, the average of the norm of the gradient generated by the iterative rounds of FedGF satisfies:

$$\mathcal{O}\left(\frac{FL}{\sqrt{RKS}} + \left(\frac{(1-c)^2}{R} + 1\right)\sigma_g^2 + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KS}}\sigma_l^2 + \frac{L^2(c^2\epsilon^2 + 1)}{R}\right), \quad (18)$$

where $S = |\mathcal{S}^r|$.

Remark 1. (Faster convergence in non-IID cases) Let us focus on two variance terms, i.e., σ_g^2 and σ_l^2 , which represent the non-IIDness of the given FL setting and the stochastic variance of local gradients. These terms are crucial in the undesired delay of convergence of FL. When FedGF utilizes a larger c , i.e., reaching 1, then the related terms in Eq. (17) and (18) can be sufficiently suppressed; FedGF can accelerate the convergence even with the strong heterogeneity. In the extensive experiments, we confirm that FedGF tends to use larger values of c in the non-IID cases, leading to a significantly faster convergence than other related baselines.

Remark 2. (Error from ϵ vanishes as round goes on) Because FedGF approximates the gradient of the global model, there exists the error term ϵ , which probably hinders the convergence (referring to ϵ -involved term in Eq. (17) and (18)). As shown in the theorem, the error term diminishes as the round goes on; it does not ruin FedGF's convergence.

Remark 3. (FedGF with $c = 0$ becomes FedSAM) As aforementioned, with a smaller c , i.e., around 0, FedGF becomes FedSAM (the convergence analysis also becomes identical to that of (Qu et al., 2022)¹). Then FedGF does not

¹We found that the analysis for partial participation in (Qu et al., 2022) has a mistake, regarding the remained constant from σ_g . The details are in Appendix B.

utilize the global model perturbation, so the last terms in the convergence analysis disappear. However, the heterogeneity and variance terms related to σ_g^2 and σ_l^2 exist in the convergence behavior. When the setting becomes IID, it means that σ_g^2 and σ_l^2 are negligibly small, so FedGF tends to use smaller c values around 0 to behave like FedSAM.

5.2. Flatness Discrepancy Analysis

We address the claims on the flatness discrepancy, $\Delta_{\mathcal{F}}$.

Theorem 2. $\Delta_{\mathcal{F}}$ is upper bounded as follows:

$$\Delta_{\mathcal{F}} \leq \rho\sigma_g^2 + L\rho \sum_{i \in [N]} \frac{m_i}{m} \|w - w_i\| \quad (19)$$

Remark 4. (Heterogeneity, model divergence, and loss smoothness bound the flatness discrepancy) As shown in Eq. (19), the term σ_g^2 , which increases when heterogeneity gets worse, directly determines the upper bound of $\Delta_{\mathcal{F}}$, coinciding with our understanding of the discrepancy. Also, when the loss is not smooth, i.e., with a larger L , then the discrepancy increases. Finally, when the local model w_i largely deviates from the global model w , then the discrepancy gets worse, which agrees with how FedGF determines interpolation coefficient c based on the model divergence.

Theorem 3. For FedGF, if we choose $\eta_l = \mathcal{O}(\frac{1}{\sqrt{RKL}})$ and $\rho = \mathcal{O}(\frac{1}{\sqrt{R}})$, $\eta_g = \sqrt{KN}$, $\Delta_{\mathcal{F}}$ is then bounded as follows:

$$\Delta_{\mathcal{F}} \leq \mathcal{O}\left(\frac{\sigma_g^2}{\sqrt{R}} + \frac{L(1-c)^2\sigma_l^2}{R^{5/2}} + \frac{\sigma_g^2}{LR^{3/2}} + \frac{Lc^2\epsilon^2}{R^{5/2}}\right) \quad (20)$$

Remark 5. ($\Delta_{\mathcal{F}}$ is suppressed as round increases) For FedGF, the upper bound of $\Delta_{\mathcal{F}}$ formalized by Eq. (19) is suppressed as the round increases. For a non-IID, FedGF prefers to use large c , reaching 1, to strongly suppress the heterogeneity-related terms (referring to the second term).

The proofs of all claims are fully presented in Appendix A.

6. Experiments

We extensively evaluate the performance of FedGF² on the FL classification benchmarks suggested by (Caldarola et al., 2022), where CIFAR-10 and CIFAR-100 datasets are distributed on clients by utilizing Dirichlet distribution.

6.1. Experimental Settings

Baselines: We compare FedGF with the following methods: FedAvg (McMahan et al., 2017), FedAvgM (Hsu et al., 2019), SCAFFOLD (Karimireddy et al., 2020), FedProx (Li

²Codes are available at github.com/hwan-sig/Official-FedGF

Table 1: Test accuracies with of the FL algorithms on the CIFAR-10 and CIFAR-100 benchmarks

Task	Algorithms	Dirichlet distribution parameter α								
		$Dir.(\alpha = 0, \text{non-IID})$			$Dir.(\alpha = 0.005)$			$Dir.(\alpha = 10, \text{IID})$		
		Number of participating clients per each round								
		5	10	20	5	10	20	5	10	20
CIFAR-10	FedAvg	63.63	65.83	68.33	67.85	71.37	73.03	82.90	82.96	82.93
	FedAvgM	62.73	65.61	68.57	67.56	71.32	75.53	82.72	83.60	83.30
	FedProx	63.13	65.95	67.98	68.06	71.42	72.87	82.72	83.19	82.92
	SCAFFOLD	(X)	(X)	(X)	57.13	56.46	45.27	82.93	83.05	83.39
	FedDyn	66.84	71.01	69.45	70.74	73.78	75.43	83.07	83.58	83.67
	FedSAM	68.11	71.17	72.49	71.87	74.31	76.07	83.78	83.88	83.82
	FedASAM	73.32	74.5	75.49	74.96	75.59	76.57	83.11	83.28	82.89
	MoFedSAM	73.1	71.08	76.66	74.43	77.53	79.27	80.9	81.01	81.02
	FedGAMMA	45.32	47.55	35.07	46.99	48.44	35.58	74.99	66.12	54.85
	FedSMOO	68.82	71.59	72.48	71.9	74.46	75.44	83.72	83.67	83.79
	FedGF	78.41	79.68	80.86	78.79	79.39	79.69	84.71	83.94	83.85
CIFAR-100	FedAvg	29.35	33.79	36.62	38.15	40.58	41.27	50.41	50.20	49.98
	FedAvgM	29.94	30.07	39.35	38.64	40.72	48.44	50.37	51.2	50.57
	FedProx	29.19	33.16	36.41	38.54	40.52	40.77	50.10	49.98	49.96
	SCAFFOLD	(X)	(X)	(X)	36.25	(X)	(X)	52.28	52.12	52.48
	FedDyn	(X)	(X)	(X)	(X)	(X)	(X)	51.74	52.41	52.59
	FedSAM	29.43	34.32	36.88	42.28	44.57	45.18	54.06	53.75	53.5
	FedASAM	34.43	37.09	38.93	44.36	45.76	46.94	54.6	54.42	54.73
	MoFedSAM	29.02	35.82	41.26	34.64	42.24	44.92	52.13	52.21	52.07
	FedGAMMA	(X)	(X)	(X)	20.52	14.76	10.33	47.43	38.18	25.06
	FedSMOO	35.35	38.78	40.82	44.39	46.03	47.5	54.31	54.89	54.65
	FedGF	45.37	46.86	47.77	46.48	46.70	46.08	54.16	54.62	54.59

(X) indicates that the method fails to train, so the results remain at the same level as the random prediction.

et al., 2020a), FedDyn (Acar et al., 2021), FedSAM (Caldarola et al., 2022; Qu et al., 2022), FedASAM (Caldarola et al., 2022), MoFedSAM (Qu et al., 2022), FedGAMMA (Dai et al., 2023), and FedSMOO³ (Sun et al., 2023a).

Model architecture: We follow the model architecture described in the prior FL works (Hsu et al., 2020; Caldarola et al., 2022), which is a variant of the LeNet architecture by (LeCun et al., 1998). For the larger architecture, such as ResNet-18, we add the results in Appendix D.6.

FL settings: For a given server and 100 clients, we test three different numbers of participating clients per round, i.e., $\{5, 10, 20\}$. We distributed 500 data samples per client, and the number of local updates per round is 8, with batch size 64. As done in (Hsu et al., 2020), the prior distribution of local data follows the Dirichlet distribution of α , i.e., $\alpha \in \{0, 0.005, 10\}$ for both CIFAR-10 and CIFAR-100 experiments. When α increases, the setting becomes a IID case. When α goes to zero, it means a non-IID case. The communication round goes up to 10,000 and 20,000 for CIFAR-10 and CIFAR-100, respectively.

Further details of the benchmarks, the hyperparameters, and the model architecture are provided in Appendix C.

³FedSMOO strongly relies on the dynamic regularization. For a fair comparison of the effectiveness on finding flatter global model, we evaluate FedSMOO without the regularizer.

6.2. Performance Evaluation

We evaluate the test accuracy of FL algorithms in Table 1 on the CIFAR-10/100 benchmarks. We here provide the following key findings based on the experimental results.

6.2.1. LARGE GAINS FOR THE NON-IID CASES

FedGF significantly outperforms the existing works in the non-IID settings, i.e., $\alpha = 0$. Specifically, it shows the gains ranging from +4.20% to +5.30% for CIFAR-10 cases and larger gains ranging from +6.51% to +10.94% for CIFAR-100 over the runner-ups. The results verify that FedGF effectively relieves the strong heterogeneity via aggregating a global model with a strong generalization across clients. We believe that the gains directly come from the efforts to search flat minima of global model by FedGF, which is to be confirmed in the following part. As the heterogeneity becomes relieved, i.e., as α increases from 0 to 10, the performance gaps between FedGF and prior works are reduced. This is due to the homogeneity of data distribution in the IID cases, which relieves the discrepancy between the local and global models. Also, we found that the regularization-based FL methods, including SCAFFOLD and FedDyn, is not successful in the non-IID case⁴, particularly in CIFAR-100.

⁴We want remind that the original evaluation in their works are done in the cases with moderate non-IIDness.

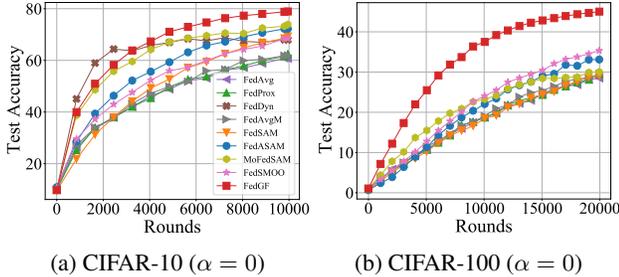


Figure 4: Convergence behaviors for non-IID

6.2.2. FASTER CONVERGENCE BEHAVIOR

We present the convergence behaviors for the non-IID case in Fig. 4a and 4b. We clearly confirm that FedGF shows significantly faster convergence than others, particularly emphasized in the CIFAR-100 case. The results verify **Remark 1** of Theorem 1, which emphasizes that FedGF can suppress the terms of the heterogeneity and the stochastic variance, so it accelerates the speed of convergence by preferring larger c . In the experiments, it indeed pushes c to be 1 for the non-IID cases, leading to faster convergence.

6.2.3. ROBUSTNESS TO PARTICIPATING CLIENTS

As shown in Table 1, when the heterogeneity gets worse, the FL baselines suffer from severe degradations when the number of participating clients is limited. In contrast, FedGF shows the robust performance for the limited participating clients. We emphasize that MoFedSAM shows a significant drop from 41.26% to 29.02% when the number of participating clients decreases from 20 to 5 for the CIFAR-100 $\alpha = 0$ case. It happens because the momentum of global models used by MoFedSAM largely fluctuates round-by-round when the number of participating clients is limited. FedSMOO suffers from 5.47% drop for the same case. We conjecture that FedSMOO struggles to find the robust perturbation correction term when the number of participations is limited. On the contrary, let us remind **Remark 2** of Theorem 1, which points out that FedGF can suppress the error, ϵ , in estimating the global perturbation as round increases.

6.2.4. FLATNESS RESULTS

To confirm the flatness of the global model in both quantitatively and qualitatively, we present various flatness results: i) loss plots, ii) flatness metrics, including flatness discrepancy, and iii) visualization of loss surface.

Loss plots along to perturbations: Fig. 5 shows the plots to confirm how the loss value increases as the perturbation of the model parameter is imposed for the CIFAR-100 experiments with 5 participating clients. In the non-IID case, FedGF shows slightly flatter loss plot than Fe-

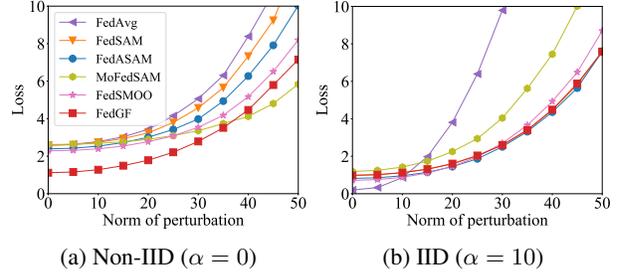


Figure 5: Loss plots along to perturbation for CIFAR-100

 Table 2: LPF, λ_{\max} , and $\Delta_{\mathcal{F}}$ results for CIFAR-100

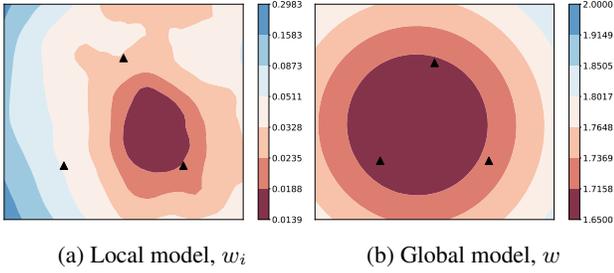
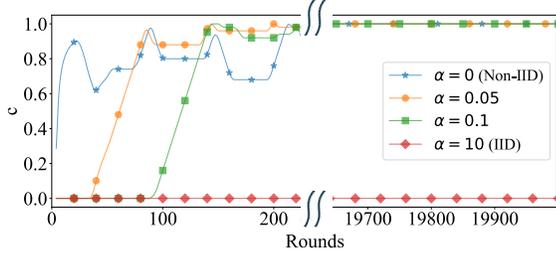
Algorithm	Non-IID ($\alpha = 0$)			IID ($\alpha = 10$)		
	LPF \downarrow	$\lambda_{\max} \downarrow$	$\Delta_{\mathcal{F}} \downarrow$	LPF \downarrow	$\lambda_{\max} \downarrow$	$\Delta_{\mathcal{F}} \downarrow$
FedAvg	2.67	81.57	0.39	1.24	103.19	0.11
FedSAM	2.67	43.32	0.33	1.10	25.36	0.04
FedASAM	2.35	25.46	0.24	0.94	14.99	0.04
MoFedSAM	2.64	15.11	0.13	1.48	27.28	0.04
FedSMOO	2.37	25.83	0.12	0.8	23.38	0.08
FedGF	1.36	14.07	0.13	1.08	23.54	0.04

\downarrow : a lower value is preferred.

dAvg, FedSAM, and FedASAM. Although MoFedSAM and FedSMOO show flatter behavior, our FedGF shows much lower loss values in the wide range of perturbation; the gap in loss is more than 1.0, which leads to the large performance gap in accuracies (as confirmed in Table 1). In the IID case, the FL algorithms show similar behavior excepting for FedAvg and MoFedSAM. We conjecture that MoFedSAM suffers from the fluctuation of global momentums caused by the limited number of clients; leading to the unexpected sharp loss plot.

Flatness metrics (LPF, λ_{\max} , and $\Delta_{\mathcal{F}}$): The loss plots show a brief understanding of loss surface, but they cannot provide quantitative measurements of flatness. We here compute various flatness metrics for an in-depth analysis: the maximum eigenvalue of the Hessian matrix, i.e., λ_{\max} , which is commonly used in prior works, Low-Pass Filter (LPF) based metric, which is recently suggested to show the robust correlation to generalization (Bisla et al., 2022), and the proposed flatness discrepancy $\Delta_{\mathcal{F}}$. While FedAvg showing the worst values, FedGF shows the best flatness for LPF and λ_{\max} in the non-IID case. FedGF is the second best for $\Delta_{\mathcal{F}}$ with a minimal gap. As noted in **Remark 4** of Theorem 2, we found that non-IIDness increases $\Delta_{\mathcal{F}}$ for all cases. Also, we confirm that FedGF sufficiently suppresses discrepancy as noted in **Remark 5** of Theorem 3.

