

METHODS WITH LOCAL STEPS AND RANDOM RESHUFFLING FOR GENERALLY SMOOTH NON-CONVEX FEDERATED OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Non-convex Machine Learning problems typically do not adhere to the standard smoothness assumption. Based on empirical findings, Zhang et al. (2020b) proposed a more realistic generalized (L_0, L_1) -smoothness assumption, though it remains largely unexplored. Many existing algorithms designed for standard smooth problems need to be revised. However, in the context of Federated Learning, only a few works address this problem but rely on additional limiting assumptions. In this paper, we address this gap in the literature: we propose and analyze new methods with local steps, partial participation of clients, and Random Reshuffling without extra restrictive assumptions beyond generalized smoothness. The proposed methods are based on the proper interplay between clients' and server's stepsizes and gradient clipping. Furthermore, we perform the first analysis of these methods under the Polyak-Łojasiewicz condition. Our theory is consistent with the known results for standard smooth problems, and our experimental results support the theoretical insights.

1 INTRODUCTION

Distributed optimization problems and distributed algorithms have gained a lot of attention in recent years in the Machine Learning (ML) community. In particular, modern problems often lead to the training of deep neural networks with billions of parameters on large datasets (Brown et al., 2020; Kolesnikov et al., 2019). To make the training time feasible (Li, 2020), it is natural to parallelize computations (e.g., stochastic gradients computations), i.e., apply *distributed training* algorithms (Goyal et al., 2017; You et al., 2019; Le Scao et al., 2023). Another motivation for the usage of distributed methods is dictated by the fact that data can be naturally distributed across multiple devices/clients and be private, which is a typical scenario in *Federated Learning* (FL) (Konečný et al., 2016; McMahan et al., 2016; Kairouz et al., 2019).

Typically, such problems are not L -smooth as indicated by Defazio & Bottou (2019) that motivated the optimization researchers to study so-called *generalized smoothness assumptions*. In particular, Zhang et al. (2020b) propose (L_0, L_1) -smoothness assumption, which allows the norm of the Hessian to grow linearly with the norm of the gradient, and empirically validate it for several problems involving the training of neural networks. In addition, Ahn et al. (2023); Crawshaw et al. (2024); Wang et al. (2024) demonstrate that linear transformers with few layers satisfy this assumption, highlighting the practical importance of (L_0, L_1) -smoothness. Moreover, the theoretical convergence of different methods is studied under (L_0, L_1) -smoothness in the literature (Zhang et al., 2020b;a; Koloskova et al., 2023a; Chen et al., 2023; Li et al., 2024a;b; Crawshaw et al., 2024). Noticeably, most of these methods utilize *gradient clipping* (Pascanu et al., 2013).

However, in the context of Distributed/Federated Learning, the theoretical convergence of methods is weakly explored under (L_0, L_1) -smoothness. In particular, only a couple of papers analyze methods with *local steps* and *Random Reshuffling* – two highly important techniques in FL – under (L_0, L_1) -smoothness but only with additional restrictive assumptions such as data homogeneity (Liu et al., 2022), bounded variance (Wang et al., 2024) or cosine relatedness (Qian et al., 2021). Moreover, to the best of our knowledge, there are no results for the methods with *partial participation* of clients under (L_0, L_1) -smoothness. Such a noticeable gap in the literature leads us to

the question: *is it possible to design methods with local steps, Random Reshuffling, and partial participation of clients with provable convergence guarantees under (L_0, L_1) -smoothness without additional restrictive assumptions?* In this paper, we give a positive answer to this question.

1.1 OUR CONTRIBUTIONS

- **New method with local steps.** We propose a new method with local steps called Clip-LocalGDJ (Algorithm 1). This method can be seen as a version of LocalGD (Mangasarian, 1995; McMahan et al., 2016) with different clients and server stepsizes and (smoothed) gradient clipping (Pascanu et al., 2013) on a server side. We also prove the convergence of Clip-LocalGDJ for distributed non-convex (L_0, L_1) -smooth problems without additional assumptions such as data homogeneity used in the previous works (Liu et al., 2022).
- **New method with local steps and Random Reshuffling.** The second method we propose – CLERR (Algorithm 2) – utilizes local steps and Random Reshuffling and clipping once-in-a-epoch. For the new method, we derive rigorous convergence bounds for distributed non-convex (L_0, L_1) -smooth problems without additional assumptions such as bounded variance (Wang et al., 2024) or cosine relatedness (Qian et al., 2021).
- **New method with local steps, Random Reshuffling, and partial participation.** We extend RR-CLI (Malinovsky et al., 2023a), utilizing Random Reshuffling of clients (as an alternative to clients’ sampling) and clients’ data at each meta-epoch, and adjust it to the case of (L_0, L_1) -smooth objectives through the usage of (smoothed) gradient clipping at the end of each meta-epoch. For the resulting method called Clipped RR-CLI (Algorithm 3), we derive a convergence rate for distributed non-convex (L_0, L_1) -smooth problems without additional restrictive assumptions. To the best of our knowledge, this is the first result for an FL method with partial participation of clients under (L_0, L_1) -smoothness assumption.
- **Results for the PL-functions.** For all three new methods, we derive new results under Polyak-Łojasiewicz condition (Polyak, 1963; Łojasiewicz, 1963) that, to the best of our knowledge, are the first results for FL methods under (L_0, L_1) -smoothness and Polyak-Łojasiewicz condition. The analysis is based on the careful consideration of two possible cases (the gradient is either “small” or “big”) and induction proof of the boundedness of certain metrics.
- **Tightness of the results.** The derived results are tight: in the special case of L -smooth functions, our results recover the known ones for the non-clipped version of the algorithms.
- **Numerical experiments.** Our numerical experiments illustrate the superiority of the proposed methods over the existing baselines.

1.2 PRELIMINARIES

In this paper, we consider a standard distributed optimization problem

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M f_m(x) \right\}, \quad (1)$$

where $[M] \stackrel{\text{def}}{=} \{1, 2, \dots, M\}$ represents the set of all workers participating in the training, and each $f_m : \mathbb{R}^d \rightarrow \mathbb{R}$ is a non-convex function corresponding to the loss computed on the data available on client m for the current model parameterized by $x \in \mathbb{R}^d$. Throughout the paper, we consider two setups: either workers can compute the full gradient $\nabla f_m(x)$ of their loss functions or they can compute only a stochastic gradient at each step. In the latter case, we will assume that functions $\{f_m\}_{m=1}^M$ have the finite-sum form

$$f_m(x) = \frac{1}{N} \sum_{j=1}^N f_{mj}(x), \quad \forall m \in [M],$$

where $f_{mj}(x)$ corresponds to the local loss of the current model parameterized by $x \in \mathbb{R}^d$, evaluated for the j -th data point on the dataset belonging to the m -th client.

1.3 RELATED WORK

Local training. Local Training (LT), where clients perform multiple optimization steps on their local data before engaging in the resource-intensive process of parameter synchronization, stands out as one of the most effective and practical techniques for training FL models. LT was proposed by Mangasarian (1995); Povey et al. (2014); Moritz et al. (2015) and later promoted by McMahan et al. (2016). Early theoretical analyses of LT methods relied on restrictive data homogeneity assumptions, which are often unrealistic in real-world federated learning (FL) settings (Stich, 2018; Li et al., 2019; Haddadpour & Mahdavi, 2019). Later, Khaled et al. (2019a;b) removed limiting data homogeneity assumptions for LocalGD (Gradient Descent (GD) with LT). Then, Woodworth et al. (2020); Glasgow et al. (2022) derived lower bounds for GD with LT and data sampling, showing that its communication complexity is no better than minibatch Stochastic Gradient Descent (SGD) in settings with heterogeneous data. Another line of works focused on the mitigating so-called client drift phenomenon, which naturally occurs in LocalGD applied to distributed problems with heterogeneous local functions (Karimireddy et al., 2020; Tran-Dinh et al., 2021; Gorbunov et al., 2021; Thapa et al., 2022; Mishchenko et al., 2022; Malinovsky et al., 2023b).

Random reshuffling. Although standard Stochastic Gradient Descent (SGD) (Robbins & Monro, 1951) is well-understood from a theoretical perspective (Rakhlin et al., 2012; Bottou et al., 2018; Nguyen et al., 2018; Gower et al., 2019; Drori & Shamir, 2020; Khaled & Richtárik, 2020; Demidovich et al., 2024), most widely-used ML frameworks rely on *sampling without replacement*, as it works better in the training neural networks (Bottou, 2009; Recht & Ré, 2013; Bengio, 2012; Sun, 2020). It leverages the finite-sum structure by ensuring each function is used once per epoch. However, this introduces bias: individual steps may not reflect full gradient descent steps on average. Thus, proving convergence requires more advanced techniques. Three popular variants of sampling without replacement are commonly used. *Random Reshuffling (RR)*, where the training data is randomly reshuffled before the start of every epoch, is an extremely popular and well-studied approach. The aim of RR is to disrupt any potentially untoward default data sequencing that could hinder training efficiency. RR works very well in practice. *Shuffle Once (SO)* is analogous to RR, however, the training data is permuted randomly only once prior to the training process. The empirical performance is similar to RR. *Incremental Gradient (IG)* is identical to SO with the difference that the initial permutation is deterministic. This approach is the simplest, however, ineffective. IG has been extensively studied over a long period (Luo, 1991; Grippo, 1994; Li et al., 2022; Ying et al., 2019; Gürbüzbalaban et al., 2019; Nguyen et al., 2021). A major challenge with IG lies in selecting a particular permutation for cycling through the iterations, a task that Nedic & Bertsekas (2001) highlight as being quite difficult. (Bertsekas, 2015) provides an example that underscores the vulnerability of IG to poor orderings, especially when contrasted with RR. Meaningful theoretical analyses of the SO method have only emerged recently (Safran & Shamir, 2020; Rajput et al., 2020). RR has been shown to outperform both SGD and IG for objectives that are twice-smooth (Gürbüzbalaban et al., 2015; Haochen & Sra, 2019). Jain et al. (2019) examine the convergence of RR for smooth objectives. Safran & Shamir (2020); Rajput et al. (2020) provide lower bounds for RR. Mishchenko et al. (2020) recently conducted a thorough analysis of IG, SO and RR using innovative and simplified proof techniques, resulting in better convergence rates. Recent advances on RR can be found in (Cha et al., 2023; Cai et al., 2023; Koloskova et al., 2023b).

Generalized smoothness. Let us remind that the function f is said to be L -smooth if there exist $L \geq 0$ such that $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$ for all $x, y \in \mathbb{R}^d$. For twice-differentiable functions, it is equivalent to $\|\nabla^2 f(x)\| \leq L$, for all $x \in \mathbb{R}^d$. This assumption is very standard in the optimization field. Recently, based on extensive experiments, Zhang et al. (2020b) introduced a generalization of this condition called (L_0, L_1) -smoothness. Namely, twice-differentiable function f is said to be (L_0, L_1) -smooth if $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$, for all $x \in \mathbb{R}^d$. Compared to the standard smoothness, this condition is its strict relaxation, and it is applied to a broader range of functions. Zhang et al. (2020b) demonstrated empirically that generalized smoothness provides a more accurate representation of real-world task objectives, especially in the context of training deep neural networks. During LSTM training, it was noted that the local Lipschitz constant L_0 near the stationary point is thousands of times smaller than the global Lipschitz constant L . Under this condition, Zhang et al. (2020b) provided a theoretical justification for the gradient clipping technique (Pascanu et al., 2013), which is considered effective in mitigating the issue of exploding gradients. Their results were improved by (Zhang et al., 2020a; Koloskova et al., 2023a). (Chen

Algorithm 1 Clip-LocalGDJ: Clipped Local Gradient Descent with Jumping

```

1: Input: Synchronization/communication times  $0 = t_0 < t_1 < t_2 < \dots < t_{P-1}$ , initial vector
    $x_0 \in \mathbb{R}^d$ , number of epochs  $P \geq 1$ , constants  $c_0, c_1 > 0$ .
2: Initialize  $x_0^m = \hat{x}_0 = x_0$  for all  $m \in [M] \stackrel{\text{def}}{=} \{1, 2, \dots, M\}$ .
3: for  $p = 0, 1, \dots, P - 1$  do
4:   Choose the server stepsize  $\gamma_p = \frac{1}{c_0 + c_1 \|\nabla f(\hat{x}_{t_p})\|}$ .
5:   Choose small inner stepsize  $\alpha_p > 0$ .
6:   for  $m = 1, \dots, M$  do
7:      $x_{t_p}^m = \hat{x}_{t_p}$ 
8:     for  $t \in \{t_p, \dots, t_{p+1} - 2\}$  do
9:        $x_{t+1}^m = x_t^m - \alpha_p \nabla f_m(x_t^m)$ 
10:    end for
11:   end for
12:    $g_p = \frac{1}{\alpha_p(t_{p+1}-1-t_p)} \left( \hat{x}_{t_p} - \frac{1}{M} \sum_{m=1}^M x_{t_{p+1}-1}^m \right)$ 
13:    $\hat{x}_{t_{p+1}} = \hat{x}_{t_p} - \gamma_p g_p$ 
14: end for

```

et al., 2023) establish various useful properties of generalized-smooth functions, propose generalizations of (L_0, L_1) -smoothness and optimal first-order algorithms for solving generalized-smooth non-convex problems. Li et al. (2024a;b) extend the (L_0, L_1) -smoothness condition, introduce a novel analysis technique that bounds gradients along the trajectory, analyze GD, SGD, Nesterov’s accelerated gradient method and Adam. (Crawshaw et al., 2024) consider a coordinate-wise version of generalized smoothness. (Ahn et al., 2023; Crawshaw et al., 2024; Wang et al., 2024) demonstrate that linear transformers with few layers satisfy generalized smoothness empirically. There are few papers on distributed algorithms that combine local steps or reshuffling with generalized smoothness. Qian et al. (2021) examined clipped IG; Wang et al. (2024) investigated Adam with RR; (Liu et al., 2022) studied LocalGD, however, all of the papers contain additional restrictive assumptions. This is a significant gap in the literature and we close it in our paper.