Visualization of loss surface: In Fig. 6, we visualize the loss surface of the local and global models of FedGF for F_i and F for the CIFAR-100 cases. It shows that local model shows a moderate flatness, and the global model shows flatter loss surface. When compared with the FedSAM’s


 Figure 6: Loss surface of FedGF for CIFAR-100 ($\alpha = 0$).

 Figure 7: Behavior of c for CIFAR-100

loss surface (as in Fig. 2), we visually confirm that FedGF finds flatter minima of global model for the global objective.

6.2.5. ANALYSIS OF COMMUNICATION COST

Herein, we analyze the communication costs of FedGF. In Table 3, the number of model transmissions is measured for the related works focusing on the flatness searching FL methods ('1' means a single transmission of model parameters). The baselines, which are FedAvg and FedSAM, upload only the locally trained model and download the averaged global model, leading to a communication cost of 2. FedSMOO and FedGAMMA, state-of-the-art flatness-searching FL methods, require higher costs, i.e., 4 model transmissions, because they transmit both parameters and perturbations between server and clients. In contrast, FedGF transmits the perturbation, which is a pseudo-gradient, from server to client for each round but does not require uploading from the client to the server, leading to the moderate cost of three model transmissions.

Table 3: Number of model transmissions per round

Algorithms	client \rightarrow server (upload)	client \leftarrow server (download)	total
FedAvg	1	1	2
FedSAM	1	1	2
MoFedSAM	1	2	3
FedGF	1	2	3
FedSMOO	2	2	4
FedGAMMA	2	2	4

In Table 4, we compare the actual communication costs, which are the measured number of model and perturbation transmissions, to reach the saturated accuracies of the most typical baseline, i.e., FedAvg. We compare FedGF with FedSMOO, which mostly follows FedGF as a runner-up for all cases (we consider the cases with 5 participating clients). We found that FedGF shows significantly faster convergence, where FedSMOO requires $\times 2.40$ and $\times 2.98$ times higher costs for the non-IID cases. This result coincides with theories (referring to **Remark 1** of Theorem 1) and the empirical results (referring to the part 6.2.2).

 Table 4: Number of model transmissions to reach the FedAvg's final performance ($\times 10^2$)

	Algorithms	$\alpha = 0$	$\alpha = 0.005$	$\alpha = 10$
CIFAR-10	FedGF	75	108	165
	FedSMOO	180 (x2.40)	228 (x2.11)	244 (x1.48)
CIFAR-100	FedGF	180	210	240
	FedSMOO	536 (x2.98)	440 (x2.10)	248 (x1.03)

6.2.6. ANALYSIS OF c

Behavior: Fig. 7 shows how FedGF utilizes c values according to the rounds. FedGF computes c based on the model divergence in Eq. (15) and (16). For the IID case, FedGF steadily uses $c = 0$ due to the minimal divergence between the local and global models, which coincides with **Remark 3** of Theorem 1. When non-IIDness gets worse, FedGF prefers to use larger c values up to 1, which strongly employs the global perturbation. This exactly agrees with the interpretation of **Remark 1** of Theorem 1, which states that FedGF relieves the heterogeneity by letting c be larger.

 Table 5: Static c vs. FedGF (adaptive c)

Dataset	IIDness	$c = 0$	$c = 0.5$	$c = 1$	FedGF
CIFAR-10	Non-IID	68.11	71.24	78.05	78.41
	IID	83.78	82.95	81.94	84.71
CIFAR-100	Non-IID	29.43	26.64	44.39	45.37
	IID	54.06	52.47	46.68	54.16

Ablation on static c : As shown in Table 5, FedGF, which adaptively computes c , is better than the cases of static c .

7. Conclusion

We rethink flat minima searching in FL with the novel perspective of flatness discrepancy. It gets worse when the heterogeneity becomes severe, leading to the deterioration of the prior flat minima searching FL algorithms. Based on this wisdom, we propose FedGF, which can relieve the discrepancy by utilizing both local and global perturbations in the SAM optimizer. FedGF largely outperforms the existing FL methods, particularly in non-IID cases.

Acknowledgements

This work was supported by the Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. RS-2020-II201336, Artificial Intelligence Graduate School Program (UNIST)), (No. RS-2021-II212201, Federated Learning for Privacy-Preserving Video Caching Networks), and (No. IITP-2024-RS-2022-00156361, Innovative Human Resource Development for Local Intellectualization program).

Impact Statement

With a broader perspective, FL would be a key technology that secures private data while training deep models. We conjecture that our method would not raise an ethical issue.

References

- Acar, D. A. E., Zhao, Y., Matas, R., Mattina, M., Whatmough, P., and Saligrama, V. Federated learning based on dynamic regularization. In *International Conference on Learning Representations (ICLR)*, 2021.
- Andriushchenko, M., Croce, F., Müller, M., Hein, M., and Flammarion, N. A modern look at the relationship between sharpness and generalization. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023.
- Bisla, D., Wang, J., and Choromanska, A. Low-pass filtering sgd for recovering flat optima in the deep learning optimization landscape. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 8299–8339. PMLR, 2022.
- Caldarola, D., Caputo, B., and Ciccone, M. Improving generalization in federated learning by seeking flat minima. In *European Conference on Computer Vision (ECCV)*, pp. 654–672. Springer, 2022.
- Cha, J., Chun, S., Lee, K., Cho, H.-C., Park, S., Lee, Y., and Park, S. Swad: Domain generalization by seeking flat minima. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:22405–22418, 2021.
- Dai, R., Yang, X., Sun, Y., Shen, L., Tian, X., Wang, M., and Zhang, Y. Fedgamma: Federated learning with global sharpness-aware minimization. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- Dinh, L., Pascanu, R., Bengio, S., and Bengio, Y. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning (ICML)*, pp. 1019–1028. PMLR, 2017.
- Du, Z., Sun, J., Li, A., Chen, P.-Y., Zhang, J., Li, H. H., and Chen, Y. Rethinking normalization methods in federated learning. In *Proceedings of the 3rd International Workshop on Distributed Machine Learning*, pp. 16–22, 2022.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations (ICLR)*, 2021.
- Gasanov, E., Khaled, A., Horváth, S., and Richtárik, P. Flix: A simple and communication-efficient alternative to local methods in federated learning. In *International Conference on Machine Learning (ICML)*, pp. 11374–11421. PMLR, 2021.
- Hochreiter, S. and Schmidhuber, J. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- Hsu, T.-M. H., Qi, H., and Brown, M. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- Hsu, T.-M. H., Qi, H., and Brown, M. Federated Visual Classification with Real-World Data Distribution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, August 2020.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning (ICML)*, pp. 448–456. pmlr, 2015.
- Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D., and Wilson, A. G. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A. T. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning (ICML)*, pp. 5132–5143. PMLR, 2020.
- Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations (ICLR)*, 2017.
- Kim, H., Park, J., Choi, Y., and Lee, J. Fantastic robustness measures: The secrets of robust generalization. In *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- Krizhevsky, A., Nair, V., and Hinton, G. Cifar-10 (canadian institute for advanced research).

- Kwon, J., Kim, J., Park, H., and Choi, I. K. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *International Conference on Machine Learning (ICML)*, pp. 5905–5914. PMLR, 2021.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems (MLSys)*, 2:429–450, 2020a.
- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations (ICLR)*, 2020b.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics (AISTATS)*, pp. 1273–1282. PMLR, 2017.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Qu, Z., Li, X., Duan, R., Liu, Y., Tang, B., and Lu, Z. Generalized federated learning via sharpness aware minimization. In *International Conference on Machine Learning (ICML)*, pp. 18250–18280. PMLR, 2022.
- Reddi, S. J., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konečný, J., Kumar, S., and McMahan, H. B. Adaptive federated optimization. In *International Conference on Learning Representations (ICLR)*, 2021.
- Sun, Y., Shen, L., Chen, S., Ding, L., and Tao, D. Dynamic regularized sharpness aware minimization in federated learning: Approaching global consistency and smooth landscape. In *International Conference on Machine Learning (ICML)*, pp. 32991–33013. PMLR, 2023a.
- Sun, Y., Shen, L., Huang, T., Ding, L., and Tao, D. Fed-speed: Larger local interval, less communication round, and higher generalization accuracy. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023b.
- Zeng, D., Liang, S., Hu, X., Wang, H., and Xu, Z. Fedlab: A flexible federated learning framework. *Journal of Machine Learning Research*, 24(100):1–7, 2023.

A. Mathematical Details of Theoretical Analysis

A.1. Preliminary Assumptions, and Lemmas

We recall the introduced assumptions as follows:

Assumption 1. (*Smoothness of loss function*) F_i is Lipschz-smooth for all $i \in [N]$, i.e.,

$$\|\nabla F_i(w) - \nabla F_i(v)\| \leq L\|w - v\|$$

for all w, v in its domain and $i \in [N]$.

Assumption 2. (*Bounds of gradients*) The global variability of the local gradient is bounded by σ_g^2 , i.e.,

$$\|\nabla F_i(w^r) - \nabla F(w^r)\|^2 \leq \sigma_g^2,$$

for all $i \in [N]$ and r .

Assumption 3. (*Bounds of the stochastic gradients*) The stochastic gradient $\nabla F_i(w, \zeta_i)$, computed by client i with model parameter w using mini-batch ζ_i is an unbiased estimator of $\nabla F_i(w)$ with variance bounded by σ_i^2 , i.e.,

$$\mathbb{E}_{\zeta_i} \left\| \frac{\nabla F_i(w, \zeta_i)}{\|\nabla F_i(w, \zeta_i)\|} - \frac{\nabla F_i(w)}{\|\nabla F_i(w)\|} \right\|^2 \leq \sigma_i^2,$$

for all $i \in [N]$.

To analyze the convergence rate of FedGF, we first state some preliminary lemmas and their proofs as follows:

Lemma 1. (*Relaxed triangle inequality*) Let $\{v_1, \dots, v_\tau\}$ be τ vectors in \mathbb{R}^d . Then, the following are true: (1) $\|v_i + v_j\|^2 \leq (1 + a)\|v_i\|^2 + (1 + \frac{1}{a})\|v_j\|^2$ for any $a > 0$, and (2) $\|\sum_{i=1}^\tau v_i\|^2 \leq \tau \sum_{i=1}^\tau \|v_i\|^2$.

Lemma 2. For random variables x_1, \dots, x_n , we have

$$\mathbb{E} \left[\|x_1 + \dots + x_n\|^2 \right] \leq n \mathbb{E} \left[\|x_1\|^2 + \dots + \|x_n\|^2 \right].$$

Lemma 3. For independent, mean 0 random variables x_1, \dots, x_n , we have

$$\mathbb{E} \left[\|x_1 + \dots + x_n\|^2 \right] = \mathbb{E} \left[\|x_1\|^2 + \dots + \|x_n\|^2 \right]$$

For the Lemmas 1, 2, and 3, they are identically introduced by the work of FedSAM (Qu et al., 2022), where they are indexed with Lemma A.1, A.2, and A.3, respectively. Thus, we skip the proofs by referring to the proofs in (Qu et al., 2022).

A.2. Main Lemmas and Theorems for Convergence Analysis, and their Proofs

Lemma 4. (*Bounded global variance of $\|\nabla F_i(w + \delta_i) - \nabla F(w + \delta)\|^2$.)* An immediate implication of Assumptions 1 and 2, the variance of local and global gradients with perturbation can be bounded as follows:

$$\|\nabla F_i(w + \delta_i) - \nabla F(w + \delta)\|^2 \leq 3\sigma_g^2 + 6L^2\rho^2.$$

Proof.

$$\begin{aligned} \|\nabla F_i(w + \delta_i) - \nabla F(w + \delta)\|^2 &= \|\nabla F_i(w + \delta_i) - \nabla F_i(w) + \nabla F_i(w) - \nabla F(w) + \nabla F(w) - \nabla F(w + \delta)\|^2 \\ &\leq 3\|\nabla F_i(w + \delta_i) - \nabla F_i(w)\|^2 + 3\|\nabla F_i(w) - \nabla F(w)\|^2 + 3\|\nabla F(w) - \nabla F(w + \delta)\|^2 \end{aligned} \quad (21)$$

$$\leq 3\sigma_g^2 + 6L^2\rho^2, \quad (22)$$

where Eq. (21) is from Lemma 2 and Eq. (22) is from Assumption 1, 2 and the perturbation is bounded by ρ . Up to this point, FedGF steps on the identical mathematical claims of FedSAM (Qu et al., 2022). \square

From the following **Lemma 5**, our analysis of FedGF deviates from the one of FedSAM in (Qu et al., 2022).

For a brief notation, we denote $\tilde{w}_{i,k,c}$ as $\bar{w}_{i,k}$ in the following analysis, i.e., $\bar{w}_{i,k}^r := \tilde{w}_{i,k,c}^r = c\tilde{w}^r + (1-c)\tilde{w}_{i,k}^r$, and denote the expectation of clients' sampling \mathbb{E}_ζ as \mathbb{E} .

Lemma 5. (Bounded \mathcal{E}_w of FedGF.) Suppose our functions satisfies Assumptions 1-2. Then, the updates of FedGF for any learning rate satisfying $(1-c)\eta_l \leq \frac{1}{8KL}$ have the drift due to $w_{i,k} - w$:

$$\mathcal{E}_w = \mathbb{E} \|w_{i,k} - w\|^2 \leq Ke(T + 8K\eta_l^2 \|\nabla f(\tilde{w})\|^2 + 16KL^2\eta_l^2 c^2 \rho^2 \epsilon^2),$$

where e is Euler's number and $T = 8K\eta_l^2 [L^2\sigma_l^2(1-c)^2\rho^2 + 8K^2L^4\eta_l^2\rho^2(1-c)^2 + 3\sigma_g^2 + 6L^2\rho^2]$

Proof.

$$\begin{aligned} \mathbb{E} \|w_{i,k} - w\|^2 &= \mathbb{E} \|w_{i,k-1} - w - \eta_l \nabla F_i(\bar{w}_{i,k-1}, \zeta_i)\|^2 \\ &\leq \left(1 + \frac{1}{2K-1}\right) \mathbb{E} \|w_{i,k-1} - w\|^2 + \underbrace{2K\eta_l^2 \mathbb{E} \|\nabla F_i(\bar{w}_{i,k-1}, \zeta_i)\|^2}_A, \end{aligned} \quad (23)$$

where Eq. (23) is from Lemma 1.

For A ,

$$\begin{aligned} &2K\eta_l^2 \mathbb{E} \|\nabla F_i(\bar{w}_{i,k-1}, \zeta_i)\|^2 \\ &= 2K\eta_l^2 \mathbb{E} \|\nabla F_i(\bar{w}_{i,k-1}, \zeta_i) - \nabla F_i(\bar{w}_{i,k-1}) + \nabla F_i(\bar{w}_{i,k-1}) - \nabla F_i(\tilde{w}) + \nabla F_i(\tilde{w}) - \nabla F(\tilde{w}) + \nabla F(\tilde{w})\|^2 \\ &\leq 8K\eta_l^2 \mathbb{E} \left[\underbrace{\|\nabla F_i(\bar{w}_{i,k-1}, \zeta_i) - \nabla F_i(\bar{w}_{i,k-1})\|^2}_{A_1} + \underbrace{\|\nabla F_i(\bar{w}_{i,k-1}) - \nabla F_i(\tilde{w})\|^2}_{A_2} + \underbrace{\|\nabla F_i(\tilde{w}) - \nabla F(\tilde{w})\|^2}_{A_3} + \|\nabla F(\tilde{w})\|^2 \right] \end{aligned} \quad (24)$$

where Eq. (24) is from Lemma 2.

For A_1 ,

$$\begin{aligned} \mathbb{E} \|\nabla F_i(\bar{w}_{i,k-1}, \zeta_i) - \nabla F_i(\bar{w}_{i,k-1})\|^2 &\leq L^2 \|\bar{w}_{i,k-1}, \zeta_i - \bar{w}_{i,k-1}\|^2 \\ &= L^2 \mathbb{E} \left\| (1-c) \left(w_{i,k-1} + \rho \frac{\nabla F_i(w_{i,k-1}, \zeta_i)}{\|\nabla F_i(w_{i,k-1}, \zeta_i)\|} \right) + c\tilde{w} - (1-c) \left(w_{i,k-1} + \rho \frac{\nabla F_i(w_{i,k-1})}{\|\nabla F_i(w_{i,k-1})\|} \right) - c\tilde{w} \right\|^2 \\ &= L^2 (1-c)^2 \rho^2 \mathbb{E} \left\| \frac{\nabla F_i(w_{i,k-1}, \zeta_i)}{\|\nabla F_i(w_{i,k-1}, \zeta_i)\|} - \frac{\nabla F_i(w_{i,k-1})}{\|\nabla F_i(w_{i,k-1})\|} \right\|^2 \\ &\leq L^2 (1-c)^2 \rho^2 \sigma_l^2, \end{aligned} \quad (25)$$

$$\leq L^2 (1-c)^2 \rho^2 \sigma_l^2, \quad (26)$$

where Eq. (25) is from $(\bar{w}_{i,k-1}, \zeta_i) = (1-c) \left(w_{i,k-1} + \frac{\nabla F_i(w_{i,k-1}, \zeta_i)}{\|\nabla F_i(w_{i,k-1}, \zeta_i)\|} \right) + c\tilde{w}$, and Eq. (26) is by Assumption 3.

For A_2 ,

$$\begin{aligned} & \mathbb{E} \|\nabla F_i(\bar{w}_{i,k-1}) - \nabla F_i(\tilde{w})\|^2 \leq L^2 \mathbb{E} \|\bar{w}_{i,k-1} - \tilde{w}\|^2 \\ & = L^2 \mathbb{E} \left\| (1-c)\tilde{w}_{i,k-1} + c \left(w^r + \rho \frac{\Delta^r}{\|\Delta^r\|} + \rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} - \rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} \right) - \tilde{w} \right\|^2 \end{aligned} \quad (27)$$

$$\begin{aligned} & = L^2 \mathbb{E} \left\| (1-c)\tilde{w}_{i,k-1} + c \left(w^r + \rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} \right) - \tilde{w} + c\rho \frac{\Delta^r}{\|\Delta^r\|} - c\rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} \right\|^2 \\ & \leq 2L^2 \mathbb{E} \|(1-c)(\tilde{w}_{i,k-1} - \tilde{w})\|^2 + 2L^2 c^2 \rho^2 \epsilon^2 \end{aligned} \quad (28)$$

$$\begin{aligned} & = 2(1-c)^2 L^2 \mathbb{E} \left\| w_{i,k-1} - w + \rho \frac{\nabla F_i(w_{i,k-1}, \zeta_i)}{\|\nabla F_i(w_{i,k-1}, \zeta_i)\|} - \rho \frac{\nabla F(w)}{\|\nabla F(w)\|} \right\|^2 + 2L^2 c^2 \rho^2 \epsilon^2 \\ & \leq 4(1-c)^2 L^2 \mathbb{E} \left[\|w_{i,k-1} - w\|^2 + \rho^2 \|\delta_{i,k-1} - \delta\|^2 \right] + 2L^2 c^2 \rho^2 \epsilon^2 \end{aligned} \quad (29)$$

$$= 4(1-c)^2 L^2 \mathbb{E} \|w_{i,k-1} - w\|^2 + 8K^2 L^4 \eta_i^2 \rho^2 (1-c)^2 + 2L^2 c^2 \rho^2 \epsilon^2 \quad (30)$$

where Eq. (27) is from $\bar{w}_{i,k-1} = (1-c)\tilde{w}_{i,k-1} + c \left(w^r + \rho \frac{\Delta^r}{\|\Delta^r\|} \right)$, Eq. (28) is from the definition of ϵ^2 , i.e., $\epsilon := \|\Delta^r/\|\Delta^r\| - \nabla F(w^r)/\|\nabla F(w^r)\|\|$, and Eq. (29) is from Lemma 2, $\delta := \nabla F(w)/\|\nabla F(w)\|$, and $\delta_{i,k-1} := \nabla F_i(w_{i,k-1}, \zeta_i)/\|\nabla F_i(w_{i,k-1}, \zeta_i)\|$. Moreover, for Eq. (30), we borrow the result of Lemma B.1 in (Qu et al., 2022).

For A_3 ,

$$\|\nabla F_i(\tilde{w}) - \nabla F(\tilde{w})\|^2 \leq 3\sigma_g^2 + 6L^2 \rho^2 \quad (31)$$

where Eq. (31) is from Lemma 4.

Therefore, when utilizing Eq. (24), (26), (30), and (31):

$$\begin{aligned} \mathbb{E} \|w_{i,k} - w\|^2 & \leq \left(1 + \frac{1}{2K-1} \right) \mathbb{E} \|w_{i,k-1} - w\|^2 + 8K\eta_i^2 \left[L^2 \sigma_i^2 (1-c)^2 \rho^2 + 4(1-c)^2 L^2 \mathbb{E} \|w_{i,k-1} - w\|^2 \right. \\ & \quad \left. + 8K^2 L^4 \eta_i^2 \rho^2 (1-c)^2 + 2L^2 c^2 \rho^2 \epsilon^2 + 3\sigma_g^2 + 6L^2 \rho^2 + \|\nabla F(\tilde{w})\|^2 \right] \\ & = \left(1 + \frac{1}{2K-1} + 32K\eta_i^2 L^2 (1-c)^2 \right) \mathbb{E} \|w_{i,k-1} - w\|^2 \\ & \quad + 8K\eta_i^2 \left[L^2 \sigma_i^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_i^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2 + 2L^2 c^2 \rho^2 \epsilon^2 + \|\nabla F(\tilde{w})\|^2 \right] \end{aligned} \quad (32)$$

If η_i satisfies

$$\begin{aligned} 1 + \frac{1}{2K-1} + 32K\eta_i^2 L^2 (1-c)^2 & \leq 1 + \frac{1}{K-1} \\ 32K\eta_i^2 L^2 (1-c)^2 & \leq \frac{1}{K-1} - \frac{1}{2K-1} = \frac{K}{(2K-1)(K-1)} \\ 32\eta_i^2 L^2 (1-c)^2 & \leq \frac{1}{(2K-1)(K-1)} \\ (1-c)\eta_i & \leq \frac{1}{4L\sqrt{2(2K-1)(K-1)}}, \end{aligned}$$

Eq. (32) will be

$$\mathbb{E} \|w_{i,k} - w\|^2 \leq \left(1 + \frac{1}{K-1} \right) \mathbb{E} \|w_{i,k-1} - w\|^2 + T + 8K\eta_i^2 \|\nabla F(\tilde{w})\|^2 + 16KL^2 \eta_i^2 c^2 \rho^2 \epsilon^2, \quad (33)$$

where $T = 8K\eta_i^2 \left[L^2 \sigma_i^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_i^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2 \right]$.

By recursively substitutes the $w_{i,k-1}$ related term via decreasing the index k down to 0 of Eq. (33), we obtain:

$$\mathbb{E} \|w_{i,k} - w\|^2 \leq \sum_{\tau=0}^{k-1} \left(1 + \frac{1}{K-1} \right)^\tau (T + 8K\eta_i^2 \|\nabla F(\tilde{w})\|^2 + 16KL^2 \eta_i^2 c^2 \rho^2 \epsilon^2) \quad (34)$$

$$\leq Ke(T + 8K\eta_i^2 \|\nabla F(\tilde{w})\|^2 + 16KL^2 \eta_i^2 c^2 \rho^2 \epsilon^2), \quad (35)$$

where Eq. (35) is from the definition of Euler's number (e) (when the k becomes infinite, the related term transforms to Euler's number, but k is finite here, so the inequality holds). □

A.3. Full Client Participating Convergence

We are ready to handle the convergence analysis for the full client participating case.

Lemma 6.

$$\begin{aligned} & \left\langle \nabla F(\tilde{w}^r), \mathbb{E} \left[-\frac{1}{N} \sum_i^N \sum_k^K \eta_l \nabla F_i(\bar{w}_{i,k}^r) \right] + K\eta_l \nabla F(\tilde{w}^r) \right\rangle \leq \\ & \left(\frac{\eta_l K}{2} + 16e(1-c)^2 K^3 \eta_l^3 L^2 \right) \|\nabla F(\tilde{w}^r)\|^2 + 2e(1-c)^2 K^2 L^2 \eta_l T + 32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \rho^2 \epsilon^2 \\ & + 4K^3 L^4 \eta_l^3 \rho^2 (1-c)^2 + K\eta_l L^2 c^2 \epsilon^2 - \frac{\eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2, \end{aligned}$$

where $T = 8K\eta_l^2 [L^2\sigma_l^2(1-c)^2\rho^2 + 8K^2L^4\eta_l^2\rho^2(1-c)^2 + 3\sigma_g^2 + 6L^2\rho^2]$, and $\langle a, b \rangle$ is the inner product of vectors a and b with the same dimensionality.

Proof.

$$\begin{aligned} & \left\langle \nabla F(\tilde{w}^r), \mathbb{E} \left[-\frac{1}{N} \sum_i^N \sum_k^K \eta_l \nabla F_i(\bar{w}_{i,k}^r) \right] + K\eta_l \nabla F(\tilde{w}^r) \right\rangle \\ & = \frac{\eta_l K}{2} \|\nabla F(\tilde{w}^r)\|^2 + \underbrace{\frac{\eta_l}{2N^2K} \mathbb{E} \left\| \sum_{i,k} [\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)] \right\|^2}_A - \frac{\eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \quad (36) \end{aligned}$$

where Eq. (36) is from $\langle a, b \rangle = \frac{1}{2}(\|a\|^2 + \|b\|^2 - \|a - b\|^2)$ with $a = \sqrt{\eta_l K} \nabla F(\tilde{w}^r)$ and $b = -\frac{\sqrt{\eta_l}}{N\sqrt{K}} \sum_{i,k} (\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r))$.

For A ,

$$\begin{aligned} & \mathbb{E} \left\| \sum_{i,k} [\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)] \right\|^2 \leq KN \sum_{i,k} \mathbb{E} \|\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)\|^2 \\ & \leq KNL^2 \sum_{i,k} \mathbb{E} \|\bar{w}_{i,k}^r - \tilde{w}^r\|^2 \quad (37) \end{aligned}$$

$$\begin{aligned}
 &= K^2 N^2 L^2 \mathbb{E} \left\| (1-c)\tilde{w}_{i,k-1} + c \left(w^r + \rho \frac{\Delta^r}{\|\Delta^r\|} + \rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} - \rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} \right) - \tilde{w} \right\|^2 \\
 &= K^2 N^2 L^2 \mathbb{E} \left\| (1-c)\tilde{w}_{i,k-1} + c \left(w^r + \rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} \right) - \tilde{w} + c\rho \frac{\Delta^r}{\|\Delta^r\|} - c\rho \frac{\nabla F(w^r)}{\|\nabla F(w^r)\|} \right\|^2 \\
 &\leq 2K^2 N^2 L^2 \mathbb{E} \|(1-c)(\tilde{w}_{i,k-1} - \tilde{w})\|^2 + 2K^2 N^2 L^2 c^2 \rho^2 \epsilon^2 \tag{38}
 \end{aligned}$$

$$\begin{aligned}
 &= 2(1-c)^2 K^2 N^2 L^2 \mathbb{E} \left\| w_{i,k-1} - w + \rho \frac{\nabla F_i(w_{i,k-1}, \zeta_i)}{\|\nabla F_i(w_{i,k-1}, \zeta_i)\|} - \rho \frac{\nabla F(w)}{\|\nabla F(w)\|} \right\|^2 + 2K^2 N^2 L^2 c^2 \rho^2 \epsilon^2 \\
 &\leq 4(1-c)^2 K^2 N^2 L^2 \mathbb{E} \left[\|w_{i,k-1} - w\|^2 + \rho^2 \|\delta_{i,k-1} - \delta\|^2 \right] + 2K^2 N^2 L^2 c^2 \rho^2 \epsilon^2 \tag{39}
 \end{aligned}$$

$$\leq 4(1-c)^2 K^2 N^2 L^2 \mathbb{E} \|w_{i,k-1} - w\|^2 + 8K^4 N^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 2K^2 N^2 L^2 c^2 \rho^2 \epsilon^2 \tag{40}$$

$$\begin{aligned}
 &\leq 4(1-c)^2 K^2 N^2 L^2 [Ke(T + 8K\eta_l^2 \mathbb{E}\|\nabla F(\tilde{w})\|^2 + 16KL^2 \eta_l^2 c^2 \rho^2 \epsilon^2)] + 8K^4 N^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 2K^2 N^2 L^2 c^2 \rho^2 \epsilon^2 \\
 &= 4e(1-c)^2 K^3 N^2 L^2 T + 32e(1-c)^2 K^4 \eta_l^2 N^2 L^2 \|\nabla F(\tilde{w}^r)\|^2 + 64ec^2(1-c)^2 N^2 K^3 L^4 \eta_l^2 \rho^2 \epsilon^2 \\
 &\quad + 8K^4 N^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 2K^2 N^2 L^2 c^2 \rho^2 \epsilon^2 \tag{41}
 \end{aligned}$$

where Eq. (37) is from Assumption 1, Eq. (38) is from the definition of ϵ , Eq. (39) is from the definition of $\delta := \nabla F(w)/\|\nabla F(w)\|$, $\delta_{i,k-1} := \nabla F_i(w_{i,k-1}, \zeta_i)/\|\nabla F_i(w_{i,k-1}, \zeta_i)\|$, Eq. (41) is from Lemma 5, and $T = 8K\eta_l^2 \mathbb{E} [L^2 \sigma_l^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2]$. For Eq. (40), we borrow the result of Lemma B.1 in (Qu et al., 2022).

By substituting Eq. (41) to Eq. (36), it becomes

$$\begin{aligned}
 &\frac{\eta_l K}{2} \|\nabla F(\tilde{w}^r)\|^2 + \frac{\eta_l}{2N^2 K} \mathbb{E} \left\| \sum_{i,k} [\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)] \right\|^2 - \frac{\eta_l}{2N^2 K} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \\
 &\leq \left(\frac{\eta_l K}{2} + 16e(1-c)^2 K^3 \eta_l^3 L^2 \right) \|\nabla F(\tilde{w}^r)\|^2 + 2e(1-c)^2 K^2 L^2 \eta_l T \\
 &\quad + 32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \rho^2 \epsilon^2 + 4K^3 L^4 \eta_l^3 \rho^2 (1-c)^2 + K\eta_l L^2 c^2 \rho^2 \epsilon^2 - \frac{\eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2.
 \end{aligned}$$

□

Lemma 7. For the full client participation scheme, we can bound $\mathbb{E} [\|\Delta^r\|^2]$ as follows:

$$\mathbb{E} [\|\Delta^r\|^2] \leq 2L^2(1-c)^2 \rho^2 \eta_l^2 \frac{K\sigma_l^2}{N} + 2\frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2$$

Proof. For the full client participation scheme, we have:

$$\mathbb{E} \left[\|\Delta^r\|^2 \right] = \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r, \zeta_i) \right\|^2 \quad (42)$$

$$\begin{aligned} &\leq 2 \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} [\nabla F_i(\bar{w}_{i,k}^r, \zeta_i) - \nabla F_i(\bar{w}_{i,k}^r)] \right\|^2 + 2 \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \\ &\leq 2 \frac{\eta_l^2}{N^2} \sum_{i,k} \mathbb{E} \|\nabla F_i(\bar{w}_{i,k}^r, \zeta_i) - \nabla F_i(\bar{w}_{i,k}^r)\|^2 + 2 \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \end{aligned} \quad (43)$$

$$\leq 2L^2 \frac{\eta_l^2}{N^2} \sum_{i,k} \mathbb{E} \|(\bar{w}_{i,k}^r, \zeta_i) - \bar{w}_{i,k}^r\|^2 + 2 \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \quad (44)$$

$$\begin{aligned} &= 2L^2(1-c)^2 \rho^2 \frac{\eta_l^2}{N^2} \sum_{i,k} \mathbb{E} \left\| \frac{\nabla F_i(w_{i,k}^r, \zeta_i)}{\|\nabla F_i(w_{i,k}^r, \zeta_i)\|} - \frac{\nabla F_i(w_{i,k}^r)}{\|\nabla F_i(w_{i,k}^r)\|} \right\|^2 + 2 \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \\ &\leq 2L^2(1-c)^2 \rho^2 \frac{\eta_l^2}{N^2} \frac{K\sigma_l^2}{N} + 2 \frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\bar{w}_{i,k}^r) \right\|^2, \end{aligned} \quad (45)$$

where Eq. (43) is from the Lemma 3, Eq. (44) is from the Assumption 1, and Eq. (45) is from Assumption 3. \square

Lemma 8. (Descent Lemma). For all $r \leq R - 1$ and $i \in \mathcal{S}^r$, with the choice of learning rate, the iterates generated by FedGF in Algorithm 1 satisfy:

$$\mathbb{E} [F(\tilde{w}^{r+1})] \leq F(\tilde{w}^r) + \|\nabla F(\tilde{w}^r)\|^2 \left(-\frac{K\eta_g\eta_l}{2} + 16eK^3\eta_g\eta_l^3L^2(1-c)^2 \right) + \Phi_0$$

where the global and local learning rate satisfies $\eta_g\eta_l \leq \frac{1}{2KL}$, $\Phi_0 = 2e(1-c)^2K^2L^2\eta_l\eta_gT + 32ec^2(1-c)^2K^2L^4\eta_l^3\eta_g\rho^2\epsilon^2 + 4K^3L^4\eta_l^3\eta_g\rho^2(1-c)^2 + K\eta_l\eta_gL^2c^2\epsilon^2 + L^3(1-c)^2\rho^2\eta_g^2\eta_l^2\frac{K\sigma_l^2}{N}$, and $T = 8K\eta_l^2 [L^2\sigma_l^2(1-c)^2\rho^2 + 8K^2L^4\eta_l^2\rho^2(1-c)^2 + 3\sigma_g^2 + 6L^2\rho^2]$.