2 NEW METHODS

In this section, we introduce the new methods – Clip-LocalGDJ (Algorithm 1), CLERR (Algorithm 2), and Clipped RR-CLI (Algorithm 3).

Clip-LocalGDJ. As standard LocalGD, the first method (Clip-LocalGDJ, Algorithm 1) alternates between local GD steps on each worker and synchronization/averaging steps. However, there are two noticeable differences between Clip-LocalGDJ and LocalGD. The first one is the usage of different clients’ and server’s stepsizes. In our method, clients’ stepsizes are typically smaller than the server’s ones, which allows us to handle the client drift. Then, on the server, the pseudogradient g_p is computed, and the server performs a Clip-GD-type step, which is a second important difference compared to LocalGD. Since the server’s stepsize is typically larger than the clients’ stepsizes, the local steps can be seen as steps determining the update direction, and the server step can be seen as a larger “jump” in the averaged update direction.

CLERR. In CLERR (Algorithm 2), each client does a full epoch of RR before between synchronization steps (similarly to (Malinovsky et al., 2023b)), and similarly to Clip-LocalGDJ, (smoothed) clipping is applied only to the averaged pseudogradient g_t once in an epoch. In contrast, a naïve combination of clipping with RR uses clipping at each step, which can amplify the bias of RR and lead to poor performance (as we illustrate in our experiments).

Clipped RR-CLI. Clipped RR-CLI (Algorithm 3) is the first FL algorithm that combines clipping, local steps, local dataset reshuffling, server and client step sizes and regularized client partial participation (sampling of clients without replacement). It is based on RR-CLI proposed by Malinovsky et al. (2023a) and leverages the core techniques proposed in FedAvg (McMahan et al., 2016). The key idea is similar to CLERR, but in addition to the reshuffling of clients’ data, Clipped RR-CLI performs a reshuffling of the groups of clients as an alternative to the standard i.i.d. sam-

Algorithm 2 CLERR: Clipped once in an Epoch Random Reshuffling

```

216 1: Input: Starting point  $x_0 \in \mathbb{R}^d$ , number of epochs  $T$ , constants  $c_0, c_1 > 0$ .
217 2: for  $t = 0, \dots, T - 1$  do
218 3:   Choose global stepsize  $\gamma_t = \frac{1}{c_0 + c_1 \|\nabla f(x_t)\|}$ .
219 4:   Choose small inner stepsize  $\alpha_t > 0$ .
220 5:   Sample a permutation  $\pi_t = \{\pi_t(1), \dots, \pi_t(N)\}$ .
221 6:   for  $m = 1, \dots, M$  do
222 7:      $x_{t,0}^m = x_t$ 
223 8:     for  $j = 0, \dots, N - 1$  do
224 9:        $x_{t,j+1}^m = x_{t,j}^m - \alpha_t \nabla f_{m,\pi_t(j)}(x_{t,j}^m)$ .
225 10:    end for
226 11:     $g_t^m = \frac{1}{\alpha_t N} (x_t - x_{t,N}^m)$ 
227 12:  end for
228 13:   $g_t = \frac{1}{M} \sum_{m=1}^M g_t^m$ .
229 14:   $x_{t+1} = x_t - \gamma_t g_t$ .
230 15: end for

```

Algorithm 3 Clipped RR-CLI: Federated optimization with server and global steps, clipping, random shuffling and partial participation with shuffling

```

234 1: Input: cohort size  $C \in \{1, 2, \dots, M\}$ ; number of rounds  $R = M/C$ ; initial iterate/model
235  $x_0 \in \mathbb{R}^d$ ; number of meta-epochs  $T \geq 1$ , constants  $c_0, c_1 > 0$ .
236 2: for meta-epoch  $t = 0, 1, \dots, T - 1$  do
237 3:   Choose global stepsize  $\theta_t = \frac{1}{c_0 + c_1 \|\nabla f(x_t)\|}$ .
238 4:   Choose small server stepsize  $\eta_t > 0$ .
239 5:   Choose small client stepsize  $\gamma_t > 0$ .
240 6:    $x_t^0 = x_t$ 
241 7:   Client-Reshuffling: sample a permutation  $\lambda = (\lambda_0, \lambda_1, \dots, \lambda_{R-1})$  of  $[R]$ 
242 8:   for communication rounds  $r = 0, \dots, R - 1$  do
243 9:     Send model  $x_t^r$  to participating clients  $m \in S_t^{\lambda_r}$  (server broadcasts  $x_t^r$  to clients  $m \in S_t^{\lambda_r}$ )
244 10:    for all clients  $m \in S_t^{\lambda_r}$ , locally in parallel do
245 11:       $x_{m,t}^{r,0} = x_t^r$  (client  $m$  initializes local training using the latest global model  $x_t^r$ )
246 12:      Data-Random-Reshuffling: sample a permutation  $\pi_m = (\pi_m^0, \pi_m^1, \dots, \pi_m^{N-1})$  of  $[N]$ 
247 13:      for all local training data points  $j = 0, 1, \dots, N - 1$  do
248 14:         $x_{m,t}^{r,j+1} = x_{m,t}^{r,j} - \gamma_t \nabla f_{m,\pi_m^j}^{r,j}(x_{m,t}^{r,j})$  (client  $m$  passes once its local data in  $\pi_m$  order)
249 15:      end for
250 16:       $g_{m,t}^r = \frac{1}{\gamma_t N} (x_t^r - x_{m,t}^{r,N})$  (client  $m$  computes local update direction  $g_{m,t}^r$ )
251 17:    end for
252 18:     $g_t^r = \frac{1}{C} \sum_{m \in S_t^{\lambda_r}} g_{m,t}^r$  (server aggregates local directions  $g_{m,t}^r$  of the clients cohort  $S_t$ )
253 19:     $x_t^{r+1} = x_t^r - \eta_t g_t^r$  (server updates the model in aggregated direction  $g_t^r$  with server stepsize  $\eta_t$ )
254 20:  end for
255 21:   $g_t = \frac{1}{R} \sum_{i=0}^{R-1} g_t^i$ 
256 22:   $x_{t+1} = x_t - \theta_t g_t$  (global step after all communication rounds during meta-epoch)
257 23: end for

```

pling of clients at each communication round. At the end of each meta-epoch, the server performs a smoothed Clip-GD-type step similar to the one used in CLERR, which allows the method to make a larger step with an accumulated pseudogradient.

When the number of workers is large, partial participation is preferable. In this case, Clipped RR-CLI (Algorithm 3) is the best option as it utilizes partial participation. Otherwise, if we have access to full gradients on the workers, then Clip-LocalGDJ (Algorithm 1) is preferable. In case when the workers can compute only a stochastic gradient, then CLERR (Algorithm 2) is recommended.

3 ASSUMPTIONS

In this section, we list assumptions adopted in the paper.

Assumption 1. *There exists $f^*, f_m^*, f_{mj}^* \in \mathbb{R}$ such that $f(x) \geq f^*$, $f_m(x) \geq f_m^*$, $f_{mj}(x) \geq f_{mj}^*$, $m \in [M]$, $j \in [N]$, for all $x \in \mathbb{R}^d$.*

The next assumption is a strict relaxation of the standard smoothness.

Assumption 2 (Asymmetric (L_0, L_1) -smoothness (Zhang et al., 2020b; Chen et al., 2023)). *The functions $f(x)$, $\{f_m(x)\}_{m=1}^M$ and $\{f_{mj}(x)\}_{m=1, j=1}^{M, N}$ are asymmetrically (L_0, L_1) -smooth:*

$$\|\nabla f(x) - \nabla f(y)\| \leq (L_0 + L_1 \|\nabla f(x)\|) \|x - y\|, \quad \forall x, y \in \mathbb{R}^d,$$

$$\|\nabla f_m(x) - \nabla f_m(y)\| \leq (L_0 + L_1 \|\nabla f_m(x)\|) \|x - y\|, \quad \forall m \in [M], x, y \in \mathbb{R}^d,$$

$$\|\nabla f_{mj}(x) - \nabla f_{mj}(y)\| \leq (L_0 + L_1 \|\nabla f_{mj}(x)\|) \|x - y\|, \quad \forall m \in [M], j \in [N], x, y \in \mathbb{R}^d.$$

Empirical findings of Zhang et al. (2020b) revealed that generalized smoothness characterizes real-world task objectives in a more precise way, particularly when applied to the training of DNNs. Moreover, the above assumption is satisfied in Distributionally Robust Optimization for some problems (Jin et al., 2021).

The assumption below generalizes the smoothness condition even further.

Assumption 3 (Symmetric (L_0, L_1) -smoothness (Chen et al., 2023)). *The functions $f(x)$, $\{f_m(x)\}_{m=1}^M$ and $\{f_{mj}(x)\}_{m=1, j=1}^{M, N}$ are symmetrically (L_0, L_1) -smooth:*

$$\|\nabla f(x) - \nabla f(y)\| \leq (L_0 + L_1 \sup_{u \in [x, y]} \|\nabla f(u)\|) \|x - y\|, \quad \forall x, y \in \mathbb{R}^d,$$

$$\|\nabla f_m(x) - \nabla f_m(y)\| \leq (L_0 + L_1 \sup_{u \in [x, y]} \|\nabla f_m(u)\|) \|x - y\|, \quad \forall m \in [M], x, y \in \mathbb{R}^d,$$

$$\|\nabla f_{mj}(x) - \nabla f_{mj}(y)\| \leq (L_0 + L_1 \sup_{u \in [x, y]} \|\nabla f_{mj}(u)\|) \|x - y\|, \quad \forall m \in [M], j \in [N], x, y \in \mathbb{R}^d.$$

A common generalization of strong convexity in the literature is the Polyak–Łojasiewicz condition.

Assumption 4 (Polyak–Łojasiewicz condition (Polyak, 1963; Łojasiewicz, 1963)). *Suppose Assumption 1 holds for the function f . There exists $\mu > 0$, such that $\|\nabla f(x)\|^2 \geq 2\mu(f(x) - f^*)$.*

4 THEORETICAL CONVERGENCE RATES

In this section, we describe our convergence results. Let us first introduce the notation. Put $\Delta^* \stackrel{\text{def}}{=} f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$, $\bar{\Delta}^* \stackrel{\text{def}}{=} f^* - \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} f_{mj}^*$. Define $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Let ζ be a constant such that $0 < \zeta \leq \frac{1}{4}$. Fix accuracy $\varepsilon > 0$. Let $P \geq 1$ be the number of epochs. For all $0 \leq p \leq P - 1$, denote

$$\hat{a}_p = L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|, \quad a_p = L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|, \quad 1 \leq t_{p+1} - t_p \leq H.$$

We start by formulating the convergence result for Clip-LocalGDJ (Algorithm 1) in non-convex asymmetric generalized-smooth case. More details can be found in Appendix B.1.

Theorem 1. *Let Assumptions 1 and 2 hold. Choose any $P \geq 1$. Choose small local stepsizes α_p , server stepsizes γ_p so that $\frac{\zeta}{\hat{a}_p} \leq \gamma_p \leq \frac{1}{4\hat{a}_p}$. Then, the iterates $\{\hat{x}_{t_p}\}_{p=0}^{P-1}$ of Algorithm 1 satisfy*

$$\min_{0 \leq p \leq P-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\} \right\} \leq \frac{\left(1 + \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p}\right)^P}{P} \delta_0 + \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p} \Delta^*.$$

Corollary 1. If $P \geq \frac{32\delta_0}{\zeta\varepsilon}$ and α_p is small enough, then $\min_{0 \leq p \leq P-1} \left\{ \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\} \right\} \leq \varepsilon$.

The rates we obtain in Corollary 1 are consistent with the previously established rates of LocalGD and GD in the standard smooth case, i.e., when $L_1 = 0$. Indeed, we recover the rate $\mathcal{O}\left(\frac{L_0\delta_0}{\varepsilon}\right)$ for LocalGD (Koloskova et al., 2020). Notice, that if $H = 1$, the Algorithm 1 reduces to vanilla GD, and we recover its rate $\mathcal{O}\left(\frac{L_0\delta_0}{\varepsilon}\right)$ (Khaled & Richtárik, 2020). In the (L_0, L_1) -smooth case, setting $H = 1$, we recover the rate $\mathcal{O}\left(\frac{L_0\delta_0}{\varepsilon}\right)$ of clipped GD from (Zhang et al., 2020b).

Below we state the convergence result for Clip-LocalGDJ (Algorithm 1) in non-convex asymmetric generalized-smooth case under the PL-condition. For more details, see Appendix B.2.

Theorem 2. Let Assumptions 1 and 2 hold. Let Assumption 4 hold. Choose any integer $P > \frac{64\delta_0 L_1^2}{\mu\zeta}$. Choose small local stepsizes α_p , server stepsizes γ_p so that $\frac{\zeta}{\hat{a}_p} \leq \gamma_p \leq \frac{1}{4\hat{a}_p}$. Let \tilde{P} be an integer such that $0 \leq \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$ be a constant, $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$. Put $\delta_P \stackrel{\text{def}}{=} f(\hat{x}_{t_P}) - f^*$. Then, the iterates $\{\hat{x}_{t_p}\}_{p=0}^P$ of Algorithm 1 satisfy

$$\delta_P \leq \left(1 - \frac{\mu\zeta}{4L_0}\right)^{P-\tilde{P}} \delta_0 + \frac{4L_0 A \alpha^2}{\mu\zeta}.$$

Corollary 2. Choose $\alpha \leq \min \left\{ \sqrt{\frac{\delta_0}{AP}}, L_1 \sqrt{\frac{8\delta_0 \varepsilon}{L_0 AP}} \right\}$. If $P \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, then $\delta_P \leq \varepsilon$.