Proof. We firstly propose the proof of full client participation scheme. Due to the smoothness in Assumption 1, $F(\tilde{w}) = \max_{\|\epsilon\| \leq \rho} F(w + \epsilon)$, and taking expectation of $F(\tilde{w}^{r+1})$ over the randomness at communication round r , we have:

$$\mathbb{E} [F(w^{r+1})] \leq \mathbb{E} [F(\tilde{w}^{r+1})] \leq F(\tilde{w}^r) + \underbrace{\mathbb{E} [\langle \nabla F(\tilde{w}^r), \tilde{w}^{r+1} - \tilde{w}^r \rangle]}_A + \underbrace{L/2 \mathbb{E} \|\tilde{w}^{r+1} - \tilde{w}^r\|_2^2}_B \quad (46)$$

For A , we borrow the approximation in (Qu et al., 2022).

$$\langle \nabla F(\tilde{w}^r), \tilde{w}^{r+1} - \tilde{w}^r \rangle \approx \langle \nabla F(\tilde{w}^r), w^{r+1} - w^r \rangle \quad (47)$$

$$\begin{aligned} \langle \nabla F(\tilde{w}^r), \mathbb{E} [w^{r+1} - w^r] \rangle &= \langle \nabla F(\tilde{w}^r), -\frac{\eta_g}{N} \sum_i \sum_k^K \eta_l \mathbb{E} [\nabla F_i(\bar{w}_{i,k}^r)] \rangle \\ &= \eta_g \langle \nabla F(\tilde{w}^r), -\frac{1}{N} \sum_i \sum_k^K \eta_l \mathbb{E} [\nabla F_i(\bar{w}_{i,k}^r)] + K\eta_l \nabla F(\tilde{w}^r) - K\eta_l \nabla F(\tilde{w}^r) \rangle \\ &= \eta_g \underbrace{\langle \nabla F(\tilde{w}^r), -\frac{1}{N} \sum_i \sum_k^K \eta_l \mathbb{E} [\nabla F_i(\bar{w}_{i,k}^r)] + K\eta_l \nabla F(\tilde{w}^r) \rangle}_{A_1} - \eta_g \underbrace{\langle \nabla F(\tilde{w}^r), K\eta_l \nabla F(\tilde{w}^r) \rangle}_{A_2} \end{aligned}$$

For A_1 , we refer to Lemma 6.

$$\begin{aligned} \left\langle \nabla F(\tilde{w}^r), -\frac{1}{N} \sum_i \sum_k^K \eta_l \mathbb{E}[\nabla F_i(\tilde{w}_{i,k}^r)] + K\eta_l \nabla F(\tilde{w}^r) \right\rangle &\leq \left(\frac{\eta_l K}{2} + 16e(1-c)^2 K^3 \eta_l^3 L^2 \right) \|\nabla F(\tilde{w}^r)\|^2 \\ &+ 2e(1-c)^2 K^2 L^2 \eta_l T + 32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \rho^2 \epsilon^2 + 4K^3 L^4 \eta_l^3 \rho^2 (1-c)^2 \\ &+ K\eta_l L^2 c^2 \rho^2 \epsilon^2 - \frac{\eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2. \end{aligned}$$

For A_2 ,

$$\langle \nabla F(\tilde{w}^r), K\eta_l \nabla F(\tilde{w}^r) \rangle = \eta_l K \|\nabla F(\tilde{w}^r)\|^2.$$

Therefore, A becomes

$$\begin{aligned} \langle \nabla F(\tilde{w}^r), \mathbb{E}[w^{r+1} - w^r] \rangle &\leq \|\nabla F(\tilde{w}^r)\|^2 \left(-\frac{K\eta_g \eta_l}{2} + 16eK^3 \eta_g \eta_l^3 L^2 (1-c)^2 \right) \\ &- \frac{\eta_g \eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 + 2e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 \\ &+ 4K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2. \end{aligned}$$

By referring the approximation in (Qu et al., 2022), B becomes,

$$\|\tilde{w}^{r+1} - \tilde{w}^r\|_2^2 \approx \eta_g^2 \|\Delta^r\|^2$$

Therefore,

$$\begin{aligned} \mathbb{E}[F(w^{r+1})] &\leq \mathbb{E}[F(\tilde{w}^{r+1})] \\ &\leq F(\tilde{w}^r) + \mathbb{E}[\langle \tilde{w}^{r+1} - \tilde{w}^r, \nabla F(\tilde{w}^r) \rangle] + L/2 \mathbb{E}\|\tilde{w}^{r+1} - \tilde{w}^r\|_2^2 \end{aligned}$$

$$\begin{aligned}
 &\leq F(\tilde{w}^r) + \|\nabla F(\tilde{w}^r)\|^2 \left(\frac{-K\eta_g\eta_l}{2} + 16eK^3\eta_g\eta_l^3L^2(1-c)^2 \right) \\
 &\quad - \frac{\eta_g\eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 + 2e(1-c)^2K^2L^2\eta_l\eta_gT + 32ec^2(1-c)^2K^2L^4\eta_l^3\eta_g\rho^2\epsilon^2 + 4K^3L^4\eta_l^3\eta_g\rho^2(1-c)^2 \\
 &\quad + K\eta_l\eta_gL^2c^2\rho^2\epsilon^2 + \frac{L\eta_g^2}{2} \mathbb{E} \|\Delta^r\|^2 \\
 &\leq F(\tilde{w}^r) + \|\nabla F(\tilde{w}^r)\|^2 \left(\frac{-K\eta_g\eta_l}{2} + 16eK^3\eta_g\eta_l^3L^2(1-c)^2 \right) \\
 &\quad - \frac{\eta_g\eta_l}{2KN^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 + 2e(1-c)^2K^2L^2\eta_l\eta_gT + 32ec^2(1-c)^2K^2L^4\eta_l^3\eta_g\rho^2\epsilon^2 + 4K^3L^4\eta_l^3\eta_g\rho^2(1-c)^2 \\
 &\quad + K\eta_l\eta_gL^2c^2\rho^2\epsilon^2 + \frac{L\eta_g^2}{2} \left[2L^2(1-c)^2\rho^2\eta_l^2\frac{K\sigma_l^2}{N} + 2\frac{\eta_l^2}{N^2} \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 \right] \tag{48}
 \end{aligned}$$

$$\begin{aligned}
 &= F(\tilde{w}^r) + \|\nabla F(\tilde{w}^r)\|^2 \left(\frac{-K\eta_g\eta_l}{2} + 16eK^3\eta_g\eta_l^3L^2(1-c)^2 \right) \\
 &\quad + 2e(1-c)^2K^2L^2\eta_l\eta_gT + 32ec^2(1-c)^2K^2L^4\eta_l^3\eta_g\rho^2\epsilon^2 + 4K^3L^4\eta_l^3\eta_g\rho^2(1-c)^2 \\
 &\quad + K\eta_l\eta_gL^2c^2\rho^2\epsilon^2 + L^3(1-c)^2\rho^2\eta_g^2\eta_l^2\frac{K\sigma_l^2}{N} + \mathbb{E} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 \left(\frac{L\eta_g^2\eta_l^2}{N^2} - \frac{\eta_g\eta_l}{2KN^2} \right) \\
 &\leq F(\tilde{w}^r) + \|\nabla F(\tilde{w}^r)\|^2 \left(\frac{-K\eta_g\eta_l}{2} + 16eK^3\eta_g\eta_l^3L^2(1-c)^2 \right) + 2e(1-c)^2K^2L^2\eta_l\eta_gT \\
 &\quad + 32ec^2(1-c)^2K^2L^4\eta_l^3\eta_g\rho^2\epsilon^2 + 4K^3L^4\eta_l^3\eta_g\rho^2(1-c)^2 + K\eta_l\eta_gL^2c^2\rho^2\epsilon^2 + L^3(1-c)^2\rho^2\eta_g^2\eta_l^2\frac{K\sigma_l^2}{N} \tag{49}
 \end{aligned}$$

where Eq. (48) is from Lemma 7, Eq. (49) is from $\eta_g\eta_l \leq \frac{1}{2KL}$, and $T = 8K\eta_l^2[L^2\sigma_l^2(1-c)^2\rho^2 + 8K^2L^4\eta_l^2\rho^2(1-c)^2 + 3\sigma_g^2 + 6L^2\rho^2]$. \square

Theorem 1. (A concretet description of the full participating client case of Theorem 1 in the main paper) Let constant local and global learning rates η_l and η_g be chosen as such that $\eta_g\eta_l \leq \frac{1}{2KL}$. Under Assumption 1-3 and with full client participation, and if we choose the learning rates $\eta_l = \mathcal{O}(\frac{1}{\sqrt{RKL}})$, $\eta_g = \sqrt{KN}$ and perturbation amplitude ρ proportional to the learning rate, i.e., $\rho = \mathcal{O}(\frac{1}{\sqrt{R}})$, the average of the norm of the gradient generated by FedGF satisfies:

$$\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} [\|\nabla F(w^{r+1})\|] = \mathcal{O} \left(\frac{FL}{\sqrt{RK\eta_l}} + \frac{(1-c)^2}{R} \sigma_g^2 + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KN}} \sigma_l^2 + \frac{L^2c^2\epsilon^2}{R} \right)$$

where $F = F(\tilde{w}^0) - F(\tilde{w}^*)$ and $F(\tilde{w}^*) = \min_{\tilde{w}} F(\tilde{w})$.

Proof. For full client participation, summing the result of Lemma 8 for $r \leq R-1$ and multiplying both sides by $\frac{1}{CK\eta_lR}$ with $(\frac{1}{2} - 16eK^2\eta_l^2L^2(1-c)^2) > C > 0$, we have

$$\begin{aligned}
 \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla F(\tilde{w}^r)\|^2 &\leq \frac{F(\tilde{w}^0) - \mathbb{E}[F(\tilde{w}^R)]}{CK\eta_g\eta_l R} + \Phi \\
 &\leq \frac{F(\tilde{w}^0) - F(\tilde{w}^*)}{CK\eta_g\eta_l R} + \Phi \\
 &= \frac{F}{CK\eta_g\eta_l R} + \Phi,
 \end{aligned}$$

where $\Phi = \frac{1}{CK\eta_g\eta_l} \left[2e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 4K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + L^3(1-c)^2 \rho^2 \eta_g^2 \eta_l^2 \frac{K\sigma_l^2}{N} \right]$.

If we choose the learning rates $\eta_l = \mathcal{O}(\frac{1}{\sqrt{RKL}})$, $\eta_g = \sqrt{KN}$ and perturbation amplitude ρ proportional to the learning rate, i.e., $\rho = \mathcal{O}(\frac{1}{\sqrt{R}})$, we have

$$\begin{aligned}
 \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla F(\tilde{w}^r)\|^2 &\leq \frac{F}{CK\eta_g\eta_l R} + \Phi \\
 &= \frac{F}{CK\eta_g\eta_l R} + \frac{1}{CK\eta_g\eta_l} \left[2e(1-c)^2 K^2 L^2 \eta_l \eta_g T \right. \\
 &\quad \left. + 32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 4K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + L^3(1-c)^2 \rho^2 \eta_g^2 \eta_l^2 \frac{K\sigma_l^2}{N} \right] \\
 &= \frac{F}{CK\eta_g\eta_l R} + \frac{16e(1-c)^2 K^2 L^2 \eta_l^2}{C} \left[L^2 \sigma_l^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2 \right] \\
 &\quad + \frac{1}{CK\eta_g\eta_l} \left[32ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 4K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + L^3(1-c)^2 \rho^2 \eta_g^2 \eta_l^2 \frac{K\sigma_l^2}{N} \right] \\
 &= \mathcal{O} \left(\frac{FL}{\sqrt{RKN}} + \frac{L^2(1-c)^4 \sigma_l^2}{R^2} + \frac{(1-c)^4 L^2}{R^3} + \frac{(1-c)^2 \sigma_g^2}{R} + \frac{(1-c)^2 L^2}{R^2} \right. \\
 &\quad \left. + \frac{L^2(1-c)^2 c^2 \epsilon^2}{KR^2} + \frac{L^2(1-c)^2}{R^2} + \frac{L^2 c^2 \epsilon^2}{R} + \frac{L^2(1-c)^2 \sigma_l^2}{R\sqrt{RKN}} \right)
 \end{aligned}$$

Note that the $\frac{(1-c)^2}{R} \sigma_g^2$ is caused by the heterogeneity between clients, $\frac{L^2(1-c)^4}{R^3} \sigma_l^2$, $\frac{L^2(1-c)^2}{R\sqrt{RKN}} \sigma_l^2$ are due to the perturbation of stochastic gradient and $\frac{(1-c)^4 L^2}{R^2}$, $\frac{L^2(1-c)^2}{R^2}$ are due to the local SAM, $\frac{L^2(1-c)^2 c^2 \epsilon^2}{KR^2}$, $\frac{L^2 c^2 \epsilon^2}{R}$ is caused by approximation of global perturbation, i.e., $\epsilon := \|\Delta^r / \|\Delta^r\| - \nabla F(w^r) / \|F(w^r)\|\|$. After omitting the higher order, i.e., $\left(\frac{L^2(1-c)^4 \sigma_l^2}{R^2}, \frac{(1-c)^4 L^2}{R^3}, \frac{(1-c)^2 L^2}{R^2}, \frac{(1-c)^2 L^2}{R^2}, \frac{L^2(1-c)^2 c^2 \epsilon^2}{KR^2} \right)$,

$$\frac{1}{R} \sum_{r=0}^{R-1} \|\nabla F(\tilde{w}^r)\|^2 = \mathcal{O} \left(\frac{FL}{\sqrt{RKN}} + \frac{(1-c)^2 \sigma_g^2}{R} + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KN}} \sigma_l^2 + \frac{L^2 c^2 \epsilon^2}{R} \right) \quad (50)$$

From the perturbation amplitude $\rho = \mathcal{O}(\frac{1}{\sqrt{R}})$, we obtain

$$\|\nabla F(\tilde{w}^r) - \nabla F(w^r)\|^2 \leq L^2 \|\tilde{w}^r - w^r\|^2 \leq L^2 \rho^2 = \mathcal{O} \left(\frac{L^2}{R} \right). \quad (51)$$

It indicates that the difference $\|\nabla F(\tilde{w}^r) - \nabla F(w^r)\|^2$ gets smaller as R increases, therefore, the convergence rate of $\|\nabla F(w^r)\|^2$ is also converge with $\|\nabla F(\tilde{w}^r)\|^2$. \square

A.4. Partial Client Participating Convergence

Lemma 9. For the partial client participation scheme, we can bound $\mathbb{E} \|\Delta^r\|^2$ as follows:

$$\mathbb{E} \left[\|\Delta^r\|^2 \right] \leq 2K\eta_l^2 L^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} + 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2.$$