In the standard smooth case, when $L_1 = 0$, we guarantee the iteration complexity $\mathcal{O}\left(\frac{L_0}{\mu} \ln \frac{2\delta_0}{\varepsilon}\right)$, which matches the LocalGD (Koloskova et al., 2020) and GD (Khaled & Richtárik, 2020) rates.

The above results can be generalized to the symmetric (L_0, L_1) -smooth case, see Theorem 5 in Appendix B.3 for details.

Let $T \geq 1$ be the number of epochs. For all $0 \leq t \leq T-1$, denote

$$\hat{a}_t = L_0 + L_1 \|\nabla f(x_t)\|, \quad a_t = L_0 + L_1 \max_m \|\nabla f_m(x_t)\|, \quad \tilde{a}_t = L_0 + L_1 \max_{m,j} \|\nabla f_{mj}(x_t)\|.$$

Further, we outline the convergence result for CLERR (Algorithm 2) in non-convex asymmetric generalized-smooth case. For more details, see Appendix C.1.

Theorem 3. Let Assumptions 1 and 2 hold. Choose any $T \geq 1$. Choose small client stepsizes α_t , global stepsizes γ_t so that $\frac{\zeta}{\hat{a}_t} \leq \gamma_t \leq \frac{1}{4\hat{a}_t}$. Then, the iterates $\{x_t\}_{t=0}^{T-1}$ of Algorithm 2 satisfy

$$\begin{aligned} & \mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \\ & \leq \frac{8 \left(1 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((N-1)(2N-1) + 2(N+1))\right)^T}{T} \delta_0 + \frac{6\alpha_t^2 \tilde{a}_t^3}{\hat{a}_t} (N+1) \Delta^*. \quad (2) \end{aligned}$$

Corollary 3. If $T \geq \frac{256\delta_0}{\zeta\varepsilon}$ and α_t is small enough, we have $\mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \leq \varepsilon$.

In the standard smooth case, we recover the rate $\mathcal{O}\left(\frac{L_0\delta_0}{\varepsilon}\right)$ of RR (Mishchenko et al., 2020).

We relegate the convergence result for CLERR (Algorithm 2) in non-convex asymmetric generalized-smooth case under the PL-condition to Appendix C.2. In the standard smooth case we recover the rate $\mathcal{O}\left(\frac{L_0}{\mu} \ln \frac{2\delta_0}{\varepsilon}\right)$ of RR (Mishchenko et al., 2020).

Further, we formulate the convergence result for Clipped RR-CLI (Algorithm 3) in non-convex asymmetric generalized-smooth case. For more details, see Appendix D.1.

Theorem 4. *Let Assumptions 1 and 2 hold for functions f , $\{f_m\}_{m=1}^M$ and $\{f_{mj}\}_{m=1, j=1}^{M, N}$. Choose any $T \geq 1$. Choose small local stepsizes γ_t , small server stepsizes η_t , global stepsizes θ_t so that $\frac{\zeta}{\hat{a}_t} \leq \theta_t \leq \frac{1}{4\hat{a}_t}$. Then, the iterates $\{x_t\}_{t=0}^{T-1}$ of Algorithm 3 satisfy*

$$\begin{aligned} \mathbb{E} \left[\min_{0 \leq t \leq T-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \\ \leq \frac{\left(1 + \frac{2\hat{a}_t\hat{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \hat{a}_t + \eta_t^2 R a_t)\right)^T}{T} \delta_0 \\ + \frac{2\hat{a}_t\hat{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \left(\eta_t^2 a_t \Delta^* + \gamma_t^2 N \hat{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^* \right). \end{aligned}$$

Corollary 4. *If $T \geq \frac{72\delta_0}{\zeta\varepsilon}$ and γ_t, η_t are small enough, then $\mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \leq \varepsilon$.*

Finally, we provide the convergence result for Clipped RR-CLI (Algorithm 3) in non-convex asymmetric generalized-smooth case under the PŁ-condition in Appendix D.2.

5 EXPERIMENTS

We split our experimental results into 3 parts. In Section 5.1, we provide results for the Algorithm 2 with random reshuffling and jumping in the end of each epoch. In Section 5.2, we consider Algorithm 1 with local steps and jumping in the end of every communication round. In Section 5.3 we consider Algorithm 3, which has local steps, uses random reshuffling of clients and client data and performs jumping in the end of every epoch. Moreover, in Section E we provide additional technical details on the experiments. Finally, in Section F we provide additional experiments, that did not fit in the main text. In Section F.1 we investigate the influence of inner step size on the convergence of Algorithm 2, and in Section F.2 we provide additional logistic regression experiments.

All the mentioned methods have a parameterized stepsize $\gamma_t = \frac{1}{c_0 + c_1 \|g_t\|}$. If we denote

$$\beta = \frac{1}{2c_0}, \quad \lambda = \frac{c_0}{c_1}, \quad (3)$$

we can estimate γ_t as stepsize multiplied by clipping coefficient: $\frac{\beta}{2} \min \left\{ 1, \frac{\lambda}{\|g_t\|} \right\} \leq \gamma_t \leq \beta \min \left\{ 1, \frac{\lambda}{\|g_t\|} \right\}$. We use this connection in the process of tuning constants c_0 and c_1 .

In our experiments, we consider the synthetic problem, a sum of shifted fourth-order functions:

$$f(x) = \frac{1}{N} \sum_{i=1}^N \|x - x_i\|^4, \quad x_i \in [-10, 10]^d. \quad (4)$$

The main reason to consider this problem is that it is (L_0, L_1) -smooth, but not L -smooth Zhang et al. (2020b). Additionally, in Section 5.1.1 we consider the problem of image classification of ResNet-20 He et al. (2016) on CIFAR-10 dataset Krizhevsky et al. (2009). All the methods and baselines were tuned with grid-search over logarithmic grid.

5.1 METHODS WITH RANDOM RESHUFFLING

We conduct this experiment on problem (4), where $d = 1$, $N = 1000$. We consider the Shuffle Once methods, which shuffle data once at the beginning of training. As baselines, we consider the following methods: regular SO method, which is just SGD with shuffling at the start of training, Nastya from Malinovsky et al. (2022) with one worker, Clipped SO (CSO), which clips stochastic gradients at every step of the method. The results are presented in Figure 1. As one can see from Figure 1, methods with clipping significantly outperform the rest. This empirical result justifies the theoretical fact of the importance of clipping for optimization of (L_0, L_1) -smooth objectives. Additionally, we see that among methods with clipping, CLERR shows better results than CSO. From this, we can conclude that clipping the final (pseudo)gradient approximation at the end of an epoch gives better results than clipping on every step.

432
433
434
435
436
437
438
439
440
441
442

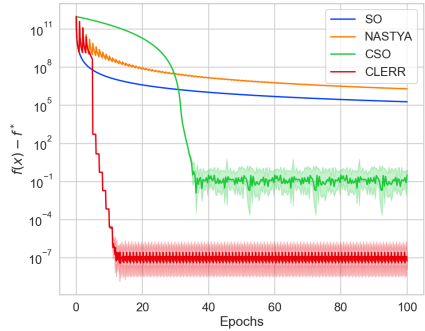


Figure 1: Function residual for (4), $\alpha_t = 10^{-7}$.

444
445
446
447
448
449
450
451
452
453

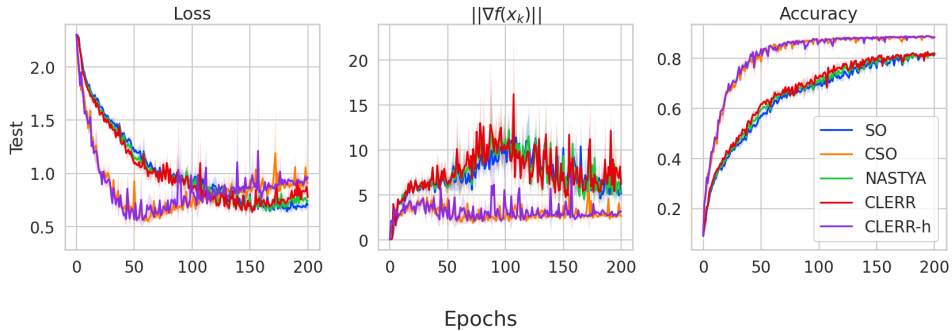


Figure 2: Loss, gradient norm and accuracy on test dataset for ResNet-18 on Cifar-10, $\alpha_t = 0.01$

454
455

5.1.1 RESNET-18 ON CIFAR-10

456
457
458
459
460
461
462
463
464
465

In Zhang et al. (2020b) the authors obtained results on a positive correlation between gradient norm and local smoothness for the problem of training neural networks in language modeling and image classification tasks. To check, whether our findings in synthetic experiments also take place for neural networks, we decided to test Algorithm 2 in the same image classification problem: train ResNet-18 He et al. (2016) on the CIFAR-10 dataset Krizhevsky et al. (2009). Additionally, we consider heuristical modification of Algorithm 2, which we call CLERR-h. The details of it we provide next. The overall results of the experiment on test data are shown in Figure 2. Additionally, we provide results on train data along with technical details in Appendix ??.

466
467
468
469
470
471
472
473

From this Figure we can see, that both jumping (Nastya and CLERR) and clipping on outer step (CLERR) does not have any impact on this problem. On the other hand, CSO shows the best results. Since in this problem regular clipping already works very well, we decided to heuristically modify our Algorithm 2: take the best clipping level and inner stepsize of CSO and use it on inner iterations, and tune c_0 with c_1 for outer stepsize. We call this method CLERR-h and also provide its results in Figure 2. CLERR-h chooses a rather big outer stepsize, while the outer clipping level is very tiny. For big clipping levels method diverges. These results show that jumping does not give performance gains when the method clips on every inner step.

5.2 METHODS WITH LOCAL STEPS

474
475
476
477
478
479
480
481
482
483
484
485

In this experiment, we aim to show the effect of the jumping technique on federated learning methods. We consider problem (4) with $d = 100, N = 1000$. To make the distributions of data on each client more distinct between each other, we sort the whole dataset at the beginning of the experiment by $\|x_i\|$. Here we consider a high-dimensional setup so that the starting point has less impact on the algorithm performance. Indeed, in one-dimensional case, if we started from $x_0 \notin [-10, 10]$, the anti-gradient of every $f_i(x) = (x - x_i)^4$ would point towards minimum. Therefore, we could find such stepsize, that method converges in one iteration. On the other hand, if we consider a high-dimensional setup, then regardless of the starting point, the gradient of each $f_i(x)$ has a different direction. In this experiment we compare Algorithm 1 (C-LGDJ) with Communication Efficient Local Gradient Clipping (CELGC) (Liu et al., 2022) and Clipping-Enabled-FedAvg (CE-FedAvg) Zhang et al. (2022). The results are shown in Figure 3.

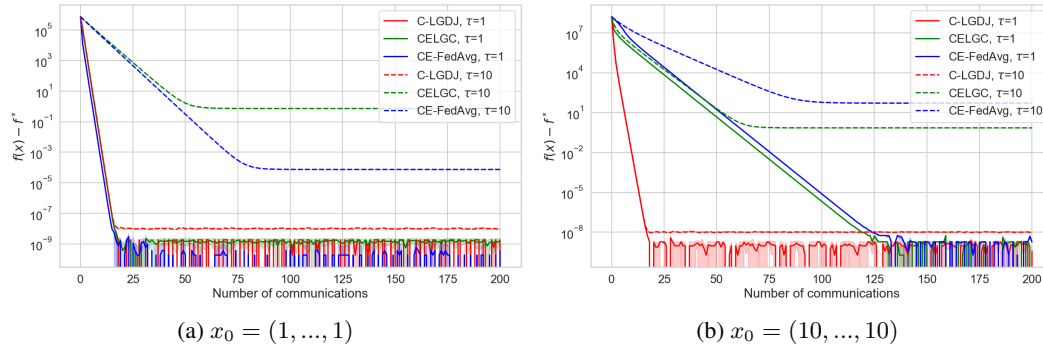


Figure 3: Function residual for (4), starting from different x_0 for different number of local steps on the client device τ .

Overall, we arrive at two conclusions. Firstly, local steps do not have any positive effect on this problem. The plots with the increased number of client steps τ only strengthen this point. Secondly, since local steps are pointless, the method works better if the server gets a better gradient approximation, which is true if the method clips gradients on the server, not on the client. This is exactly the reason why C-LGDJ has better performance in Figure 3b.

5.3 METHODS WITH LOCAL STEPS, RANDOM RESHUFFLING AND PARTIAL PARTICIPATION

In the final experiment, we consider methods with partial participation. The goal of this experiment is to investigate how clipping, local steps, partial participation and random reshuffling of both clients and client data works together. We compare Algorithm 3 with CE-FedAvg Zhang et al. (2022) with partial participation (CE-FedAvg-PP) on problem (4) with $d = 100, N = 1000$. Again, to make the distributions of data on each client more distinct between each other, we sort the whole dataset at the beginning of the experiment by $\|x_i\|$. The results are presented in Figure 4.

Since CRR-CLI uses random reshuffling of the data instead of sampling with replacement, and clips only in the end of meta-epoch, it has better gradient approximation on the global step, which results in better performance, than CE-FedAvg-PP.

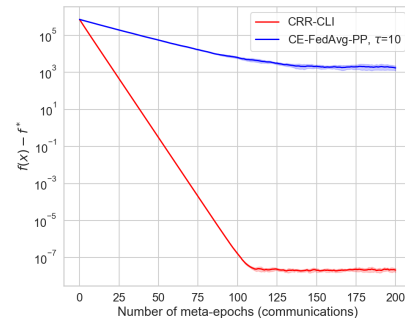


Figure 4: Function residual for (4), starting from $x_0 = (1, \dots, 1)$ with batch size 16.

6 DISCUSSION

In this paper, we consider a more general smoothness assumption and propose three new distributed methods for Federated Learning with local steps under this setting. Specifically, we analyze local gradient descent (GD) steps, local steps with Random Reshuffling, and a method that combines local steps with Random Reshuffling and Partial Participation. We provide a tight analysis for general non-convex and Polyak-Łojasiewicz settings, recovering previous results as special cases. Furthermore, we present numerical results to support our theoretical findings.

For future work, it would be valuable to explore local methods with communication compression under the generalized smoothness assumption, as well as methods incorporating incomplete local epochs. Additionally, investigating local methods with client drift reduction mechanisms to address the effects of heterogeneity, along with potentially parameter-free approaches, represents a promising direction.

REFERENCES

Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Linear attention is (maybe) all you need (to understand transformer optimization). *ArXiv*, abs/2310.01082,

- 540 2023. URL <https://api.semanticscholar.org/CorpusID:263605847>.
- 541
- 542 Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures.
543 In *Neural Networks, 2012*. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:10808461)
544 10808461.
- 545
- 546 Dimitri P. Bertsekas. Incremental gradient, subgradient, and proximal methods for convex optimiza-
547 tion: A survey. *ArXiv*, abs/1507.01030, 2015.
- 548
- 549 Léon Bottou. Curiously fast convergence of some stochastic gradient descent algorithms. 2009.
550 URL <https://api.semanticscholar.org/CorpusID:16822133>.
- 551
- 552 Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine
553 learning. *SIAM Review*, 60(2):223–311, 2018. doi: 10.1137/16M1080173. URL [https://](https://doi.org/10.1137/16M1080173)
554 doi.org/10.1137/16M1080173.
- 555
- 556 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
557 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
558 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
559 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz
560 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
561 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In
562 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-
563 ral Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc.,
564 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf)
565 [file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf).
- 566
- 567 Xufeng Cai, Cheuk Yin Lin, and Jelena Diakonikolas. Empirical risk minimization with shuffled
568 SGD: A primal-dual perspective and improved bounds. *CoRR*, abs/2306.12498, 2023. doi: 10.
569 48550/ARXIV.2306.12498. URL <https://doi.org/10.48550/arXiv.2306.12498>.
- 570
- 571 Jaeyoung Cha, Jaewook Lee, and Chulhee Yun. Tighter lower bounds for shuffling sgd: random per-
572 mutations and beyond. In *Proceedings of the 40th International Conference on Machine Learning*,
573 ICML’23. JMLR.org, 2023.
- 574
- 575 Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM*
576 *Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at
577 <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 578
- 579 Ziyi Chen, Yi Zhou, Yingbin Liang, and Zhaosong Lu. Generalized-smooth nonconvex optimiza-
580 tion is as efficient as smooth nonconvex optimization. In *Proceedings of the 40th International*
581 *Conference on Machine Learning*, ICML’23. JMLR.org, 2023.
- 582
- 583 Michael Crawshaw, Mingrui Liu, Francesco Orabona, Wei Zhang, and Zhenxun Zhuang. Robust-
584 ness to unbounded smoothness of generalized signsgd. In *Proceedings of the 36th International*
585 *Conference on Neural Information Processing Systems*, NIPS ’22, Red Hook, NY, USA, 2024.
586 Curran Associates Inc. ISBN 9781713871088.
- 587
- 588 Aaron Defazio and Léon Bottou. On the ineffectiveness of variance reduced optimization for deep
589 learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- 590
- 591 Yury Demidovich, Grigory Malinovsky, Igor Sokolov, and Peter Richtárik. A guide through the
592 zoo of biased sgd. In *Proceedings of the 37th International Conference on Neural Information*
593 *Processing Systems*, NIPS ’23, Red Hook, NY, USA, 2024. Curran Associates Inc.
- 594
- 595 Yoel Drori and Ohad Shamir. The complexity of finding stationary points with stochastic gradient
596 descent. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Con-
597 ference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp.
598 2658–2667. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.press/v119/](https://proceedings.mlr.press/v119/drori20a.html)
599 [drori20a.html](https://proceedings.mlr.press/v119/drori20a.html).

- 594 Margalit R. Glasgow, Honglin Yuan, and Tengyu Ma. Sharp bounds for federated averaging (local
595 sgd) and continuous perspective. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera
596 (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*,
597 volume 151 of *Proceedings of Machine Learning Research*, pp. 9050–9090. PMLR, 28–30 Mar
598 2022. URL <https://proceedings.mlr.press/v151/glasgow22a.html>.
- 599 Eduard Gorbunov, Filip Hanzely, and Peter Richtarik. Local sgd: Unified theory and new ef-
600 ficient methods. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceedings of The 24th*
601 *International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings*
602 *of Machine Learning Research*, pp. 3556–3564. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/gorbunov21a.html>.
- 603
604 Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter
605 Richtárik. SGD: General analysis and improved rates. In Kamalika Chaudhuri and Ruslan
606 Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*,
607 volume 97 of *Proceedings of Machine Learning Research*, pp. 5200–5209. PMLR, 09–15 Jun
608 2019. URL <https://proceedings.mlr.press/v97/qian19b.html>.
- 609
610 Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola,
611 Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training ima-
612 genet in 1 hour. *ArXiv*, abs/1706.02677, 2017.
- 613 Luigi Grippo. A class of unconstrained minimization methods for neural network training. *Opti-*
614 *mization Methods & Software*, 4:135–150, 1994.
- 615
616 M. Gürbüzbalaban, A. Ozdaglar, and P. A. Parrilo. Convergence rate of incremental gradient and
617 incremental newton methods. *SIAM Journal on Optimization*, 29(4):2542–2565, 2019. doi: 10.
618 1137/17M1147846. URL <https://doi.org/10.1137/17M1147846>.
- 619
620 Mert Gürbüzbalaban, Asuman E. Ozdaglar, and Pablo A. Parrilo. Why random reshuffling beats
621 stochastic gradient descent. *Mathematical Programming*, 186:49 – 84, 2015.
- 622
623 Farzin Haddadpour and Mehrdad Mahdavi. On the convergence of local descent methods in feder-
624 ated learning. *ArXiv*, abs/1910.14425, 2019.
- 625
626 Jeff Haochen and Suvrit Sra. Random shuffling beats SGD after finite epochs. In Kama-
627 lika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Con-*
628 *ference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp.
629 2624–2633. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/haochen19a.html>.
- 630
631 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
632 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
633 770–778, 2016.
- 634
635 Prateek Jain, Dheeraj M. Nagaraj, and Praneeth Netrapalli. Sgd without replacement: Sharper rates
636 for general smooth convex functions. *ArXiv*, abs/1903.01463, 2019.
- 637
638 Jikai Jin, Bohang Zhang, Haiyang Wang, and Liwei Wang. Non-convex distributionally robust
639 optimization: Non-asymptotic analysis. *Advances in Neural Information Processing Systems*, 34:
640 2771–2782, 2021.
- 641
642 Peter Kairouz, H. B. McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
643 Bhagoji, Keith Bonawitz, Zachary B. Charles, Graham Cormode, Rachel Cummings, Rafael G. L.
644 D’Oliveira, Salim Y. El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón,
645 Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He,
646 Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Kho-
647 dak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Oluwasanmi Koyejo, Tancrede
Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, R. Pagh, Mari-
ana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Xiaodong Song, Weikang Song,
Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma,
Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances
and open problems in federated learning. *Found. Trends Mach. Learn.*, 14:1–210, 2019.

- 648 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
649 Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learn-
650 ing. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Confer-*
651 *ence on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp.
652 5132–5143. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.press/v119/
653 karimireddy20a.html](https://proceedings.mlr.press/v119/karimireddy20a.html).
- 654 Ahmed Khaled and Peter Richtárik. Better theory for sgd in the nonconvex world. *ArXiv*,
655 abs/2002.03329, 2020. URL [https://api.semanticscholar.org/CorpusID:
656 211069380](https://api.semanticscholar.org/CorpusID:211069380).
- 657 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. First analysis of local gd on hetero-
658 geneous data. *ArXiv*, abs/1909.04715, 2019a.
- 659 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on identi-
660 cal and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*,
661 2019b.
- 662 Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly,
663 and Neil Houlsby. Big transfer (bit): General visual representation learning. In *European Con-*
664 *ference on Computer Vision*, 2019.
- 665 Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian U. Stich. A unified
666 theory of decentralized sgd with changing topology and local updates. In *Proceedings of the 37th
667 International Conference on Machine Learning*, ICML’20. JMLR.org, 2020.
- 668 Anastasia Koloskova, Hadrien Hendrikx, and Sebastian U. Stich. Revisiting gradient clipping:
669 stochastic bias and tight convergence guarantees. In *Proceedings of the 40th International Con-*
670 *ference on Machine Learning*, ICML’23. JMLR.org, 2023a.
- 671 Anastasia Koloskova, Ryan McKenna, Zachary B. Charles, Keith Rush, and Brendan McMahan.
672 Convergence of gradient descent with linearly correlated noise and applications to differentially
673 private learning. *ArXiv*, abs/2302.01463, 2023b. URL [https://api.semanticscholar.
674 org/CorpusID:256598052](https://api.semanticscholar.org/CorpusID:256598052).
- 675 Jakub Konečný, H. B. McMahan, Felix X. Yu, Peter Richtárik, Ananda Theertha Suresh, and
676 Dave Bacon. Federated learning: Strategies for improving communication efficiency. *ArXiv*,
677 abs/1610.05492, 2016.
- 678 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
679 2009.
- 680 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman
681 Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-
682 parameter open-access multilingual language model. 2023.
- 683 Chuan Li. Demystifying gpt-3 language model: A technical overview, 2020. URL [https://
684 lambdalabs.com/blog/demystifying-gpt-3](https://lambdalabs.com/blog/demystifying-gpt-3).
- 685 Haochuan Li, Jian Qian, Yi Tian, Alexander Rakhlin, and Ali Jadbabaie. Convex and non-convex
686 optimization under generalized smoothness. NIPS ’23, Red Hook, NY, USA, 2024a. Curran
687 Associates Inc.
- 688 Haochuan Li, Alexander Rakhlin, and Ali Jadbabaie. Convergence of adam under relaxed assump-
689 tions. In *Proceedings of the 37th International Conference on Neural Information Processing
690 Systems*, NIPS ’23, Red Hook, NY, USA, 2024b. Curran Associates Inc.
- 691 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of
692 fedavg on non-iid data. *ArXiv*, abs/1907.02189, 2019.
- 693 Xiao Li, Zhihui Zhu, Anthony Man-Cho So, and Jason D Lee. Incremental methods for weakly
694 convex optimization, 2022. URL <https://arxiv.org/abs/1907.11687>.

- 702 Mingrui Liu, Zhenxun Zhuang, Yunwen Lei, and Chunyang Liao. A communication-efficient dis-
703 tributed gradient clipping algorithm for training deep neural networks. *Advances in Neural Infor-*
704 *mation Processing Systems*, 35:26204–26217, 2022.
- 705 Stanislaw Lojasiewicz. A topological property of real analytic subsets. *Coll. du CNRS, Les équations*
706 *aux dérivées partielles*, 117(87-89):2, 1963.
- 707
708 Zhi-Quan Tom Luo. On the convergence of the lms algorithm with adaptive learning rate for linear
709 feedforward networks. *Neural Computation*, 3:226–245, 1991.
- 710
711 Grigory Malinovsky, Konstantin Mishchenko, and Peter Richtárik. Server-side stepsizes and sam-
712 pling without replacement provably help in federated optimization. *Proceedings of the 4th Inter-*
713 *national Workshop on Distributed Machine Learning*, 2022.
- 714 Grigory Malinovsky, Samuel Horv’ath, Konstantin Burlachenko, and Peter Richt’arik. Federated
715 learning with regularized client participation. *ArXiv*, abs/2302.03662, 2023a. URL <https://api.semanticscholar.org/CorpusID:256627753>.
- 716
717 Grigory Malinovsky, Konstantin Mishchenko, and Peter Richtárik. Server-side stepsizes and sam-
718 pling without replacement provably help in federated optimization. In *Proceedings of the 4th*
719 *International Workshop on Distributed Machine Learning*, DistributedML ’23, pp. 85–104, New
720 York, NY, USA, 2023b. Association for Computing Machinery. ISBN 9798400704475. doi:
721 10.1145/3630048.3630187. URL <https://doi.org/10.1145/3630048.3630187>.
- 722
723 LO Mangasarian. Parallel gradient distribution in unconstrained optimization. *SIAM Journal on*
724 *Control and Optimization*, 1995.
- 725
726 H. B. McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *International*
727 *Conference on Artificial Intelligence and Statistics*, 2016.
- 728
729 Konstantin Mishchenko, Ahmed Khaled, and Peter Richtárik. Random reshuffling: Simple analysis
730 with vast improvements. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin
731 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17309–17320. Cur-
732 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/c8cc6e90ccbff44c9cee23611711cdc4-Paper.pdf)
733 [paper/2020/file/c8cc6e90ccbff44c9cee23611711cdc4-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/c8cc6e90ccbff44c9cee23611711cdc4-Paper.pdf).
- 734 Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtarik. ProxSkip:
735 Yes! Local gradient steps provably lead to communication acceleration! Finally! In Ka-
736 malika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato
737 (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
738 *Proceedings of Machine Learning Research*, pp. 15750–15769. PMLR, 17–23 Jul 2022. URL
739 <https://proceedings.mlr.press/v162/mishchenko22b.html>.
- 740 Philipp Moritz, Robert Nishihara, Ion Stoica, and Michael I. Jordan. Sparknet: Training deep net-
741 works in spark. *CoRR*, abs/1511.06051, 2015.
- 742
743 Angelia Nedic and Dimitri P. Bertsekas. Incremental subgradient methods for nondifferen-
744 tiable optimization. *SIAM Journal on Optimization*, 12(1):109–138, 2001. doi: 10.1137/
745 S1052623499362111. URL <https://doi.org/10.1137/S1052623499362111>.
- 746
747 Lam M. Nguyen, Quoc Tran-Dinh, Dzung T. Phan, Phuong Ha Nguyen, and Marten Van Dijk. A
748 unified convergence analysis for shuffling-type gradient methods. *J. Mach. Learn. Res.*, 22(1),
749 jan 2021. ISSN 1532-4435.
- 750
751 Phuong Ha Nguyen, Lam M. Nguyen, and Marten van Dijk. Tight dimension independent lower
752 bound on the expected convergence rate for diminishing step sizes in sgd. In *Neural Information*
753 *Processing Systems*, 2018. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:52965883)
754 52965883.
- 755 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural
756 networks. In *Proceedings of the 30th International Conference on International Conference on*
757 *Machine Learning-Volume 28*, 2013.