Proof. For the partial client participation scheme, we have:

$$\begin{aligned} \mathbb{E} \|\Delta^r\|^2 &= \frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r, \zeta_i) \right\|^2 \\ &\leq 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k [\nabla F_i(\bar{w}_{i,k}^r, \zeta_i) - \nabla F_i(\bar{w}_{i,k}^r)] \right\|^2 + 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \end{aligned} \quad (52)$$

$$= 2\frac{\eta_l^2}{S^2} \sum_{i \in S^r} \sum_k \mathbb{E} \|\nabla F_i(\bar{w}_{i,k}^r, \zeta_i) - \nabla F_i(\bar{w}_{i,k}^r)\|^2 + 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \quad (53)$$

$$\leq 2\frac{\eta_l^2 L^2}{S^2} \sum_{i \in S^r} \sum_k \mathbb{E} \|(\bar{w}_{i,k}^r, \zeta_i) - \bar{w}_{i,k}^r\|^2 + 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \quad (54)$$

$$\begin{aligned} &= 2\frac{\eta_l^2 L^2 \rho^2 (1-c)^2}{S^2} \sum_{i \in S^r} \sum_k \mathbb{E} \left\| \frac{\nabla F_i(w_{i,k}^r, \zeta_i)}{\|\nabla F_i(w_{i,k}^r, \zeta_i)\|} - \frac{\nabla F_i(w_{i,k}^r)}{\|\nabla F_i(w_{i,k}^r)\|} \right\|^2 + 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \\ &\leq 2K\eta_l^2 L^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} + 2\frac{\eta_l^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2, \end{aligned} \quad (55)$$

where Eq. (52) is from Lemma 2, Eq. (53) is from Lemma 3, Eq. (54) is from Assumption 1 and Eq. (55) is from Assumption 3. \square

Lemma 10.

$$\begin{aligned} &\langle \nabla F(\bar{w}^r), -\frac{1}{S} \sum_{i \in S^r} \sum_k \eta_l \nabla F_i(\bar{w}_{i,k}^r) + K\eta_l \nabla f(\bar{w}^r) \rangle \leq \\ &\left(\frac{\eta_l K}{2} + 32e(1-c)^2 K^3 \eta_l^3 L^2 \right) \|\nabla F(\bar{w}^r)\|^2 + 4e(1-c)^2 K^2 L^2 \eta_l T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \epsilon^2 \\ &\quad + 8K^3 L^4 \eta_l^3 \rho^2 (1-c)^2 + 2K\eta_l L^2 c^2 \rho^2 \epsilon^2 + 2\eta_l K(3\sigma_g^2 + 6L^2 \rho^2) - \frac{\eta_l}{2KS^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \end{aligned}$$

where $T = 8K\eta_l^2 [L^2 \sigma_l^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2]$.

Proof.

$$\begin{aligned} &\langle \nabla F(\bar{w}^r), -\frac{1}{S} \sum_{i \in S^r} \sum_k \eta_l \mathbb{E}[\nabla F_i(\bar{w}_{i,k}^r)] + K\eta_l \nabla f(\bar{w}^r) \rangle \\ &= \frac{\eta_l K}{2} \|\nabla F(\bar{w}^r)\|^2 + \frac{\eta_l}{2S^2 K} \underbrace{\mathbb{E} \left\| \sum_{i \in S^r} \sum_k [\nabla F_i(\bar{w}_{i,k}^r) - \nabla F(\bar{w}^r)] \right\|^2}_A - \frac{\eta_l}{2S^2 K} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \end{aligned} \quad (56)$$

where Eq. (36) is from $\langle a, b \rangle = \frac{1}{2}(\|a\|^2 + \|b\|^2 - \|a - b\|^2)$ with $a = \sqrt{\eta_l K} \nabla F(\bar{w}^r)$ and $b = -\frac{\sqrt{\eta_l}}{S\sqrt{K}} \sum_{i \in S^r} \sum_k (\nabla F_i(\bar{w}_{i,k}^r) - \nabla F(\bar{w}^r))$.

For A,

$$\begin{aligned}
 & \mathbb{E} \left\| \sum_{i \in S^r} \sum_k [\nabla F_i(\bar{w}_{i,k}^r) - \nabla F(\bar{w}^r)] \right\|^2 \\
 & \leq SK \sum_{i \in S^r} \sum_k \mathbb{E} \|\nabla F_i(\bar{w}_{i,k}^r) - \nabla F(\bar{w}^r)\|^2 \\
 & = SK \sum_{i \in S^r} \sum_k \mathbb{E} \|\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\bar{w}^r) + \nabla F_i(\bar{w}^r) - \nabla F(\bar{w}^r)\|^2 \\
 & \leq 2SK \sum_{i \in S^r} \sum_k \mathbb{E} [\|\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\bar{w}^r)\|^2 + \|\nabla F_i(\bar{w}^r) - \nabla F(\bar{w}^r)\|^2] \\
 & \leq 2SKL^2 \sum_{i \in S^r} \sum_k \mathbb{E} \|\bar{w}_{i,k}^r - \bar{w}^r\|^2 + 2SK \sum_{i \in S^r} \sum_k \|\nabla F_i(\bar{w}^r) - \nabla F(\bar{w}^r)\|^2 \\
 & \leq 8e(1-c)^2 K^3 S^2 L^2 T + 64e(1-c)^2 K^4 \eta_l^2 S^2 L^2 \|\nabla F(\bar{w}^r)\|^2 + 128ec^2(1-c)^2 S^2 K^3 L^4 \eta_l^2 \rho^2 \epsilon^2 \\
 & \quad + 16K^4 S^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 4K^2 S^2 L^2 c^2 \rho^2 \epsilon^2 + 2SK \sum_{i \in S^r} \sum_k \|\nabla F_i(\bar{w}^r) - \nabla F(\bar{w}^r)\|^2 \tag{57}
 \end{aligned}$$

$$\begin{aligned}
 & \leq 8e(1-c)^2 K^3 S^2 L^2 T + 64e(1-c)^2 K^4 \eta_l^2 S^2 L^2 \|\nabla F(\bar{w}^r)\|^2 + 128ec^2(1-c)^2 S^2 K^3 L^4 \eta_l^2 \rho^2 \epsilon^2 \\
 & \quad + 16K^4 S^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 4K^2 S^2 L^2 c^2 \rho^2 \epsilon^2 + 2S^2 K^2 (3\sigma_g^2 + 6L^2 \rho^2) \tag{58}
 \end{aligned}$$

where Eq. (57) is from the Eq. (37) in Lemma 6, Eq. (58) is from Lemma 4, $T = 8K\eta_l^2 \mathbb{E}[L^2 \sigma_l^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2]$.

By substitute Eq. (58) in Eq. (56), it becomes

$$\begin{aligned}
 & \frac{\eta_l K}{2} \|\nabla f(\bar{w}^r)\|^2 + \frac{\eta_l}{2S^2 K} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k [\nabla F_i(\bar{w}_{i,k}^r) - \nabla F_i(\bar{w}^r)] \right\|^2 - \frac{\eta_l}{2S^2 K} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2 \\
 & \leq \left(\frac{\eta_l K}{2} + 32e(1-c)^2 K^3 \eta_l^3 L^2 \right) \|\nabla F(\bar{w}^r)\|^2 + 4e(1-c)^2 K^2 L^2 \eta_l T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \rho^2 \epsilon^2 \\
 & \quad + 8K^3 L^4 \eta_l^3 \rho^2 (1-c)^2 + 2K\eta_l L^2 c^2 \rho^2 \epsilon^2 + \eta_l K (3\sigma_g^2 + 6L^2 \rho^2) - \frac{\eta_l}{2KS^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k \nabla F_i(\bar{w}_{i,k}^r) \right\|^2
 \end{aligned}$$

□

Lemma 11. (Descent Lemma). For all $r \leq R - 1$ and $i \in S^r$, with the choice of learning rate, the iterates generated by FedGF in Algorithm 1 satisfy:

$$\mathbb{E} [F(\bar{w}^{r+1})] \leq F(\bar{w}^r) + \|\nabla F(\bar{w}^r)\|^2 \left(\frac{-K\eta_g \eta_l}{2} + 24eK^3 \eta_g \eta_l^3 L^2 (1-c)^2 \right) + \Phi_0$$

where the global and local learning rate satisfies $\eta_g \eta_l \leq \frac{1}{4KL}$, the expectation is w.r.t. the stochasticity of the algorithm, and $\Phi_0 = 4e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + \eta_l \eta_g K (3\sigma_g^2 + 6L^2 \rho^2) + 2K\eta_l^2 L^3 \eta_g^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S}$.

Proof. We propose the proof of a partial client participation scheme. We denote briefly because it has a similar flow with proof of full participation. Due to the smoothness in Assumption 1, $F(\bar{w}) = \max_{\|\epsilon\| \leq \rho} F(\bar{w} + \epsilon)$, and taking expectation of $F(\bar{w}^{r+1})$ over the randomness at communication round r , we have:

$$\mathbb{E} [F(\bar{w}^{r+1})] \leq \mathbb{E} [F(\bar{w}^{r+1})] \leq F(\bar{w}^r) + \underbrace{\mathbb{E} \langle \bar{w}^{r+1} - \bar{w}^r, \nabla F(\bar{w}^r) \rangle}_A + \underbrace{L/2 \mathbb{E} \|\bar{w}^{r+1} - \bar{w}^r\|_2^2}_B \tag{59}$$

For A , from the approximation in (Qu et al., 2022),

$$\langle \nabla F(\tilde{w}^r), \tilde{w}^{r+1} - \tilde{w}^r \rangle \approx \langle \nabla F(\tilde{w}^r), w^{r+1} - w^r \rangle$$

$$\langle \nabla F(\tilde{w}^r), \mathbb{E}[w^{r+1} - w^r] \rangle = \underbrace{\eta_g \langle \nabla F(\tilde{w}^r), -\frac{1}{S} \sum_{i \in S^r} \sum_k^K \eta_l \mathbb{E}[\nabla F_i(\tilde{w}_{i,k}^r)] + K\eta_l \nabla F(\tilde{w}^r) \rangle}_{A_1} - \underbrace{\eta_g \langle \nabla F(\tilde{w}^r), K\eta_l \nabla F(\tilde{w}^r) \rangle}_{A_2}$$

For A_1 , we refer the Lemma 10.

$$\begin{aligned} & \langle \nabla F(\tilde{w}^r), -\frac{1}{S} \sum_{i \in S^r} \sum_k^K \eta_l \mathbb{E}[\nabla F_i(\tilde{w}_{i,k}^r)] + K\eta_l \nabla F(\tilde{w}^r) \rangle \leq \\ & \left(\frac{\eta_l K}{2} + 32e(1-c)^2 K^3 \eta_l^3 L^2 \right) \|\nabla F(\tilde{w}^r)\|^2 + 4e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 \\ & + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + \eta_l \eta_g K(3\sigma_g^2 + 6L^2 \rho^2) - \frac{\eta_l \eta_g}{2KS^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k^K \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 \end{aligned}$$

For A_2 ,

$$\langle \nabla F(\tilde{w}^r), K\eta_l \nabla F(\tilde{w}^r) \rangle = \eta_l K \|\nabla F(\tilde{w}^r)\|^2.$$

Therefore, A becomes

$$\begin{aligned} & \langle \nabla F(\tilde{w}^r), \mathbb{E}[w^{r+1} - w^r] \rangle \leq \\ & \left(\frac{-\eta_l \eta_g K}{2} + 32eK^3 L^2 \eta_l^3 \eta_g (1-c)^2 \right) \|\nabla F(\tilde{w}^r)\|^2 + 4e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 \\ & + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + \eta_l \eta_g K(3\sigma_g^2 + 6L^2 \rho^2) - \frac{\eta_l}{2KS^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k^K \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2. \end{aligned}$$

B becomes,

$$\begin{aligned} L/2\mathbb{E}\|\tilde{w}^{r+1} - \tilde{w}^r\|_2^2 & \approx L\eta_g^2 \mathbb{E}\|\Delta^r\|^2 \\ & \leq 2K\eta_l^2 L^3 \eta_g^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} + 2\frac{L\eta_l^2 \eta_g^2}{S^2} \mathbb{E} \left\| \sum_{i \in S^r} \sum_k^K \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 \end{aligned} \quad (60)$$

where Eq. (60) is from Lemma 9

Therefore,

$$\begin{aligned} \mathbb{E}[F(w^{r+1})] & \leq \mathbb{E}[F(\tilde{w}^{r+1})] \\ & \leq F(\tilde{w}^r) + \langle \nabla F(\tilde{w}^r), \mathbb{E}[\tilde{w}^{r+1} - \tilde{w}^r] \rangle + L/2\mathbb{E}\|\tilde{w}^{r+1} - \tilde{w}^r\|_2^2 \\ & \leq \|\nabla F(\tilde{w}^r)\|^2 \left(\frac{-K\eta_g \eta_l}{2} + 32eK^3 \eta_g \eta_l^3 L^2 (1-c)^2 \right) + \left(2\frac{L\eta_l^2 \eta_g^2}{S^2} - \frac{\eta_g \eta_l}{2KS^2} \right) \mathbb{E} \left\| \sum_{i \in S^r} \sum_k^K \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2 \\ & \quad + 4e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 \\ & \quad + \eta_l \eta_g K(3\sigma_g^2 + 6L^2 \rho^2) + 2K\eta_l^2 L^3 \eta_g^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} \\ & \leq \|\nabla F(\tilde{w}^r)\|^2 \left(\frac{-K\eta_g \eta_l}{2} + 32eK^3 \eta_g \eta_l^3 L^2 (1-c)^2 \right) \\ & \quad + 4e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 \\ & \quad + \eta_l \eta_g K(3\sigma_g^2 + 6L^2 \rho^2) + 2K\eta_l^2 L^3 \eta_g^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} \end{aligned} \quad (61)$$

where Eq. (61) is from $\eta_g \eta_l \leq \frac{1}{4KL}$, and $T = 8K\eta_l^2 \mathbb{E}[L^2 \sigma_l^2 (1-c)^2 \rho^2 + 8K^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 3\sigma_g^2 + 6L^2 \rho^2]$. \square

Theorem 1. (A concretet description of the partial participating client case of Theorem 1 in the main paper) Let constant size of perturbation, local and global learning rates, i.e., ρ, η_l and η_g , be chosen as such that $(1-c)\eta_l \leq \frac{1}{8KL}, \eta_g \eta_l \leq \frac{1}{4KL}$. Under Assumption 1-3 and with partial client participation, and if we choose the learning rates $\eta_l = \frac{1}{\sqrt{RKL}}, \eta_g = \sqrt{KS}$ and perturbation amplitude ρ proportional to the learning rate, e.g., $\rho = \mathcal{O}\left(\frac{1}{\sqrt{R}}\right)$, the average of norm of gradient generated by FedGF satisfies:

$$\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} [\|\nabla F(w^{r+1})\|] = \mathcal{O}\left(\frac{FL}{\sqrt{RKS}} + \left(\frac{(1-c)^2}{R} + 1\right) \sigma_g^2 + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KS}} \sigma_l^2 + \frac{L^2(c^2\epsilon^2 + 1)}{R}\right)$$

Proof. For partial client participation, summing the result of Lemma 11 for $r \leq R-1$ and multiplying both sides by $\frac{1}{CK\eta_l R}$ with $\left(\frac{1}{2} - 32eK^2\eta_l^2 L^2(1-c)^2\right) > C > 0$, we have

$$\begin{aligned} \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla F(\tilde{w}^r)\|^2 &\leq \frac{F(\tilde{w}^0) - \mathbb{E}[F(\tilde{w}^R)]}{CK\eta_g\eta_l R} + \Phi \\ &\leq \frac{F(\tilde{w}^0) - F^*}{CK\eta_g\eta_l R} + \Phi, \end{aligned}$$

where $\Phi = \frac{1}{CK\eta_g\eta_l} \left[4e(1-c)^2 K^2 L^2 \eta_l \eta_g T + 64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \rho^2 \epsilon^2 + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 + \eta_l \eta_g K(3\sigma_g^2 + 6L^2 \rho^2) + 2K\eta_l^2 L^3 \eta_g^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} \right]$.