- 756 Boris Polyak. Gradient methods for the minimisation of functionals. *Ussr Computational Mathe-*
757 *matics and Mathematical Physics*, 3:864–878, 1963.
- 758 Daniel Povey, Xiaohui Zhang, and Sanjeev Khudanpur. Parallel training of deep neural networks
759 with natural gradient and parameter averaging. In *International Conference on Learning Repre-*
760 *sentations*, 2014.
- 761 Jiang Qian, Yuren Wu, Bojin Zhuang, Shaojun Wang, and Jing Xiao. Understanding gradient clip-
762 ping in incremental gradient methods. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceed-*
763 *ings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130
764 of *Proceedings of Machine Learning Research*, pp. 1504–1512. PMLR, 13–15 Apr 2021. URL
765 <https://proceedings.mlr.press/v130/qian21a.html>.
- 766 Shashank Rajput, Anant Gupta, and Dimitris Papailiopoulos. Closing the convergence gap of SGD
767 without replacement. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th Inter-*
768 *national Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*
769 *Research*, pp. 7964–7973. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/rajput20a.html)
770 [press/v119/rajput20a.html](https://proceedings.mlr.press/v119/rajput20a.html).
- 771 Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Making gradient descent optimal for
772 strongly convex stochastic optimization. In *Proceedings of the 29th International Conference on*
773 *International Conference on Machine Learning*, ICML’12, pp. 1571–1578, Madison, WI, USA,
774 2012. Omnipress. ISBN 9781450312851.
- 775 Benjamin Recht and Christopher Ré. Parallel stochastic gradient algorithms for large-scale matrix
776 completion. *Mathematical Programming Computation*, 5:201 – 226, 2013. URL [https://](https://api.semanticscholar.org/CorpusID:17109415)
777 api.semanticscholar.org/CorpusID:17109415.
- 778 Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Math-*
779 *ematical Statistics*, 22(3):400 – 407, 1951. doi: 10.1214/aoms/1177729586. URL [https://](https://doi.org/10.1214/aoms/1177729586)
780 doi.org/10.1214/aoms/1177729586.
- 781 Itay Safran and Ohad Shamir. How good is sgd with random shuffling? In Jacob Abernethy and
782 Shivani Agarwal (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125
783 of *Proceedings of Machine Learning Research*, pp. 3250–3284. PMLR, 09–12 Jul 2020. URL
784 <https://proceedings.mlr.press/v125/safran20a.html>.
- 785 Sebastian U Stich. Local sgd converges fast and communicates little. *arXiv preprint*
786 *arXiv:1805.09767*, 2018.
- 787 Ruoyu Sun. Optimization for deep learning: An overview. *Journal of the Operations Research*
788 *Society of China*, 8:249 – 294, 2020. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:220511793)
789 [CorpusID:220511793](https://api.semanticscholar.org/CorpusID:220511793).
- 790 Chandra Thapa, Mahawaga Arachchige Pathum Chamikara, Seyit Camtepe, and Lichao Sun.
791 Splitfed: When federated learning meets split learning. In *AAAI*, pp. 8485–8493. AAAI Press,
792 2022. ISBN 978-1-57735-876-3.
- 793 Quoc Tran-Dinh, Nhan H. Pham, D. Phan, and Lam M. Nguyen. Feddr - randomized douglas-
794 rachford splitting algorithms for nonconvex federated composite optimization. In *Neural*
795 *Information Processing Systems*, 2021. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:235376727)
796 [CorpusID:235376727](https://api.semanticscholar.org/CorpusID:235376727).
- 797 Bohan Wang, Yushun Zhang, Huishuai Zhang, Qi Meng, Ruoyu Sun, Zhi-Ming Ma, Tie-Yan
798 Liu, Zhi-Quan Luo, and Wei Chen. Provable adaptivity of adam under non-uniform smooth-
799 ness. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and*
800 *Data Mining*, KDD ’24, pp. 2960–2969, New York, NY, USA, 2024. Association for Com-
801 puting Machinery. ISBN 9798400704901. doi: 10.1145/3637528.3671718. URL [https://](https://doi.org/10.1145/3637528.3671718)
802 doi.org/10.1145/3637528.3671718.
- 803 Blake Woodworth, Kumar Kshitij Patel, and Nathan Srebro. Minibatch vs local sgd for hetero-
804 geneous distributed learning. In *Proceedings of the 34th International Conference on Neural*
805 *Information Processing Systems*, NIPS ’20, Red Hook, NY, USA, 2020. Curran Associates Inc.
806 ISBN 9781713829546.

810 Bicheng Ying, Kun Yuan, Stefan Vlaski, and Ali H. Sayed. Stochastic learning under random
811 reshuffling with constant step-sizes. *IEEE Transactions on Signal Processing*, 67(2):474–489,
812 2019. doi: 10.1109/TSP.2018.2878551.

813
814 Yang You, Jing Li, Sashank J. Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan
815 Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep
816 learning: Training bert in 76 minutes. *arXiv: Learning*, 2019.

817 Bohang Zhang, Jikai Jin, Cong Fang, and Liwei Wang. Improved analysis of clipping algorithms
818 for non-convex optimization. In *Proceedings of the 34th International Conference on Neural
819 Information Processing Systems, NIPS '20*, Red Hook, NY, USA, 2020a. Curran Associates Inc.
820 ISBN 9781713829546.

821 Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates
822 training: A theoretical justification for adaptivity. In *8th International Conference on Learning
823 Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020b.
824 URL <https://openreview.net/forum?id=BJgnXpVYwS>.

825
826 Xinwei Zhang, Xiangyi Chen, Mingyi Hong, Zhiwei Steven Wu, and Jinfeng Yi. Understanding
827 clipping for federated learning: Convergence and client-level differential privacy. In *International
828 Conference on Machine Learning, ICML 2022*, 2022.

829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864	CONTENTS	
865		
866		
867	1 Introduction	1
868	1.1 Our contributions	2
869	1.2 Preliminaries	2
870	1.3 Related Work	3
871		
872		
873	2 New methods	4
874		
875	3 Assumptions	6
876		
877		
878	4 Theoretical convergence rates	6
879		
880	5 Experiments	8
881		
882	5.1 Methods with random reshuffling	8
883	5.1.1 ResNet-18 on CIFAR-10	9
884	5.2 Methods with local steps	9
885	5.3 Methods with local steps, random reshuffling and partial participation	10
886		
887		
888	6 Discussion	10
889		
890		
891	A Implications of generalized smoothness	18
892		
893	B Local gradient descent	19
894		
895	B.1 Asymmetric generalized-smooth non-convex functions	19
896	B.2 Asymmetric generalized-smooth functions under $\mathcal{P}\mathcal{L}$ -condition	23
897	B.3 Symmetric generalized-smooth non-convex functions	25
898	B.4 Symmetric generalized-smooth functions under $\mathcal{P}\mathcal{L}$ -condition	30
899		
900		
901	C Random reshuffling	32
902		
903	C.1 Asymmetric generalized-smooth non-convex functions	33
904	C.2 Asymmetric generalized-smooth functions under $\mathcal{P}\mathcal{L}$ -condition	38
905		
906		
907	D Partial participation	40
908		
909	D.1 Asymmetric generalized-smooth non-convex functions	40
910	D.2 Asymmetric generalized-smooth functions under $\mathcal{P}\mathcal{L}$ -condition	46
911		
912	E Additional experimental details for main part	49
913		
914	E.1 Methods with random reshuffling	49
915	E.1.1 ResNet-18 on CIFAR-10	50
916	E.2 Methods with local steps	51
917	E.3 Methods with local steps, random reshuffling and partial participation	52

918	F Additional experiments	53
919	F.1 How the inner step size affects convergence of the method	53
920	F.2 Logistic regression experiments	54
921		
922		
923		
924	A IMPLICATIONS OF GENERALIZED SMOOTHNESS	
925		
926	Lemma 1. <i>Let f satisfy Assumption 2. Then, for any $x, y \in \mathbb{R}^d$ we have</i>	
927		
928	$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L_0 + L_1 \ \nabla f(x)\ }{2} \ x - y\ ^2.$	
929		
930	<i>Moreover, if $f^* := \inf_{x \in \mathbb{R}^d} f(x) > -\infty$, then, for all $x \in \mathbb{R}^d$, we obtain</i>	
931		
932	$\frac{\ \nabla f(x)\ ^2}{2(L_0 + L_1 \ \nabla f(x)\)} \leq f(x) - f^*.$	
933		
934		
935	<i>Proof of Lemma 1.</i> The first statement of the lemma is proven in (Zhang et al., 2020b, Ap-	
936	pendix A.1). The second statement follows from the first statement, if one substitutes y for	
937	$x - \frac{\ \nabla f(x)\ }{L_0 + L_1 \ \nabla f(x)\ } \nabla f(x)$ and uses the fact that $f^* \leq f(y)$. □	
938		
939		
940	Lemma 2. <i>Let f satisfy Assumption 3. Then, for any $x, y \in \mathbb{R}^d$ we have</i>	
941		
942	$f(y) - f(x) \leq \langle \nabla f(x), y - x \rangle + \frac{L_0 + L_1 \ \nabla f(x)\ }{2} \exp(L_1 \ x - y\) \ x - y\ ^2.$	
943		
944	<i>Moreover, if $f^* := \inf_{x \in \mathbb{R}^d} f(x) > -\infty$, then, for $\eta > 0$, such that $\eta \exp \eta \leq 1$, for all $x \in \mathbb{R}^d$, we</i>	
945	<i>obtain</i>	
946	$\frac{\eta \ \nabla f(x)\ ^2}{2(L_0 + L_1 \ \nabla f(x)\)} \leq f(x) - f^*.$	
947		
948		
949		
950	<i>Proof of Lemma 2.</i> The first part of this lemma is one of the results of (Chen et al., 2023, Propo-	
951	sition 3.2). To deal with the second statement, let us substitute y in the first statement with	
952	$x - \frac{\eta \ \nabla f(x)\ }{L_0 + L_1 \ \nabla f(x)\ } \nabla f(x)$:	
953		
954	$f^* \leq f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L_0 + L_1 \ \nabla f(x)\ }{2} \exp(L_1 \ x - y\) \ x - y\ ^2$	
955	$= f(x) - \frac{\eta \ \nabla f(x)\ ^2}{L_0 + L_1 \ \nabla f(x)\ }$	
956	$+ \frac{L_0 + L_1 \ \nabla f(x)\ }{2} \cdot \exp\left(\frac{L_1 \eta \ \nabla f(x)\ }{L_0 + L_1 \ \nabla f(x)\ }\right) \cdot \frac{\eta^2 \ \nabla f(x)\ ^2}{(L_0 + L_1 \ \nabla f(x)\)^2}$	
957	$\leq f(x) - \frac{\eta \ \nabla f(x)\ ^2}{L_0 + L_1 \ \nabla f(x)\ } + \frac{\eta \ \nabla f(x)\ ^2}{2(L_0 + L_1 \ \nabla f(x)\)} \cdot \eta \exp(\eta)$	
958	$\leq f(x) - \frac{\eta \ \nabla f(x)\ ^2}{2(L_0 + L_1 \ \nabla f(x)\)}.$	
959		
960		
961		
962		
963		
964		
965	Rearranging the terms, we get the second statement of the lemma. □	
966		
967	Lemma 3. <i>Assumption 3 holds for the function f if and only if, for any $x, y \in \mathbb{R}^d$,</i>	
968	$\ \nabla f(x) - \nabla f(y)\ \leq (L_0 + L_1 \ \nabla f(y)\) \exp(L_1 \ x - y\) \ x - y\ .$	
969		
970		
971	<i>Proof of Lemma 3.</i> This lemma is one of the results of (Chen et al., 2023, Proposition 3.2) □	

B LOCAL GRADIENT DESCENT

B.1 ASYMMETRIC GENERALIZED-SMOOTH NON-CONVEX FUNCTIONS

Theorem 1 (non-convex asymmetric generalized-smooth convergence analysis of Algorithm 1). *Let Assumptions 1 and 2 hold for functions f and $\{f_m\}_{m=1}^M$. Choose any $P \geq 1$. For all $0 \leq p \leq P-1$, denote*

$$\hat{a}_p = L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|, \quad a_p = L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|, \quad 1 \leq t_{p+1} - t_p \leq H.$$

Put $\Delta^* = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$. *Impose the following conditions on the local stepsizes α_p and server stepsizes γ_p :*

$$\alpha_p \leq \min \left\{ \frac{1}{2Ha_p}, \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}} \right\}, \quad \frac{\zeta}{\hat{a}_p} \leq \gamma_p \leq \frac{1}{4\hat{a}_p}, \quad 0 \leq p \leq P-1,$$

where $0 < \zeta \leq \frac{1}{4}$, $c \geq \sqrt{P}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Then, the iterates $\{\hat{x}_{t_p}\}_{p=0}^{P-1}$ of Algorithm 1 satisfy

$$\min_{0 \leq p \leq P-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\} \right\} \leq \frac{\left(1 + \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p}\right)^P}{P} \delta_0 + \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p} \Delta^*.$$

Put $v_p \stackrel{\text{def}}{=} t_{p+1} - 1$.

Lemma 4. *Assume that f and each f_m satisfy Assumptions 1 and 2. Then we have the following bound:*

$$\frac{1}{M} \sum_{m=1}^M \sum_{t=t_p+1}^{v_p} \|x_t^m - \hat{x}_{t_p}\|^2 \leq 8(v_p - t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).$$

Proof of Lemma 4. We have

$$\begin{aligned} \|x_t^m - \hat{x}_{t_p}\|^2 &= \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(x_j^m) \right\|^2 \\ &\leq 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p (\nabla f_m(x_j^m) - \nabla f_m(\hat{x}_{t_p})) \right\|^2 + 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\ &\leq 2(t - t_p) \sum_{j=t_p}^{t-1} (\alpha_p)^2 (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \|x_j^m - \hat{x}_{t_p}\|^2 \\ &\quad + 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2. \end{aligned}$$

Averaging, we get

$$\begin{aligned}
\frac{1}{M} \sum_{m=1}^M \|x_t^m - \hat{x}_{t_p}\|^2 &\leq \frac{2(t-t_p)}{M} \sum_{m=1}^M \sum_{j=t_p}^{t-1} (\alpha_p)^2 (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\quad + \frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\
&\leq \frac{2(t-t_p)}{M} (a_p)^2 \sum_{m=1}^M \sum_{j=t_p}^{t-1} (\alpha_p)^2 \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\quad + \frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2.
\end{aligned}$$

Recall that $\alpha_p \leq \frac{1}{2H(L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|)}$. Then we have

$$\begin{aligned}
\frac{1}{M} \sum_{m=1}^M \|x_t^m - \hat{x}_{t_p}\|^2 &\leq \frac{t-t_p}{2H^2 M} \sum_{m=1}^M \sum_{j=t_p}^{t-1} \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\quad + \frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2. \tag{5}
\end{aligned}$$

Let us bound the last term:

$$\begin{aligned}
\frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 &\leq \frac{2}{M} \sum_{m=1}^M \|\nabla f_m(\hat{x}_{t_p})\|^2 (t-t_p)^2 \alpha_p^2 \\
&\leq \frac{4}{M} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|) (f_m(\hat{x}_{t_p}) - f_m^*) \\
&\quad \times (t-t_p)^2 \alpha_p^2 \\
&\leq \frac{4(t-t_p)^2 a_p \alpha_p^2}{M} \sum_{m=1}^M (f_m(\hat{x}_{t_p}) - f_m^*) \\
&= 4(t-t_p)^2 a_p \alpha_p^2 \left(f(\hat{x}_{t_p}) - f^* + \left(f^* - \frac{1}{M} \sum_{m=1}^M f_m^* \right) \right) \\
&= 4(t-t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).
\end{aligned}$$

Further, summing (7) with respect to t , we obtain

$$\begin{aligned}
\frac{1}{M} \sum_{m=1}^M \sum_{t=t_p+1}^v \|x_t^m - \hat{x}_{t_p}\|^2 &\leq \frac{1}{2H^2M} \sum_{m=1}^M \sum_{t=t_p+1}^v (t-t_p) \sum_{j=t_p}^{t-1} \|x_j^m - \hat{x}_{t_p}\|^2 \\
&+ \sum_{t=t_p+1}^v 4(t-t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*) \\
&\leq \frac{v-t_p}{2H^2M} \sum_{m=1}^M \sum_{t=t_p+1}^v \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&+ 4 \sum_{t=t_p+1}^v (v-t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*) \\
&\leq \frac{(v-t_p)^2}{2H^2M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&+ 4(v-t_p)^3 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).
\end{aligned}$$

Using the fact that $v-t_p \leq H-1 < H$, we obtain that

$$\frac{1}{M} \sum_{m=1}^M \sum_{t=t_p+1}^v \|x_t^m - \hat{x}_{t_p}\|^2 \leq 8(v-t_p)^3 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).$$

□

Proof of Theorem 1. Applying Lemma 1, we obtain that

$$f(\hat{x}_{t_{p+1}}) \leq f(\hat{x}_{t_p}) - \gamma_p \langle \nabla f(\hat{x}_{t_p}), g_p \rangle + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2}.$$

Additionally, from the fact that $2\langle a, b \rangle = -\|a-b\|^2 + \|a\|^2 + \|b\|^2$

$$\begin{aligned}
f(\hat{x}_{t_{p+1}}) &\leq f(\hat{x}_{t_p}) - \gamma_p \langle \nabla f(\hat{x}_{t_p}), g_p \rangle + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2} \\
&\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} (-\|\nabla f(\hat{x}_{t_p}) - g_p\|^2 + \|\nabla f(\hat{x}_{t_p})\|^2 + \|g_p\|^2) \\
&+ (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2} \\
&\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p}) - g_p\|^2 + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2}.
\end{aligned}$$

Consider $\frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p}) - g_p\|^2$. We have

$$\begin{aligned}
\frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p}) - g_p\|^2 &= \frac{\gamma_p}{2} \left\| \frac{1}{M} \sum_{m=1}^M \left(\nabla f_m(\hat{x}_{t_p}) - \frac{1}{v-t_p} \sum_{j=t_p}^v \nabla f_m(x_j^m) \right) \right\|^2 \\
&\leq \frac{\gamma_p}{2} \frac{1}{M(v-t_p)^2} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\leq \frac{\gamma_p}{2(v-t_p)^2} (L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|)^2 \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&= \frac{\gamma_p a_p^2}{2(v-t_p)^2} \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2.
\end{aligned}$$

1134 Notice that

$$\begin{aligned}
1135 & \\
1136 & \frac{\gamma_p^2 \|g_p\|^2}{2} = \frac{\gamma_p^2}{2} \left\| \frac{1}{M(v-t_p)} \sum_{m=1}^M \sum_{j=t_p+1}^v \nabla f_m(x_j^m) \right\|^2 \\
1137 & \\
1138 & \\
1139 & \\
1140 & \leq \frac{\gamma_p^2}{(v-t_p)^2} \left\| \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p+1}^v (\nabla f_m(x_j^m) - \nabla f_m(\hat{x}_{t_p})) \right\|^2 \\
1141 & \\
1142 & + \frac{\gamma_p^2}{(v-t_p)^2} \left\| \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p+1}^v \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\
1143 & \\
1144 & \leq \frac{\gamma_p^2}{(v-t_p)^2} \left(L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\| \right)^2 \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
1145 & \\
1146 & + \gamma_p^2 \|\nabla f(\hat{x}_{t_p})\|^2 \\
1147 & \\
1148 & = \frac{\gamma_p^2 a_p^2}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 + \gamma_p^2 \|\nabla f(\hat{x}_{t_p})\|^2. \\
1149 & \\
1150 & \\
1151 & \\
1152 & \\
1153 &
\end{aligned}$$

1154 Therefore, we obtain

$$\begin{aligned}
1155 & \\
1156 & f(\hat{x}_{t_{p+1}}) \leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{\gamma_p a_p^2}{2M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 + \frac{\hat{a}_p \gamma_p^2 \|g_p\|^2}{2} \\
1157 & \\
1158 & \leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{\gamma_p a_p^2}{2M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
1159 & \\
1160 & + \frac{\hat{a}_p a_p^2 \gamma_p^2}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 + \hat{a}_p \gamma_p^2 \|\nabla f(\hat{x}_{t_p})\|^2 \\
1161 & \\
1162 & = f(\hat{x}_{t_p}) + \left(\hat{a}_p \gamma_p^2 - \frac{\gamma_p}{2} \right) \|\nabla f(\hat{x}_{t_p})\|^2 \\
1163 & \\
1164 & + \left(\hat{a}_p a_p^2 \gamma_p^2 + \frac{\gamma_p a_p^2}{2} \right) \frac{1}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2. \\
1165 & \\
1166 & \\
1167 & \\
1168 & \\
1169 &
\end{aligned}$$

1170 Recall that $\gamma_p \leq \frac{1}{4\hat{a}_p}$. Then, using Lemma 4, we have

$$\begin{aligned}
1171 & \\
1172 & f(\hat{x}_{t_{p+1}}) \leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \\
1173 & \\
1174 & + \left(\hat{a}_p a_p^2 \gamma_p^2 + \frac{\gamma_p a_p^2}{2} \right) \frac{1}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
1175 & \\
1176 & \leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \\
1177 & \\
1178 & + \left(\hat{a}_p a_p^2 \gamma_p^2 + \frac{\gamma_p a_p^2}{2} \right) 8(v-t_p) a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*) \\
1179 & \\
1180 & \leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{3(H-1) a_p^3 \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}. \\
1181 & \\
1182 & \\
1183 & \\
1184 &
\end{aligned}$$

1185 Let us rewrite the inequality in the following way:

$$1186 \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1) a_p^3 \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}. \quad (6)$$

Since $\gamma_p \geq \frac{\zeta}{\hat{a}_p}$, we get that

$$\frac{\gamma_p \|\nabla f(\hat{x}_{t_p})\|^2}{4} \geq \frac{\zeta \|\nabla f(\hat{x}_{t_p})\|^2}{4\hat{a}_p}.$$

Therefore,

$$\frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \geq \begin{cases} \frac{\zeta \|\nabla f(\hat{x}_{t_p})\|^2}{8L_0}, & \|\nabla f(\hat{x}_{t_p})\| \leq \frac{L_0}{L_1}, \\ \frac{\zeta \|\nabla f(\hat{x}_{t_p})\|^2}{8L_1}, & \|\nabla f(\hat{x}_{t_p})\| > \frac{L_0}{L_1}, \end{cases} = \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_1} \right\}.$$

Denote $\delta_p \stackrel{\text{def}}{=} f(\hat{x}_{t_p}) - f^*$. Then we have

$$\frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_1} \right\} \leq \delta_p - \delta_{p+1} + \frac{3(H-1)\alpha_p^2 a_p^3 (\delta_p + \Delta^*)}{2\hat{a}_p}.$$

Let $\alpha_p \leq \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}$, where $c \geq \sqrt{P}$. Applying the result of Mishchenko et al. (2020, Lemma 6), we appear at

$$\begin{aligned} \min_{0 \leq p \leq P-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_1} \right\} \right\} &\leq \frac{\left(1 + \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p}\right)^P}{P} \delta_0 \\ &+ \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p} \Delta^*. \end{aligned}$$

□

Corollary 1. Fix $\varepsilon > 0$. Choose $c = \sqrt{3(H-1)P}$. Let $\alpha_p \leq 2\sqrt{\frac{\hat{a}_p \delta_0}{3P(H-1)a_p^3 \Delta^*}}$. Then, if $P \geq \frac{32\delta_0}{\zeta\varepsilon}$, we have $\min_{0 \leq p \leq P-1} \left\{ \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_1} \right\} \right\} \leq \varepsilon$.

Proof of Corollary 1. Since $c = \sqrt{3(H-1)P}$ and $\alpha_p \leq \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}$, $\alpha_p \leq 2\sqrt{\frac{\hat{a}_p \delta_0}{3P(H-1)a_p^3 \Delta^*}}$, due to the choice of $P \geq \frac{32\delta_0}{\zeta\varepsilon}$, we obtain that

$$\frac{\left(1 + \frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p}\right)^P}{P} \delta_0 \leq \frac{\sqrt{\varepsilon} \delta_0}{P} \leq \frac{2\delta_0}{P} \leq \frac{\zeta\varepsilon}{16},$$

and that

$$\frac{3(H-1)\alpha_p^2 a_p^3}{2\hat{a}_p} \Delta^* \leq \frac{\zeta\varepsilon}{16}.$$

Therefore, $\min_{0 \leq p \leq P-1} \left\{ \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_1} \right\} \right\} \leq \varepsilon$. □

B.2 ASYMMETRIC GENERALIZED-SMOOTH FUNCTIONS UNDER PŁ-CONDITION

Theorem 2 (Asymmetric generalized-smooth convergence analysis of Algorithm 1 in PŁ-case). *Let Assumptions 1 and 2 hold for functions f and $\{f_m\}_{m=1}^M$. Let Assumption 4 hold. Choose $0 < \zeta \leq \frac{1}{4}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Choose any integer $P > \frac{64\delta_0 L_1^2}{\mu\zeta}$. For all $0 \leq p \leq P-1$, denote*

$$\hat{a}_p = L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|, \quad a_p = L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|, \quad 1 \leq t_{p+1} - t_p \leq H.$$

Put $\Delta^* = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$. Impose the following conditions on the local stepsizes α_p and server stepsizes γ_p :

$$\alpha_p \leq \min \left\{ \frac{1}{2Ha_p}, \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}, \sqrt{\frac{\mu\zeta\hat{a}_p}{48L_1^2(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}} \right\},$$

$$\gamma_p \leq \min \left\{ \frac{1}{2Ha_p}, \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}, \sqrt{\frac{2\delta_0\hat{a}_p}{3P(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}} \right\},$$

$$\frac{\zeta}{\hat{a}_p} \leq \gamma_p \leq \frac{1}{4\hat{a}_p}, \quad 0 \leq p \leq P-1,$$

where $c \geq \sqrt{P}$. Let \tilde{P} be an integer such that $0 \leq \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$ be a constant, $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$.