If we choose the learning rates $\eta_l = \mathcal{O}\left(\frac{1}{\sqrt{RKL}}\right), \eta_g = \sqrt{KS}$ and perturbation amplitude ρ proportional to the learning rate, e.g., $\rho = \mathcal{O}\left(\frac{1}{\sqrt{R}}\right)$, we have

$$\begin{aligned} \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla F(\tilde{w}^r)\|^2 &\leq \frac{F(\tilde{w}^0) - \mathbb{E}[F(\tilde{w}^R)]}{CK\eta_g\eta_l R} + \Phi \\ &= \frac{F(\tilde{w}^0) - \mathbb{E}[F(\tilde{w}^R)]}{CK\eta_g\eta_l R} + \frac{4e(1-c)^2 K^2 L^2 \eta_l \eta_g T}{CK\eta_g\eta_l} \\ &\quad + \frac{1}{CK\eta_g\eta_l} \left[64ec^2(1-c)^2 K^2 L^4 \eta_l^3 \eta_g \epsilon^2 + 8K^3 L^4 \eta_l^3 \eta_g \rho^2 (1-c)^2 + 2K\eta_l \eta_g L^2 c^2 \rho^2 \epsilon^2 \right. \\ &\quad \left. + \eta_l \eta_g K(3\sigma_g^2 + 6L^2 \rho^2) + 2K\eta_l^2 L^3 \eta_g^2 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} \right] \\ &= \frac{F}{CK\eta_g\eta_l R} + \frac{4e(1-c)^2 KL^2}{C} (8K\eta_l^2 [L^2 \sigma_l^2 (1-c)^2 \rho^2 + 8K^2 (1-c)^2 L^4 \eta_l^2 \rho^2 + 3\sigma_g^2 + 6L^2 \rho^2]) \\ &\quad + \frac{1}{C} \left[64ec^2(1-c)^2 KL^4 \eta_l^2 \rho^2 \epsilon^2 + 8K^2 L^4 \eta_l^2 \rho^2 (1-c)^2 + 2L^2 c^2 \rho^2 \epsilon^2 + 3\sigma_g^2 + 6L^2 \rho^2 + 2\eta_l \eta_g L^3 \rho^2 (1-c)^2 \frac{\sigma_l^2}{S} \right] \\ &= \mathcal{O}\left(\frac{FL}{\sqrt{RKS}} + \frac{L^2(1-c)^4 \sigma_l^2}{R^2} + \frac{L^2(1-c)^4}{R^3} + \frac{(1-c)^2 \sigma_g^2}{R} + \frac{(1-c)^2 L^2}{R^2} \right. \\ &\quad \left. + \frac{c^2(1-c)^2 L^2 \epsilon^2}{KR^2} + \frac{L^2(1-c)^2}{R^2} + \frac{L^2 c^2 \epsilon^2}{R} + \sigma_g^2 + \frac{L^2}{R} + \frac{L^2(1-c)^2 \sigma_l^2}{R^{3/2}\sqrt{KS}}\right) \end{aligned}$$

Note that $\left(\frac{(1-c)^2}{R} + 1\right) \sigma_g^2$ are caused by the heterogeneity between clients, $\frac{L^2(1-c)^4}{R^2} \sigma_l^2, \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KS}} \sigma_l^2$ are due to the perturbation of stochastic gradient and $\frac{L^2(1-c)^4}{R^3}, \frac{L^2(1-c)^2}{R^2}, \frac{L^2}{R}$ are due to the local SAM, $\frac{c^2(1-c)^2 L^2 \epsilon^2}{KR^2}, \frac{L^2 c^2 \epsilon^2}{R}$ is caused by approximation of global perturbation, i.e., $\epsilon := \|\Delta^r / \|\Delta^r\| - \nabla F(w^r) / \|\nabla F(w^r)\|\|$. After omitting the higher order, i.e., $\left(\frac{L^2(1-c)^4 \sigma_l^2}{R^2}, \frac{L^2(1-c)^4}{R^3}, \frac{(1-c)^2 L^2}{R^2}, \frac{c^2(1-c)^2 L^2 \epsilon^2}{KR^2}\right)$,

$$\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla F(\tilde{w}^r)\|^2 = O\left(\frac{FL}{\sqrt{RKS}} + \left(\frac{(1-c^2)}{R} + 1\right) \sigma_g^2 + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KS}} \sigma_l^2 + \frac{L^2(c^2\epsilon^2 + 1)}{R}\right) \quad (62)$$

□

A.5. Proofs of Theorems for Flatness Discrepancy

We here provide the proofs of Theorem 2 and 3.

Theorem 2. $\Delta_{\mathcal{F}}$ is upper bounded as follows:

$$\Delta_{\mathcal{F}} \leq \rho \sigma_g^2 + L\rho \sum_{i \in [N]} \|w - w_i\| \quad (63)$$

Proof.

$$\left| \max_{\|\delta\| \leq \rho} F(w + \delta) - F(w) - \left[\sum_{i \in [N]} \frac{m_i}{m} \max_{\|\delta_i\| \leq \rho} F_i(w_i + \delta_i) - F_i(w_i) \right] \right| \approx \left| \rho \|\nabla F(w)\| - \rho \sum_{i \in [N]} \frac{m_i}{m} \|\nabla F_i(w_i)\| \right| \quad (64)$$

$$\begin{aligned} &\leq \rho \sum_{i \in [N]} \frac{m_i}{m} \|\nabla F(w) - \nabla F_i(w_i)\| \\ &= \rho \sum_{i \in [N]} \frac{m_i}{m} \|\nabla F(w) - \nabla F_i(w) + \nabla F_i(w) - \nabla F_i(w_i)\| \\ &\leq \rho \sum_{i \in [N]} \frac{m_i}{m} \|\nabla F(w) - \nabla F_i(w)\| + \|\nabla F_i(w) - \nabla F_i(w_i)\| \\ &\leq \rho \sigma_g^2 + L\rho \sum_{i \in [N]} \frac{m_i}{m} \|w - w_i\| \end{aligned} \quad (65)$$

where Eq. (64) is from first-order approximation, and Eq. (65) is from Assumption (2) □

Therefore, as the non-IIDness increases, i.e., σ_g^2 increases, the upper bound of the sharpness of the global model also increases.

Theorem 3. For FedGF, if we choose $\eta_l = \mathcal{O}\left(\frac{1}{\sqrt{RKL}}\right)$ and $\rho = \mathcal{O}\left(\frac{1}{\sqrt{R}}\right)$, $\eta_g = \sqrt{KN}$, $\Delta_{\mathcal{F}}$ is then bounded as follows:

$$\Delta_{\mathcal{F}} \leq \rho \sigma_g^2 + L\rho \sum_{i \in [N]} \frac{m_i}{m} \|w - w_i\|_2 = \mathcal{O}\left(\frac{\sigma_g^2}{\sqrt{R}} + \frac{L(1-c)^2 \sigma_l^2}{R^{5/2}} + \frac{\sigma_g^2}{LR^{3/2}} + \frac{Lc^2\epsilon^2}{R^{5/2}}\right)$$

Proof. First, we derive the upper bound of $\|w - w_i\|_2$

$$\|w - w_i\|_2 \leq \|w - w_i\|_2^2 \quad (66)$$

$$\leq Ke(T + 8K\eta_l^2 \|\nabla F(\tilde{w})\|^2 + 16KL^2\eta_l^2 c^2 \rho^2 \epsilon^2) \quad (67)$$

If we set $\eta_l = \mathcal{O}\left(\frac{1}{\sqrt{RKL}}\right)$, $\rho = \mathcal{O}\left(\frac{1}{\sqrt{R}}\right)$ used in Theorem 1.

The first term is

$$T = 8K\eta_l^2 [L^2\sigma_l^2(1-c^2)\rho^2 + 8K^2L^4\eta_l^2\rho^2(1-c)^2 + 3\sigma_g^2 + 6L^2\rho^2] \quad (68)$$

$$= \mathcal{O}\left(\frac{1}{RKL^2}\right) \cdot \mathcal{O}\left(\frac{L^2\sigma_l^2(1-c)^2}{R} + \frac{(1-c)^2L^2}{R^2} + \sigma_g^2 + \frac{L^2}{R}\right) \quad (69)$$

$$= \mathcal{O}\left(\frac{(1-c)^2\sigma_l^2}{R^2K} + \frac{(1-c)^2}{R^3K} + \frac{\sigma_g^2}{RKL^2} + \frac{1}{R^2K}\right) \quad (70)$$

Second term,

$$8K\eta_l^2 \|\nabla F(\tilde{w}^r)\|^2 \quad (71)$$

$$8K \cdot \mathcal{O}\left(\frac{1}{RK^2L^2}\right) \cdot \mathcal{O}\left(\frac{FL}{\sqrt{RKN}} + \frac{(1-c)^2}{R}\sigma_g^2 + \frac{L^2(1-c)^2}{R^{3/2}\sqrt{KN}}\sigma_l^2 + \frac{L^2c^2\epsilon^2}{R}\right) \quad (72)$$

$$= \mathcal{O}\left(\frac{F}{LR^{3/2}K^{3/2}\sqrt{N}} + \frac{(1-c)^2\sigma_g^2}{KR^2L^2} + \frac{(1-c)^2\sigma_l^2}{R^{5/2}K^{3/2}\sqrt{N}} + \frac{c^2\epsilon^2}{KR^2}\right) \quad (73)$$

Eq. (72) is from Theorem 1.

Last term

$$16KL^2\eta_l^2c^2\rho^2\epsilon^2 \quad (74)$$

$$= \mathcal{O}\left(KL^2\frac{1}{RK^2L^2}c^2\frac{1}{R}\epsilon^2\right) \quad (75)$$

$$= \mathcal{O}\left(\frac{c^2\epsilon^2}{KR^2}\right) \quad (76)$$

If we summarize and omit higher order of R ,

$$Ke(T + 8K\eta_l \|\nabla F(\tilde{w})\|^2 + 24KL^2\eta_l c^2\epsilon^2) \quad (77)$$

$$= \mathcal{O}\left(\frac{(1-c)^2\sigma_l^2}{R^2} + \frac{(1-c)^2}{R^3} + \frac{\sigma_g^2}{RL^2} + \frac{1}{R^2} + \frac{F}{LR^{3/2}\sqrt{KN}} + \frac{(1-c)^2\sigma_g^2}{R^2L^2} + \frac{(1-c)^2\sigma_l^2}{R^{5/2}\sqrt{KN}} + \frac{c^2\epsilon^2}{R^2}\right) \quad (78)$$

where Eq. (78) is from the related term with $\sigma_l^2, \sigma_g^2, \epsilon$. Therefore, by substituting the Eq. (78) in $\Delta_{\mathcal{F}}$, we get

$$\begin{aligned} \Delta_{\mathcal{F}} &\leq \mathcal{O}\left(\frac{\sigma_g^2}{\sqrt{R}} + \frac{L}{\sqrt{R}}\left(\frac{(1-c)^2\sigma_l^2}{R^2} + \frac{(1-c)^2}{R^3} + \frac{\sigma_g^2}{RL^2} + \frac{1}{R^2}\right.\right. \\ &\quad \left.\left.+ \frac{F}{LR^{3/2}\sqrt{KN}} + \frac{(1-c)^2\sigma_g^2}{R^2L^2} + \frac{(1-c)^2\sigma_l^2}{R^{5/2}\sqrt{KN}} + \frac{c^2\epsilon^2}{R^2}\right)\right) \\ &= \mathcal{O}\left(\frac{\sigma_g^2}{\sqrt{R}} + \frac{L(1-c)^2\sigma_l^2}{R^{5/2}} + \frac{L(1-c)^2}{R^{7/2}} + \frac{\sigma_g^2}{LR^{3/2}} + \frac{L}{R^{5/2}}\right. \\ &\quad \left.+ \frac{F}{R^2\sqrt{KN}} + \frac{(1-c)^2\sigma_g^2}{LR^{5/2}} + \frac{L(1-c)^2\sigma_l^2}{R^3\sqrt{KN}} + \frac{Lc^2\epsilon^2}{R^{5/2}}\right) \\ &= \mathcal{O}\left(\frac{\sigma_g^2}{\sqrt{R}} + \frac{L(1-c)^2\sigma_l^2}{R^{5/2}} + \frac{\sigma_g^2}{LR^{3/2}} + \frac{Lc^2\epsilon^2}{R^{5/2}}\right), \end{aligned} \quad (79)$$

where Eq. (79) is from omitting higher order $\left(\frac{L(1-c)^2}{R^{7/2}}, \frac{L}{R^{5/2}}, \frac{(1-c)^2\sigma_g^2}{LR^{5/2}}, \frac{L(1-c)^2\sigma_l^2}{R^3\sqrt{KN}}\right)$ and irrelevant term $\left(\frac{F}{R^2\sqrt{KN}}\right)$. \square

B. Corrections on Convergence Analysis for Partial Client of FedSAM

FedGF shows a similar convergence behavior as that of FedSAM (Qu et al., 2022). We found that there exists an error in the proof of convergence for partial client participation of (Qu et al., 2022)

B.1. Full client participation: Lemma B.3 in FedSAM (Qu et al., 2022) and Lemma 6 in our FedGF

$$\begin{aligned}
 & \langle \nabla F(\tilde{w}^r), -\frac{1}{N} \sum_i^N \sum_k^K \eta_l \nabla F_i(\tilde{w}_{i,k}^r) + K\eta_l \nabla F(\tilde{w}^r) \rangle \\
 &= \frac{\eta_l K}{2} \|\nabla F(\tilde{w}^r)\|^2 + \underbrace{\frac{\eta_l}{2N^2 K} \left\| \sum_{i,k} [\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)] \right\|^2}_A - \frac{\eta_l}{2N^2 K} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2
 \end{aligned} \tag{80}$$

where the equality is from $\langle a, b \rangle = \frac{1}{2}(\|a\|^2 + \|b\|^2 - \|a - b\|^2)$ with $a = \sqrt{\eta_l K} \nabla F(\tilde{w}^r)$ and $b = -\frac{\sqrt{\eta_l}}{N\sqrt{K}} \sum_{i,k} (\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r))$.

- **Blue term**, i.e., $-\frac{1}{N} \sum_i^N \sum_k^K \eta_l \nabla F_i(\tilde{w}_{i,k}^r)$, is the updated gradients from the full clients, N local clients.
- **Red term**, i.e., $\nabla F(\tilde{w}^r), \nabla F_i(\tilde{w}^r)$, is from the definition of $f(w) = \sum_i^N \frac{m_i}{m} F_i(w)$.

B.2. Partial client participation: Lemma 10 in FedGF

In the partial client participation settings, i.e., $S < N$ where S is the number of participating client, Eq. (80) becomes as follows:

$$\begin{aligned}
 & \langle \nabla F(\tilde{w}^r), -\frac{1}{S} \sum_{i \in S^r} \sum_k^K \eta_l \nabla F_i(\tilde{w}_{i,k}^r) + K\eta_l \nabla F(\tilde{w}^r) \rangle \\
 &= \frac{\eta_l K}{2} \|\nabla F(\tilde{w}^r)\|^2 + \underbrace{\frac{\eta_l}{2S^2 K} \left\| \sum_{i \in S^r} \sum_k^K [\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)] \right\|^2}_A - \frac{\eta_l}{2N^2 K} \left\| \sum_{i,k} \nabla F_i(\tilde{w}_{i,k}^r) \right\|^2
 \end{aligned}$$

- **Blue term**, i.e., $-\frac{1}{N} \sum_i^N \sum_k^K \eta_l \nabla F_i(\tilde{w}_{i,k}^r)$, is the updated gradients from the partial clients, S local clients.
- **Red term**: $\nabla F(\tilde{w}^r)$ must be not represented as $\nabla F_i(w)$, because there are more terms of not included terms, i.e., $\sum_{i \notin S^r} \sum_k^K F_i(\tilde{w}^r)$. But the proof of (Qu et al., 2022) treats them to be the same. In FedSAM (Qu et al., 2022), the error is ignored in the convergence of partial client participation, in other words, they use the same result in partial client participation, which is used for full partial client participation. In FedGF, **ours**, we deal with the **red term** in the convergence of FedGF as follows:
(the detailed proof is in Lemma. 10):

$$\begin{aligned}
 & \left\| \sum_{i \in S^r} \sum_k^K [\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)] \right\|^2 \\
 & \leq SK \sum_{i \in S^r} \sum_k^K \|\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)\|^2 \\
 & = SK \sum_{i \in S^r} \sum_k^K \|\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r) + \nabla F_i(\tilde{w}^r) - \nabla F(\tilde{w}^r)\|^2 \\
 & \leq 2SK \sum_{i \in S^r} \sum_k^K [\|\nabla F_i(\tilde{w}_{i,k}^r) - \nabla F_i(\tilde{w}^r)\|^2 + \|\nabla F_i(\tilde{w}^r) - \nabla F(\tilde{w}^r)\|^2]
 \end{aligned}$$

Algorithm 1 FedGF

Input: $w^0, \Delta^r = 0$, local learning rate η_l , global learning rate η_g , and the number of epochs K
Parameter: T_D is threshold for D^r , W is length of the windowing.

```

1: for  $r = 0, \dots, R - 1$  do
2:   Sample subset  $\mathcal{S}^r \subseteq [N]$  of clients.
3:    $\tilde{w}^r = w^r + \rho \frac{\Delta^r}{\|\Delta^r\|}$ 
4:   for each client  $i \in \mathcal{S}^r$  in parallel (Local training) do
5:     send the  $w^r, \tilde{w}^r$  to the active clients
6:      $w_{i,0}^r = w^r$ 
7:     for  $k = 0, \dots, K - 1$  do
8:        $g_{i,k}^r = \nabla F_i(w_{i,k}^r, \zeta_{i,k})$ 
9:        $\tilde{w}_{i,k}^r = w_{i,k}^r + \rho \frac{g_{i,k}^r}{\|g_{i,k}^r\|}$ 
10:       $\tilde{w}_{i,k,c}^r = \tilde{w}^r c + \tilde{w}_{i,k}^r (1 - c)$ 
11:       $w_{i,k+1}^r = w_{i,k}^r - \eta \nabla F_i(\tilde{w}_{i,k,c}^r, \zeta_{i,k})$ 
12:    end for
13:     $\Delta_i^r = w^r - w_{i,K}^r$ 
14:    send the  $\Delta_i^r$  to the global server
15:  end for
16:   $\Delta^{r+1} = \frac{1}{|\mathcal{S}^r|} \sum_{i \in \mathcal{S}^r} \Delta_i^r$ 
17:   $D^{r+1} = \frac{1}{|\mathcal{S}^r|} \sum_{i \in \mathcal{S}^r} \|\Delta_i^r\|$ 
18:   $w^{r+1} = w^r - \eta_g \Delta^{r+1}$  (Aggregation)
19:   $I^{r+1} = \mathbf{I}[D^{r+1} > T_D]$ 
20:   $c = \text{Avg}[I^{\max(r-W+1,0)}, \dots, I^{r+1}]$ 
21: end for
    
```

C. Details of Experimental Setup

We implemented the models based on the PyTorch framework (Paszke et al., 2019) and ran the experiments with NVIDIA A5000 and A6000 processors. Specifically, our implementation of the federated learning scenario is based on the public implementation by (Zeng et al., 2023). Here, we provide the detailed settings of the experimental results in the main paper, including datasets, the model architecture, pre-processing, and hyperparameters. Also, we provide the pseudocode of FedGF in Algorithm 1.

C.1. Datasets and Model Architectures

Table 6: Datasets and Clients

Dataset	Task	Total clients	Total samples	Training samples	Test samples
CIFAR-10	Image classification	100	60,000	50,000	10,000
CIFAR-100	Image classification	100	60,000	50,000	10,000

Dataset: We consider two datasets that are widely used in federated learning: CIFAR-10 and CIFAR-100 (Krizhevsky et al.). To construct the benchmarks for the federated learning setting, we borrowed the federated learning benchmarks based on CIFAR-10/100 which are proposed by (Hsu et al., 2020). Specifically, each dataset is distributed across 100 clients, where each one receives 500 images according to the prior distribution of labels sampled from the Dirichlet distribution. When focusing on the Dirichlet distribution-based method to distribute the data samples, each client’s prior distribution is selected following a multinomial distribution drawn from a symmetric Dirichlet distribution with parameter α . If $\alpha \rightarrow \infty$, the distribution of data across clients becomes an IID distribution. Otherwise, it becomes a non-IID distribution when α decreases, i.e., $\alpha \rightarrow 0$. We test the various cases with $\alpha \in \{0, 0.005, 10\}$ on CIFAR-10 and CIFAR-100. We emphasize that $\alpha = 10$ is sufficiently large to represent the IID case. Also, $\alpha = 0.005$ looks too small, but it is appropriate to provide a moderate scenario in between the IID and non-IID cases.

Model architecture: By following the backbone architecture of the most-related work, i.e., FedSAM by (Caldarola et al., 2022), we utilize a Convolutional Neural Network (CNN) that is similar to LeNet5 (LeCun et al., 1998). Specifically, the backbone model has two 64-channel convolutional layers with a kernel size of 5×5 , each followed by a 2×2 max-pooling layer, ended by two fully connected layers with 384 and 192 channels and a linear classifier. For a fair comparison, we use the same backbone architecture for all different types of methods in Table 1 at the main paper. Also, the same architecture is identically used for the two CIFAR-10/100 benchmarks. Noteworthy, we added ResNet backbone results in Appendix D.6.

Data pre-processing: Every 32×32 input image in the datasets is pre-processed through the following standard pipeline: each training image is randomly cropped, applying padding 4 to show the final size of 32×32 , randomly horizontally flipped with probability 0.5, and finally, the pixel values are normalized with the dataset’s mean and standard deviation; normalization is applied to test images as well.

C.2. Hyperparameters

For each case of algorithm and its evaluation on the benchmarks, we tuned the hyperparameters: μ for FedProx is tuned among two choices $\{0.1, 1\}$, the global learning rate η for SCAFFOLD is selected among $\{0.1, 0.01\}$, the momentum ratio β of FedAvgM is chosen among $\{0.1, 0.9\}$. For the parameters ρ, η of FedSAM and FedASAM, we borrow the same values that were used in the original paper by (Caldarola et al., 2022). For MoFedSAM, we tuned the hyperparameter $\beta \in \{0.1, 0.9\}$, and we also use the same ρ with FedSAM. For FedGAMMA and FedSMOO, we also use the same ρ with FedSAM. For our FedGF, we tuned the hyperparameter T_D , which is the threshold for c value, among the choices: $\{0.1, 0.2, 0.3\}$. For each dataset, we have chosen $T_D = 0.2$ for CIFAR-10 and $T_D = 0.3$ for CIFAR-100, and used the same values for $\rho \in \{0.02, 0.05, 0.1\}$, as suggested in (Caldarola et al., 2022). Following the benchmark suggested in (Caldarola et al., 2022), we set the local learning rate as $\eta_l = 0.01$, global learning rate as $\eta_g = 1$, weight decay as 0.0004, and batch size as 64. Finally, the window size W is set among the choices: $\{10, 30, 50\}$.

C.3. Pseudocode of FedGF

Algorithm 1 presents the algorithmic pseudocode of the training procedures of FedGF.

D. Additional Experimental Results

Here, we provide the additional experimental results that are dropped due to the limited space of the main paper. It includes the standard deviations of the main test accuracies, the plots of test accuracies for the IID cases, the loss plots along perturbations for the CIFAR-10 dataset, the flatness metrics (LPF, eigenvalue of Hessian, Flatness Discrepancy), and the analysis of c for the CIFAR-10 experiment. Moreover, we evaluate the related algorithms in the CIFAR-100 dataset using a larger backbone, i.e., the ResNet-18 architecture, to confirm the validity of our algorithm in the larger architecture. We provide all related results, including test accuracy curves along the rounds, the flatness metrics, and the loss plots along perturbations. To smooth the convergence behavior, we average the 100 test accuracies of its vicinity every 800 rounds for CIFAR-10 and every 1000 rounds for CIFAR-100 datasets. We apply the moving average to smooth the behavior of c .

D.1. Standard Deviation of Test Accuracy for Table 1 of the Main Paper

We provide the standard deviation of test accuracies of Table 1 in the main paper. We measure the deviations by using the last 100 values of test accuracies as shown in Table 7. Most of the algorithms show an increasing trend of deviations as the number of the participating clients decreases. It confirms that the learning becomes unreliable when the number of active clients drops. In contrast, as shown in Table 7, our FedGF is shown to be robust to the number of participating clients. Also, as the non-IIDness becomes severe, i.e., $\alpha \rightarrow 0$, the deviation value has tendency to increase for most of the prior methods. However, FedGF shows consistently low deviations even with strong non-IIDness, i.e., $\pm 1.35\%$ and $\pm 0.83\%$ for CIFAR-10 and CIFAR-100 cases.

D.2. Convergence Behavior Curves of the IID cases for CIFAR-10 and CIFAR-100

Figure 8 shows the test accuracies curves for the IID cases of CIFAR-10/100 with 5 participating clients per round. When dividing the methods into three groups: i) FedAvgM and FedDyn, ii) FedGF, FedSAM, FedASAM, MoFedSAM, and FedSMOO, and the remaining group iii) FedAvg and FedProx, the first group shows the fastest convergence than other algorithms. However, the second group with flat minima searching shows a better final performance. The third group

Table 7: Standard Deviation of Test accuracy

Task	Algorithms	Dirichlet distribution parameter α								
		$Dir.(\alpha = 0)$			$Dir.(\alpha = 0.005)$			$Dir.(\alpha = 10)$		
		The number of participating clients								
		5	10	20	5	10	20	5	10	20
CIFAR-10	FedAvg	4.17	2.79	1.52	3.37	2.02	1.34	0.25	0.19	0.13
	FedAvgM	3.64	2.98	1.6	3.68	2.14	1.02	0.23	0.15	0.12
	FedProx	5.5	2.41	1.84	3.19	2	1.42	0.23	0.16	0.12
	SCAFFOLD	(X)	(X)	(X)	4.79	3.64	0.56	0.07	0.03	0.03
	FedDyn	3.32	1.95	1.69	2.78	1.65	1.11	0.16	0.13	0.12
	FedSAM	3.44	2.67	1.42	2.67	1.74	0.77	0.18	0.11	0.08
	FedASAM	3.01	2.65	1.36	3.45	1.79	1.06	0.18	0.11	0.11
	MoFedSAM	3.15	1.1	0.65	1.88	0.89	0.51	1.21	1.78	1.13
	FedGAMMA	0.1	0.09	0.06	0.11	0.09	0.06	0.19	0.17	0.15
	FedSMOO	2.86	1.7	0.96	2.85	1.18	0.51	0.14	0.12	0.09
FedGF	1.35	1.55	1.16	1.50	1.29	1.11	0.17	0.12	0.10	
CIFAR-100	FedAvg	1.75	1.2	0.77	0.92	0.66	0.37	0.19	0.16	0.14
	FedAvgM	1.79	1.19	1.22	1.05	0.52	0.35	0.21	0.17	0.15
	FedProx	1.41	1.18	0.93	0.89	0.59	0.4	0.24	0.19	0.19
	SCAFFOLD	(X)	(X)	(X)	2.17	(X)	(X)	0.07	0.07	0.08
	FedDyn	(X)	(X)	(X)	(X)	(X)	(X)	0.17	0.21	0.16
	FedSAM	1.54	0.86	0.5	1.54	0.73	0.33	0.16	0.12	0.11
	FedASAM	1.66	1.04	0.58	1.05	0.49	0.36	0.20	0.18	0.13
	MoFedSAM	2.36	2.13	1.59	2	1.3	1.02	0.9	0.76	0.66
	FedGAMMA	(X)	(X)	(X)	1.32	0.77	0.44	0.21	0.17	0.16
	FedSMOO	1.13	0.7	0.39	1.01	0.46	0.36	0.13	0.12	0.08
FedGF	0.83	0.99	0.63	0.78	0.68	0.88	0.17	0.15	0.11	

In the corresponding performance table of the main paper, we have some cases that cannot converge or fail to train (denoted as (X)). For these cases, it is not meaningful to measure the standard deviation, so we collectively denote the cases as (X).

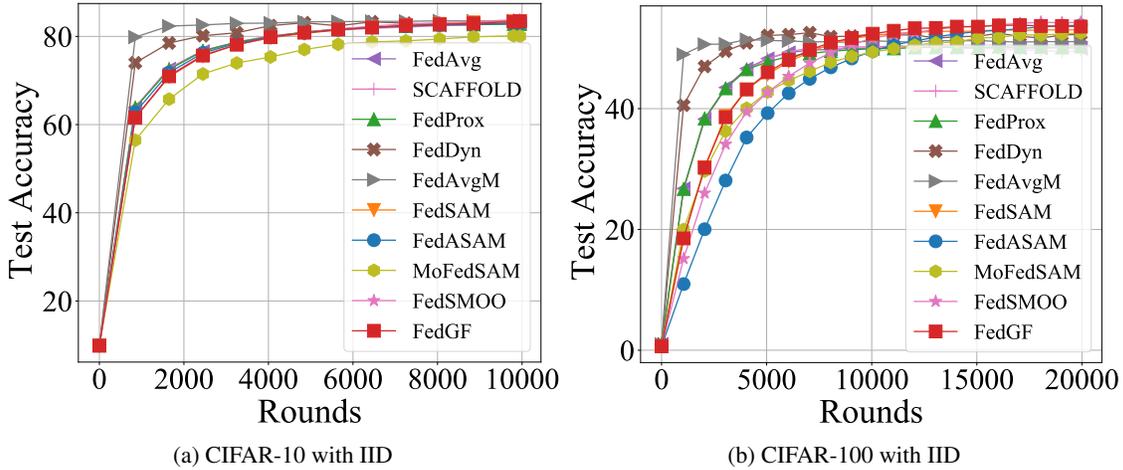


Figure 8: Convergence behaviors for the IID cases

shows moderate convergence speed but is limited to the degraded final performance. As aforementioned, IID cases do not show meaningful differences between the FL algorithms. When reminding the results for the non-IID cases, FedGF shows significantly faster convergence for the non-IID cases which are more important and crucial in the practical FL settings, where the strong heterogeneity is trivial.

D.3. Flatness Results for the IID cases for CIFAR-10

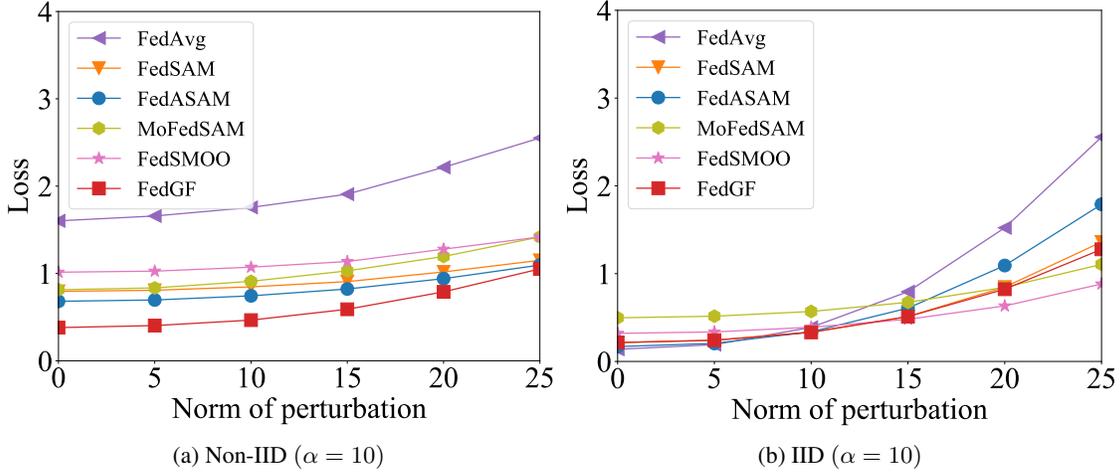


Figure 9: Loss plots along to perturbation for CIFAR-10

Loss plots along the perturbation for CIFAR-10: In Figure 9, we present the plot of loss values according to the increasing perturbation radius. It shows the family of sharpness-aware federated learning algorithms is shown to find quite flat minima in both IID and Non-IID experiments. MoFedSAM and FedSMOO show more flatter in IID case, but FedGF shows lower loss values within the perturbation range of $0 \sim 15$. It confirms the competitive performance of FedGF in the IID cases. Also, we highlight that FedGF consistently finds flatter minima even in the non-IID cases as verified in the main paper, but others suffer from the significant degradation of flatness and accuracies.