Then, the iterates $\{\hat{x}_{t_p}\}_{p=0}^P$ of Algorithm 1 satisfy

$$\delta_P \leq \left(1 - \frac{\mu\zeta}{4L_0}\right)^{P-\tilde{P}} \delta_0 + \frac{4L_0 A \alpha^2}{\mu\zeta},$$

where $\delta_P \stackrel{\text{def}}{=} f(\hat{x}_{t_P}) - f^*$.

Proof of Theorem 2. Let us follow the first steps of the proof of Theorem 1. Consider (6):

$$\frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}.$$

Since $\gamma_p \geq \frac{\zeta}{\hat{a}_p}$, and f satisfies Polyak–Łojasiewicz Assumption 4, we obtain that

$$\frac{\mu\zeta(f(\hat{x}_{t_p}) - f^*)}{2\hat{a}_p} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}.$$

1. Let \tilde{P} be the number of steps p , so that $\|\nabla f(\hat{x}_{t_p})\| \geq \frac{L_0}{L_1}$. For such p , we have $L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\| = \hat{a}_p \leq 2L_1 \|\nabla f(\hat{x}_{t_p})\|$. Therefore, we get

$$\frac{\mu\zeta(f(\hat{x}_{t_p}) - f^*)}{4L_1 \|\nabla f(\hat{x}_{t_p})\|} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}.$$

Notice that the relation $\hat{a}_p \leq 2L_1 \|\nabla f(\hat{x}_{t_p})\|$ and Lemma 1 together imply

$$\frac{\|\nabla f(\hat{x}_{t_p})\|}{4L_1} \leq \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{2\hat{a}_p} \leq f(\hat{x}_{t_p}) - f^*.$$

Hence, we have

$$\frac{\mu\zeta}{16L_1^2} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}.$$

Subtracting f^* on both sides and introducing $\delta_p \stackrel{\text{def}}{=} f(\hat{x}_{t_p}) - f^*$, we obtain

$$\delta_{p+1} \leq \delta_p - \frac{\mu\zeta}{16L_1^2} + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}.$$

As $\alpha_p \leq \sqrt{\frac{\mu\zeta\hat{a}_p}{48L_1^2(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}}$, it follows that $\frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p} \leq \frac{\mu\zeta}{32L_1^2}$.

Therefore, we get

$$\delta_{p+1} \leq \delta_p - \frac{\mu\zeta}{32L_1^2}.$$

2. Suppose now that $\|\nabla f(\hat{x}_{t_p})\| \leq \frac{L_0}{L_1}$. For such p , we have $L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\| = \hat{a}_p \leq 2L_0$. Hence,

$$\frac{\mu\zeta(f(\hat{x}_{t_p}) - f^*)}{4L_0} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}.$$

Subtracting f^* on both sides and introducing $\delta_p \stackrel{\text{def}}{=} f(\hat{x}_{t_p}) - f^*$, we obtain

$$\delta_{p+1} \leq \delta_p \rho + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{2\hat{a}_p}, \quad \text{where } \rho \stackrel{\text{def}}{=} 1 - \frac{\mu\zeta}{4L_0}.$$

Let $\alpha_p \stackrel{\text{def}}{=} \alpha \hat{\alpha}_p$ and $\hat{\alpha}_p \leq \sqrt{\frac{2A\hat{a}_p}{3(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}}$ for some constant $A > 0$. Then,

$$\delta_{p+1} \leq \rho \delta_p + A\alpha^2.$$

Unrolling the recursion, we derive

$$\begin{aligned} \delta_P &\leq \rho^{P-\tilde{P}} \delta_0 + A\alpha^2 \sum_{i=0}^{\infty} \rho^i - \frac{\mu\zeta}{32L_1^2} \sum_{i=0}^{N-1} \rho^i \\ &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho} - \frac{1-\rho^{\tilde{P}}}{1-\rho} \frac{\mu\zeta}{32L_1^2}. \end{aligned}$$

Notice that $\delta_{p+1} \leq \delta_p + A\alpha^2$, which implies

$$\delta_P \leq \delta_0 + (P - \tilde{P}) A\alpha^2 - \tilde{P} \frac{\mu\zeta}{32L_1^2}.$$

Since $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$, we conclude that

$$0 \leq \delta_P \leq 2\delta_0 - \tilde{P} \frac{\mu\zeta}{32L_1^2}, \quad \Rightarrow \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}.$$

Therefore, for $P > \frac{64\delta_0 L_1^2}{\mu\zeta}$ we can guarantee that $P - \tilde{P} > 0$ and

$$\begin{aligned} \delta_P &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho} - \tilde{P} \rho^{\tilde{P}} \frac{\mu\zeta}{32L_1^2} \\ &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho}. \end{aligned}$$

□

Corollary 2. Fix $\varepsilon > 0$. Choose $\alpha \leq \min \left\{ \sqrt{\frac{\delta_0}{AP}}, L_1 \sqrt{\frac{8\delta_0\varepsilon}{L_0AP}} \right\}$. Then, if $P \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we have $\delta_P \leq \varepsilon$.

Proof of Corollary 2. Since $0 \leq \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$, $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$, $\alpha \leq L_1 \sqrt{\frac{8\delta_0\varepsilon}{L_0AP}}$, due to the choice of $P \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we obtain that

$$\left(1 - \frac{\mu\zeta}{4L_0}\right)^{P-\tilde{P}} \delta_0 \leq e^{-\frac{\mu\zeta}{4L_0}(P-\tilde{P})} \delta_0 \leq \frac{\varepsilon}{2},$$

and that

$$\frac{4L_0 A}{\mu\zeta} \cdot \frac{\delta_0}{AP} \leq \frac{\varepsilon}{2}.$$

Therefore, $\delta_P \leq \varepsilon$. □

B.3 SYMMETRIC GENERALIZED-SMOOTH NON-CONVEX FUNCTIONS

Theorem 5. Let Assumptions 1 and 2 hold for functions f and $\{f_m\}_{m=1}^M$. Choose any $P \geq 1$. For all $0 \leq p \leq P-1$, denote

$$\hat{a}_p = L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|, \quad a_p = L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|, \quad 1 \leq t_{p+1} - t_p \leq H.$$

Put $\Delta^* = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$. Impose the following conditions on the local stepsizes α_p and server stepsizes γ_p :

$$\alpha_p \leq \min \left\{ \frac{1}{2Ha_p}, \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}, \frac{C}{L_0 + L_1 A t_p} \right\}, \quad \frac{\zeta}{\hat{a}_p} \leq \gamma_p \leq \frac{1}{4\hat{a}_p}, \quad 0 \leq p \leq P-1,$$

1350 where $0 < \zeta \leq \frac{1}{4}$, $c \geq \sqrt{P}$, $C \leq \frac{\ln 1.5}{H}$, $A_{t_p} \stackrel{\text{def}}{=} \max \left\{ \sqrt{\frac{2L_0 M(\delta_{t_p} + \Delta^*)}{\nu}}, \frac{2L_1 M(\delta_{t_p} + \Delta^*)}{\nu} \right\}$, ν such
 1351 that $\nu \exp \nu = 1$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Then, the iterates $\{\hat{x}_{t_p}\}_{p=0}^{P-1}$ of Algorithm 1 satisfy
 1352
 1353

$$1354 \min_{0 \leq p \leq P-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\} \right\} \leq \frac{\left(1 + \frac{3(H-1)\alpha_p^2 a_p^3}{\hat{a}_p}\right)^P}{P} \delta_0$$

$$1355 + \frac{3(H-1)\alpha_p^2 a_p^3}{\hat{a}_p} \Delta^*.$$

1360 Let us remind that $\hat{a}_p = L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|$, $a_p = L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|$ and $\Delta^* = f^* -$
 1361 $\frac{1}{M} \sum_{m=1}^M f_m^*$. Put $v_p \stackrel{\text{def}}{=} t_{p+1} - 1$.

1362 **Lemma 5.** Assume that f and each f_m satisfy Assumptions 1 and 3. Then we have the following
 1363 bound:
 1364

$$1365 \frac{1}{M} \sum_{m=1}^M \sum_{t=t_p+1}^v \|x_t^m - \hat{x}_{t_p}\|^2 \leq 16(v - t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).$$

1366
 1367
 1368 *Proof of Lemma 5.* We have

$$1369 \begin{aligned} 1370 \|x_t^m - \hat{x}_{t_p}\|^2 &= \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(x_j^m) \right\|^2 \\ 1371 &\leq 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p (\nabla f_m(x_j^m) - \nabla f_m(\hat{x}_{t_p})) \right\|^2 + 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\ 1372 &\leq 2(t - t_p) \sum_{j=t_p}^{t-1} (\alpha_p)^2 (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \\ 1373 &\quad \times \exp \{L_1 \|x_j^m - \hat{x}_{t_p}\|\} \|x_j^m - \hat{x}_{t_p}\|^2 + 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\ 1374 &\leq 2(t - t_p) \sum_{j=t_p}^{t-1} (\alpha_p)^2 (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \\ 1375 &\quad \times \exp \left\{ L_1 \left\| \sum_{\ell=t_p}^{j-1} \alpha_p \nabla f_m(x_\ell^m) \right\| \right\} \|x_j^m - \hat{x}_{t_p}\|^2 + 2 \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2. \end{aligned}$$

1376
 1377
 1378 Let us show that if $\alpha_p \leq \frac{1}{L_0 + L_1 \max_m \max_{t_p \leq \ell \leq t_{p+1}} \|\nabla f_m(x_\ell^m)\|}$, then $f_m(x_\ell^m) \leq f_m(x_{t_p}^m)$ for $t_p \leq$
 1379 $\ell \leq t_{p+1} - 1$. Notice that locally we perform the iterations of gradient descent. It means, that
 1380

$$1381 f_m(x_{\ell+1}^m) \leq f_m(x_\ell^m) - \alpha_p \|\nabla f_m(x_\ell^m)\|^2 + \frac{L_0 + L_1 \|\nabla f_m(x_\ell^m)\|}{2} \alpha_p^2 \|\nabla f_m(x_\ell^m)\|^2$$

$$1382 \leq f_m(x_\ell^m) - \alpha_p \|\nabla f_m(x_\ell^m)\|^2 + \frac{\alpha_p}{2} \|\nabla f_m(x_\ell^m)\|^2$$

$$1383 = f_m(x_\ell^m) - \frac{\alpha_p}{2} \|\nabla f_m(x_\ell^m)\|^2.$$

1384
 1385 Then $f_m(x_\ell^m) \leq f_m(x_{t_p}^m) = f_m(\hat{x}_{t_p})$ for $t_p \leq \ell \leq t_{p+1} - 1$ follows. Therefore, for such α_p we
 1386 have that
 1387

$$1388 f_m(x_\ell^m) - f_m^* \leq f_m(x_{t_p}^m) - f_m^* \leq \sum_{m=1}^M (f_m(x_{t_p}^m) - f_m^*) = M\delta_{t_p} + M\Delta^*.$$

From Lemma 2 we have

$$\min \left\{ \frac{\nu \|\nabla f_m(x_\ell^m)\|^2}{L_0}, \frac{\nu \|\nabla f_m(x_\ell^m)\|}{L_1} \right\} \leq \frac{\nu \|\nabla f_m(x_\ell^m)\|^2}{2(L_0 + L_1 \|\nabla f_m(x_\ell^m)\|)} \leq f_m(x_\ell^m) - f_m^* \leq M(\delta_{t_p} + \Delta^*).$$

For every $t_p \leq \ell \leq t_{p+1} - 1$, for every m , we establish

$$\|\nabla f_m(x_\ell^m)\| \leq \max \left\{ \sqrt{\frac{2L_0 M(\delta_{t_p} + \Delta^*)}{\nu}}, \frac{2L_1 M(\delta_{t_p} + \Delta^*)}{\nu} \right\} \stackrel{\text{def}}{=} A_{t_p}.$$

Let us choose $\alpha_p \leq \frac{C}{L_0 + L_1 A_{t_p}}$ for some $C \leq \frac{\ln 1.5}{H}$ and show by induction that for such local processes $\max_m \|\nabla f_m(x_\ell^m)\| \leq A_{t_p}$, for all $t_p \leq \ell \leq t_{p+1} - 1$. Indeed, for $\ell = t_p$ it holds trivially. Suppose it holds for all ℓ such that $t_p \leq \ell \leq \ell'$ for some ℓ' . Then, $f_m(x_{\ell'+1}^m) \leq f_m(x_{\ell'}^m)$ holds for any $\alpha_p \leq \frac{C}{L_0 + L_1 \|\nabla f_m(x_{\ell'}^m)\|}$, including the chosen stepsize. Hence, $f_m(x_{\ell'+1}^m) \leq f_m(x_{t_p}^m)$. Therefore, $\max_m \|\nabla f_m(x_\ell^m)\| \leq A_{t_p}$, for all $t_p \leq \ell \leq t_{p+1} - 1$. Then, $\frac{C}{L_0 + L_1 A_{t_p}} \leq \frac{1}{L_0 + L_1 \max_m \max_{t_p \leq \ell \leq t_{p+1}} \|\nabla f_m(x_\ell^m)\|}$.

It means that $\exp \left\{ L_1 \left\| \sum_{\ell=t_p}^{j-1} \alpha_p \nabla f_m(x_\ell^m) \right\| \right\} \leq e^{\ln 1.5} = 1.5$.