Table 8: LPF, σ_{\max} , and $\Delta_{\mathcal{F}}$ results for CIFAR-10

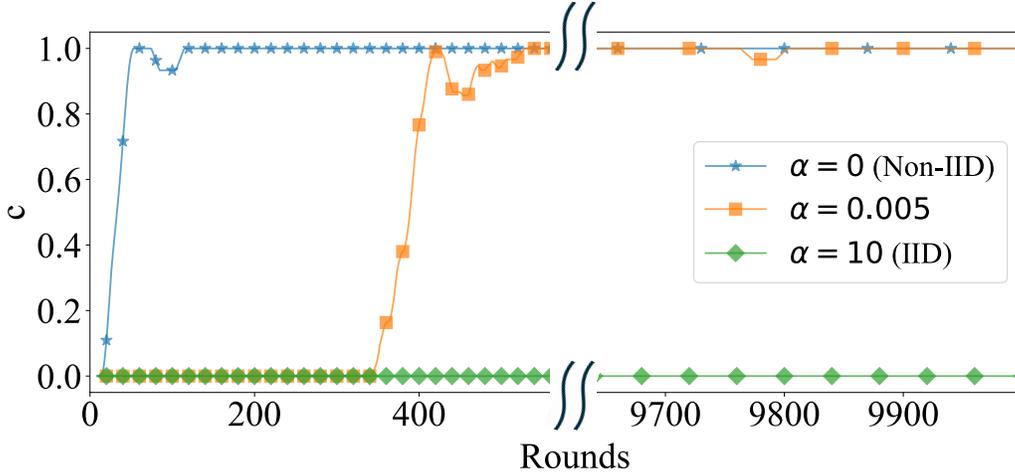
Algorithm	Non-IID ($\alpha = 0$)			IID ($\alpha = 10$)		
	LPF \downarrow	$\lambda_{\max} \downarrow$	$\Delta_{\mathcal{F}} \downarrow$	LPF \downarrow	$\lambda_{\max} \downarrow$	$\Delta_{\mathcal{F}} \downarrow$
FedAvg	1.03	81.57	0.083	0.33	103.19	0.003
FedSAM	0.79	43.32	0.023	0.34	25.30	0.002
FedASAM	0.72	25.46	0.029	0.34	14.99	0.003
MoFedSAM	0.79	15.11	0.028	0.73	27.28	0.003
FedSMOO	1.05	36.33	0.019	0.37	24.45	0.009
FedGF (ours)	0.68	14.07	0.023	0.30	21.32	0.003

\downarrow means that a lower value is preferred.

Flatness metrics: Table 8 presents the flatness metrics of the global model for the Non-IID and IID cases of the CIFAR-10 benchmark. In the IID case, there do not exist large differences between algorithms with flat minima searching. In contrast, the Non-IID results show the huge gap for LPF and λ_{\max} between FedGF and other algorithms. Although $\Delta_{\mathcal{F}}$ of FedGF shows similar value with other flat minima searching algorithms, considering LPF and loss plot in Fig. 9, FedGF finds flat minima with lower loss values. It verifies our conjecture, where the flatness of the local model, is not connected to the flatness of the global model. The behaviors are exactly coincide with the results on CIFAR-100.

D.4. Analysis of c for CIFAR-10

We also present the c values according to communication rounds for CIFAR-10 in Figure 10. It shows an identical trend to the CIFAR-100 dataset. As $\alpha \rightarrow 0$, i.e., for the non-IID case, the c values quickly get close to 1. Moreover, as the distribution gets close to IID, i.e., $\alpha = 10$, c stays near 0 during training.


 Figure 10: Behavior of c for CIFAR-10

D.5. Impact of Parameters

Fig. 11 and Fig. 12 shows the impact of parameters, i.e., the size of local perturbation ρ , the size of global perturbation ρ_g , the length of windowing W , and the threshold parameter T_D , on the CIFAR-10 and CIFAR-100 cases with $\alpha = 0$ and 5 active clients. For ρ and ρ_g , we optimize it in the range of $[0.02, 0.5]$. As shown in the following table, it has a minimal impact to test accuracy when it is smaller than 0.1. Therefore, we set the same value for ρ and ρ_g . For T_D , we searched it in the range of $[0.1, 1]$. It shows stable performance in the range of $[0.1, 0.3]$, which is not much more sensitive. For W , which is the window of computing c , we found that it has a small impact on accuracy in experiments.

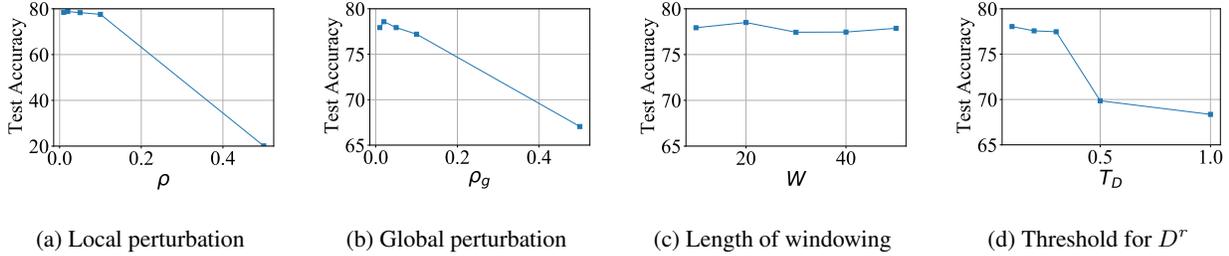


Figure 11: Impacts of parameters on CIFAR-10 dataset.

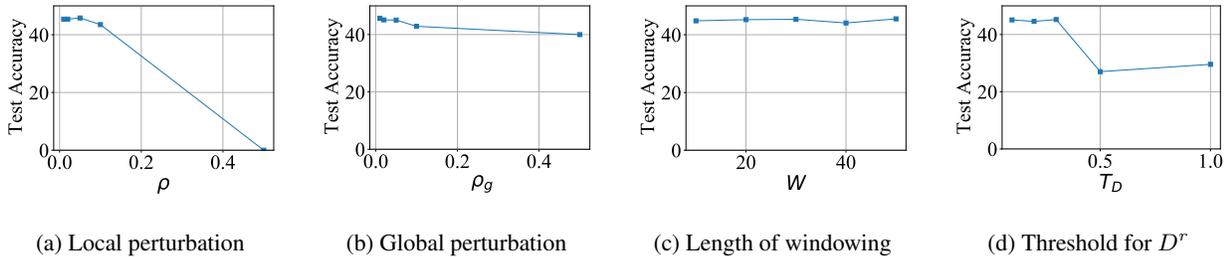


Figure 12: Impacts of parameters on CIFAR-100 dataset.

D.6. Experiments with Larger Network Architecture

We conduct training using a modified ResNet-18 architecture to verify the effectiveness of FedGF for the regime of larger models. In the experiments, we have omitted Batch Normalization (BN) layers (Ioffe & Szegedy, 2015) because the

Table 9: Evaluation on a ResNet-18 backbone

Task	Algorithms	$Dir.(\alpha = 0)$		$Dir.(\alpha = 0.005)$		$Dir.(\alpha = 10)$	
		Number of participating clients per each round					
		5	20	5	20	5	20
CIFAR-100	FedAvg	17.01 (1.62)	35.1(0.7)	34.52 (0.81)	35.24(0.47)	45.25 (0.47)	43.8 (0.2)
	FedSAM	16.04 (1.58)	26.21 (0.63)	39.97 (0.79)	42.62(0.29)	57.92 (0.15)	57.75 (0.09)
	FedASAM	17.24 (1.66)	34.87 (0.71)	39.09 (0.8)	39.82(0.3)	51.07 (0.14)	50.15 (0.12)
	MoFedSAM	27.48 (1.75)	35.42 (2.81)	44.27 (2.79)	21.96(0.97)	48.52 (1.81)	50.14 (1.11)
	FedSMOO	22.3 (1.71)	26.08 (0.52)	41.64 (0.95)	44.58(0.27)	57.29 (0.12)	56.94 (0.1)
	FedGF (ours)	42.85 (1.25)	48.2 (0.17)	48.69 (0.34)	44.29 (0.18)	57.98 (0.13)	58.16 (0.08)

Numbers shown in the (·) is the standard deviation of the test accuracies.

performance is maximized without the BN layers (we follow the wisdom in (Du et al., 2022)). As done in the main experiments of the main paper, we consider the CIFAR-100 dataset with heterogeneity $\alpha = \{0, 0.005, 10\}$, and the number of the participating clients ranges between 5 and 20.

Hyperparameters: We follow the same setting of hyperparameters as shown in Appendix C, excepting for the following setups: For SAM-based algorithms, i.e., FedSAM, FedASAM, MoFedSAM, and FedGF, we use the $\rho \in \{0.02, 0.05, 0.1\}$. For FedASAM, we set the hyperparameter $\eta = 0.2$. The local learning rate is $\eta_l = 0.01$, and the batch size is 64.

D.6.1. TEST ACCURACIES

As shown in Table 9, FedGF shows outstanding performance over the baselines, as similarly observed in the LeNet-based results. As we expect, the performance gains are remarkable for the non-IID cases. Also, we emphasize that FedGF also outperforms the baselines for the IID case with considerable gains.

D.6.2. CONVERGENCE BEHAVIOR CURVES OF THE RESNET-18 EXPERIMENTS

Fig. 13b and 13a show the convergence behavior of the FL algorithms for the ResNet-based experiments. For the IID case, FedGF shows slightly slower behavior, but it reaches the top. For the non-IID case, FedGF shows significantly faster performance and obtains the top accuracy. The behavior shown in the non-IID case is analogous to our observation for the LeNet-based experiments in the main paper.

D.6.3. FLATNESS RESULTS OF THE RESNET-18 EXPERIMENTS

Loss plots along the perturbation: Fig. 14a and 14b show the loss plots along the perturbations. For both the IID and non-IID cases, FedGF shows the flatter minima and the lowest loss values. The results coincide with the results in the main paper, verifying the effectiveness of FedGF with even larger model architectures.

Flatness metrics: We compute LPF, λ_{\max} , and $\Delta_{\mathcal{F}}$ as done in the main paper. FedGF shows the best LPF values for both IID and non-IID cases.

Analysis of c : Fig. 15 shows the behavior of c along the rounds for the ResNet-18 experiments. As expected, for the non-IID case, FedGF gradually pushes c to the larger value as the round goes on. As aforementioned, larger c is preferred when the heterogeneity gets worse. Interestingly, when the model is sufficiently trained, FedGF is shown to prefer smaller c , so that it drops to zero at the last moment of the rounds. We conjecture that the result is due to the much larger capacity of ResNet over LeNet architecture. When sufficient model capacity is given, deep training has more chances to find a much flatter region, which simultaneously covers the global and local models. For that case, sufficiently flat minima can be found even when FedGF uses a smaller c value. At the early stage of rounds, FedGF struggles to suppress the heterogeneity by using a larger c ; then it tunes to use smaller c when sufficiently high-performance models are found.

E. Formal Descriptions for Flatness Metrics

In this paper, we have measured the flatness of models by using Low-Pass Filter (LPF)-based flatness metrics suggested in (Bisla et al., 2022), as a great tool well correlates to the generalization performance of models. We borrow the formal

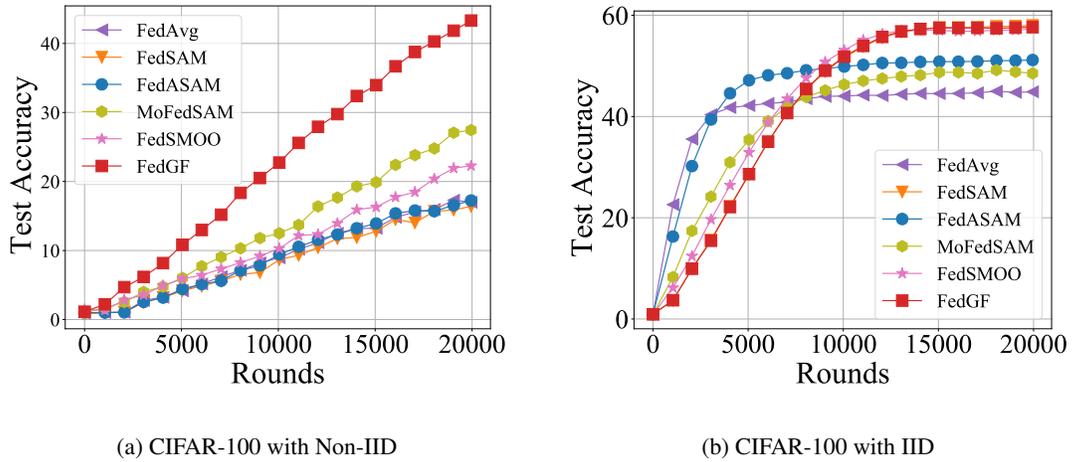


Figure 13: Test Accuracy for ResNet-18 Backbone.

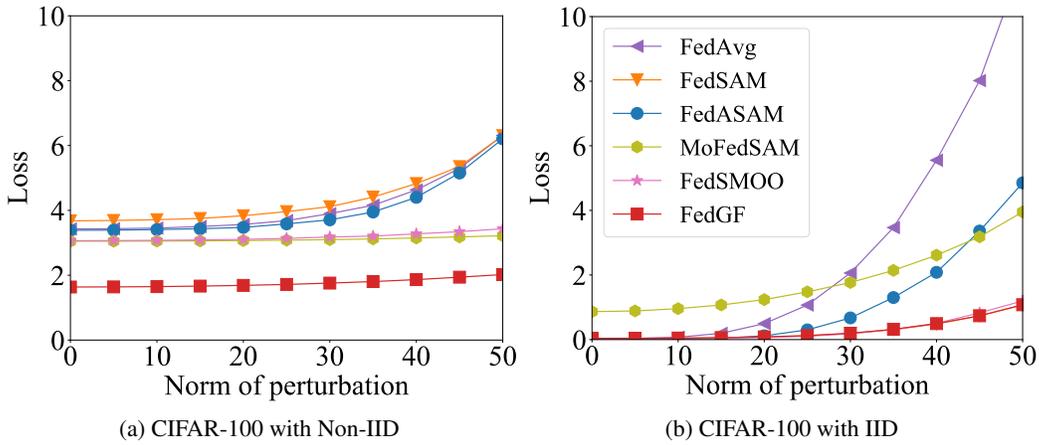


Figure 14: Loss plots along the perturbations for ResNet-18 backbone

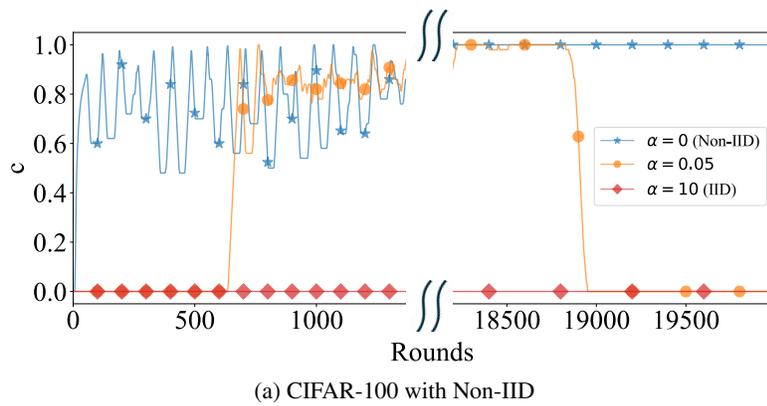


Figure 15: Behavior of c of the ResNet-18 backbone for CIFAR-100

definition of LPF metrics as follows. For an intuition, LPF becomes large when the loss varies rapidly.

Definition (from (Bisla et al., 2022)) (LPF) Let $K \sim \mathcal{N}(0, \sigma^2 I)$ be a kernel of a Gaussian filter. LPF-based sharpness

Table 10: LPF, λ_{\max} , and $\Delta_{\mathcal{F}}$ results of the ResNet-18 backbone for the CIFAR-100 experiments

Algorithm	Non-IID ($\alpha = 0$)			IID ($\alpha = 10$)		
	LPF \downarrow	λ_{\max} \downarrow	$\Delta_{\mathcal{F}}$ \downarrow	LPF \downarrow	λ_{\max} \downarrow	$\Delta_{\mathcal{F}}$ \downarrow
FedAvg	3.38	95.71	0.039	3.0	30.56	0.001
FedSAM	3.35	44.42	0.028	0.26	7.04	0.001
FedASAM	3.86	113.39	0.032	1.07	5.89	0.001
MoFedSAM	2.69	3.83	0.024	1.35	25.13	0.001
FedSMOO	3.20	29.62	0.026	0.26	5.45	0.004
FedGF (ours)	1.79	18.91	0.027	0.27	7.82	0.001

\downarrow means that a lower value is preferred.

measure at w^* is defined as the convolution of loss function F with the Gaussian filter:

$$LPF(w^*) := (F \circledast K) = \int F(w^* - \tau)K(\tau)d\tau.$$