Averaging, we get

$$\begin{aligned} \frac{1}{M} \sum_{m=1}^M \|x_t^m - \hat{x}_{t_p}\|^2 &\leq \frac{3(t-t_p)}{M} \sum_{m=1}^M \sum_{j=t_p}^{t-1} (\alpha_p)^2 (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \|x_j^m - \hat{x}_{t_p}\|^2 \\ &\quad + \frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\ &\leq \frac{3(t-t_p)}{M} (a_p)^2 \sum_{m=1}^M \sum_{j=t_p}^{t-1} (\alpha_p)^2 \|x_j^m - \hat{x}_{t_p}\|^2 \\ &\quad + \frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2. \end{aligned}$$

Recall that $\alpha_p \leq \frac{1}{2Ha_p}$. Then we have

$$\frac{1}{M} \sum_{m=1}^M \|x_t^m - \hat{x}_{t_p}\|^2 \leq \frac{1.5(t-t_p)}{2H^2M} \sum_{m=1}^M \sum_{j=t_p}^{t-1} \|x_j^m - \hat{x}_{t_p}\|^2 \quad (7)$$

$$+ \frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2. \quad (8)$$

Let us bound the last term:

$$\begin{aligned}
\frac{2}{M} \sum_{m=1}^M \left\| \sum_{j=t_p}^{t-1} \alpha_p \nabla f_m(\hat{x}_{t_p}) \right\|^2 &\leq \frac{2}{M} \sum_{m=1}^M \|\nabla f_m(\hat{x}_{t_p})\|^2 (t-t_p)^2 \alpha_p^2 \\
&\leq \frac{8}{M} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|) (f_m(\hat{x}_{t_p}) - f_m^*) \\
&\quad \times (t-t_p)^2 \alpha_p^2 \\
&\leq \frac{8(t-t_p)^2 a_p \alpha_p^2}{M} \sum_{m=1}^M (f_m(\hat{x}_{t_p}) - f_m^*) \\
&= 8(t-t_p)^2 a_p \alpha_p^2 \left(f(\hat{x}_{t_p}) - f^* + \left(f^* - \frac{1}{M} \sum_{m=1}^M f_m^* \right) \right) \\
&= 8(t-t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).
\end{aligned}$$

Further, summing (7) with respect to t , we obtain

$$\begin{aligned}
\frac{1}{M} \sum_{m=1}^M \sum_{t=t_p+1}^v \|x_t^m - \hat{x}_{t_p}\|^2 &\leq \frac{1.5}{2H^2 M} \sum_{m=1}^M \sum_{t=t_p+1}^v (t-t_p) \sum_{j=t_p}^{t-1} \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\quad + \sum_{t=t_p+1}^v 8(t-t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*) \\
&\leq \frac{1.5(v-t_p)}{2H^2 M} \sum_{m=1}^M \sum_{t=t_p+1}^v \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\quad + 8 \sum_{t=t_p+1}^v (v-t_p)^2 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*) \\
&\leq \frac{1.5(v-t_p)^2}{2H^2 M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\quad + 8(v-t_p)^3 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).
\end{aligned}$$

Using the fact that $v-t_p \leq H-1 < H$, we obtain that

$$\frac{1}{M} \sum_{m=1}^M \sum_{t=t_p+1}^v \|x_t^m - \hat{x}_{t_p}\|^2 \leq 32(v-t_p)^3 a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*).$$

□

Proof of Theorem 5. Applying Lemma 1, we obtain that

$$f(\hat{x}_{t_{p+1}}) \leq f(\hat{x}_{t_p}) - \gamma_p \langle \nabla f(\hat{x}_{t_p}), g_p \rangle + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2}.$$

Additionally, from the fact that $2\langle a, b \rangle = -\|a-b\|^2 + \|a\|^2 + \|b\|^2$

$$\begin{aligned}
f(\hat{x}_{t_{p+1}}) &\leq f(\hat{x}_{t_p}) - \gamma_p \langle \nabla f(\hat{x}_{t_p}), g_p \rangle + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2} \\
&\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} (-\|\nabla f(\hat{x}_{t_p}) - g_p\|^2 + \|\nabla f(\hat{x}_{t_p})\|^2 + \|g_p\|^2) \\
&\quad + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2} \\
&\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p}) - g_p\|^2 + (L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|) \frac{\gamma_p^2 \|g_p\|^2}{2}.
\end{aligned}$$

1512 Consider $\frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p}) - g_p\|^2$. We have
 1513
 1514

$$\begin{aligned}
 1515 \quad \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p}) - g_p\|^2 &= \frac{\gamma_p}{2} \left\| \frac{1}{M} \sum_{m=1}^M \left(\nabla f_m(\hat{x}_{t_p}) - \frac{1}{v-t_p} \sum_{j=t_p}^v \nabla f_m(x_j^m) \right) \right\|^2 \\
 1516 &\leq \frac{\gamma_p}{2} \frac{1}{M(v-t_p)^2} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(\hat{x}_{t_p})\|)^2 \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
 1517 &\leq \frac{\gamma_p}{2(v-t_p)^2} (L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|)^2 \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
 1518 &= \frac{\gamma_p a_p^2}{2(v-t_p)^2} \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2.
 \end{aligned}$$

1519 Notice that
 1520
 1521

$$\begin{aligned}
 1522 \quad \frac{\gamma_p^2 \|g_p\|^2}{2} &= \frac{\gamma_p^2}{2} \left\| \frac{1}{M(v-t_p)} \sum_{m=1}^M \sum_{j=t_p+1}^v \nabla f_m(x_j^m) \right\|^2 \\
 1523 &\leq \frac{\gamma_p^2}{(v-t_p)^2} \left\| \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p+1}^v (\nabla f_m(x_j^m) - \nabla f_m(\hat{x}_{t_p})) \right\|^2 \\
 1524 &\quad + \frac{\gamma_p^2}{(v-t_p)^2} \left\| \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p+1}^v \nabla f_m(\hat{x}_{t_p}) \right\|^2 \\
 1525 &\leq \frac{\gamma_p^2}{(v-t_p)^2} (L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|)^2 \frac{1}{M} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
 1526 &\quad + \gamma_p^2 \|\nabla f(\hat{x}_{t_p})\|^2 \\
 1527 &= \frac{\gamma_p^2 a_p^2}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 + \gamma_p^2 \|\nabla f(\hat{x}_{t_p})\|^2.
 \end{aligned}$$

1528 Therefore, we obtain
 1529
 1530

$$\begin{aligned}
 1531 \quad f(\hat{x}_{t_{p+1}}) &\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{\gamma_p a_p^2}{2M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 + \frac{\hat{a}_p \gamma_p^2 \|g_p\|^2}{2} \\
 1532 &\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{2} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{\gamma_p a_p^2}{2M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
 1533 &\quad + \frac{\hat{a}_p a_p^2 \gamma_p^2}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 + \hat{a}_p \gamma_p^2 \|\nabla f(\hat{x}_{t_p})\|^2 \\
 1534 &= f(\hat{x}_{t_p}) + \left(\hat{a}_p \gamma_p^2 - \frac{\gamma_p}{2} \right) \|\nabla f(\hat{x}_{t_p})\|^2 \\
 1535 &\quad + \left(\hat{a}_p a_p^2 \gamma_p^2 + \frac{\gamma_p a_p^2}{2} \right) \frac{1}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2.
 \end{aligned}$$

Recall that $\gamma_p \leq \frac{1}{4\hat{a}_p}$. Then, using Lemma 5, we have

$$\begin{aligned}
f(\hat{x}_{t_{p+1}}) &\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \\
&\quad + \left(\hat{a}_p a_p^2 \gamma_p^2 + \frac{\gamma_p a_p^2}{2} \right) \frac{1}{M(v-t_p)^2} \sum_{m=1}^M \sum_{j=t_p}^v \|x_j^m - \hat{x}_{t_p}\|^2 \\
&\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \\
&\quad + \left(\hat{a}_p a_p^2 \gamma_p^2 + \frac{\gamma_p a_p^2}{2} \right) 16(v-t_p) a_p \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*) \\
&\leq f(\hat{x}_{t_p}) - \frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 + \frac{3(H-1) a_p^3 \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.
\end{aligned}$$

Let us rewrite the inequality in the following way:

$$\frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1) a_p^3 \alpha_p^2 (f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}. \quad (9)$$

Since $\gamma_p \geq \frac{\zeta}{\hat{a}_p}$, we get that

$$\frac{\gamma_p \|\nabla f(\hat{x}_{t_p})\|^2}{4} \geq \frac{\zeta \|\nabla f(\hat{x}_{t_p})\|^2}{4\hat{a}_p}.$$

Therefore,

$$\frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \geq \begin{cases} \frac{\zeta \|\nabla f(\hat{x}_{t_p})\|^2}{4L_0}, & \|\nabla f(\hat{x}_{t_p})\| \leq \frac{L_0}{L_1}, \\ \frac{\zeta \|\nabla f(\hat{x}_{t_p})\|^2}{4L_1}, & \|\nabla f(\hat{x}_{t_p})\| > \frac{L_0}{L_1}, \end{cases} = \frac{\zeta}{4} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\}.$$

Denote $\delta_p \stackrel{\text{def}}{=} f(\hat{x}_{t_p}) - f^*$. Then we have

$$\frac{\zeta}{4} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\} \leq \delta_p - \delta_{p+1} + \frac{3(H-1) a_p^2 \alpha_p^3 (\delta_p + \Delta^*)}{\hat{a}_p}.$$

Let $\alpha_p \leq \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}$, where $c \geq \sqrt{P}$. Applying the result of Mishchenko et al. (2020, Lemma 6), we appear at

$$\begin{aligned}
\min_{0 \leq p \leq P-1} \left\{ \frac{\zeta}{4} \min \left\{ \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{L_0}, \frac{\|\nabla f(\hat{x}_{t_p})\|}{L_1} \right\} \right\} &\leq \frac{\left(1 + \frac{3(H-1) a_p^2 \alpha_p^3}{\hat{a}_p}\right)^P}{P} \delta_0 \\
&\quad + \frac{3(H-1) a_p^2 \alpha_p^3 \Delta^*}{\hat{a}_p}.
\end{aligned}$$

□

B.4 SYMMETRIC GENERALIZED-SMOOTH FUNCTIONS UNDER PŁ-CONDITION

Theorem 6 (Symmetric generalized-smooth convergence analysis of Algorithm 1 in PŁ-case). *Let Assumptions 1 and 3 hold for functions f and $\{f_m\}_{m=1}^M$. Let Assumption 4 hold. Choose $0 < \zeta \leq \frac{1}{4}$.*

Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^$. Choose any integer $P > \frac{64\delta_0 L_1^2}{\mu\zeta}$. For all $0 \leq p \leq P-1$, denote*

$$\hat{a}_p = L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\|, \quad a_p = L_0 + L_1 \max_m \|\nabla f_m(\hat{x}_{t_p})\|, \quad 1 \leq t_{p+1} - t_p \leq H.$$

Put $\Delta^ = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$. Impose the following conditions on the local stepsizes α_p and server stepsizes γ_p :*

$$\alpha_p \leq \min \left\{ \frac{1}{2Ha_p}, \frac{1}{ca_p} \sqrt{\frac{\hat{a}_p}{a_p}}, \sqrt{\frac{\mu\zeta\hat{a}_p}{96L_1^2(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}} \right\},$$

$$\sqrt{\frac{\delta_0 \hat{a}_p}{3P(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}} \Big\},$$

1620

1621

$$\frac{\zeta}{\hat{a}_p} \leq \gamma_p \leq \frac{1}{4\hat{a}_p}, \quad 0 \leq p \leq P-1,$$

1622

1623

where $c \geq \sqrt{P}$. Let \tilde{P} be an integer such that $0 \leq \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$ be a constant, $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$.

1624

1625

Then, the iterates $\{\hat{x}_{t_p}\}_{p=0}^P$ of Algorithm 1 satisfy

1626

1627

$$\delta_P \leq \left(1 - \frac{\mu\zeta}{4L_0}\right)^{P-\tilde{P}} \delta_0 + \frac{4L_0 A \alpha^2}{\mu\zeta},$$

1628

1629

1630

where $\delta_P \stackrel{\text{def}}{=} f(\hat{x}_{t_P}) - f^*$.

1631

1632

Proof of Theorem 6. Let us follow the first steps of the proof of Theorem 5. Consider (9):

1633

1634

$$\frac{\gamma_p}{4} \|\nabla f(\hat{x}_{t_p})\|^2 \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.$$

1635

1636

Since $\gamma_p \geq \frac{\zeta}{\hat{a}_p}$, and f satisfies Polyak–Łojasiewicz Assumption 4, we obtain that

1637

1638

1639

$$\frac{\mu\zeta(f(\hat{x}_{t_p}) - f^*)}{2\hat{a}_p} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.$$

1640

1641

1642

1. Let \tilde{P} be the number of steps p , so that $\|\nabla f(\hat{x}_{t_p})\| \geq \frac{L_0}{L_1}$. For such p , we have $L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\| = \hat{a}_p \leq 2L_1 \|\nabla f(\hat{x}_{t_p})\|$. Therefore, we get

1643

1644

1645

$$\frac{\mu\zeta(f(\hat{x}_{t_p}) - f^*)}{4L_1 \|\nabla f(\hat{x}_{t_p})\|} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.$$

1646

1647

1648

1649

Notice that the relation $\hat{a}_p \leq 2L_1 \|\nabla f(\hat{x}_{t_p})\|$ and Lemma 1 together imply

$$\frac{\|\nabla f(\hat{x}_{t_p})\|}{4L_1} \leq \frac{\|\nabla f(\hat{x}_{t_p})\|^2}{2\hat{a}_p} \leq f(\hat{x}_{t_p}) - f^*.$$

1650

1651

Hence, we have

1652

1653

$$\frac{\mu\zeta}{16L_1^2} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.$$

1654

1655

Subtracting f^* on both sides and introducing $\delta_p \stackrel{\text{def}}{=} f(\hat{x}_{t_p}) - f^*$, we obtain

1656

1657

1658

$$\delta_{p+1} \leq \delta_p - \frac{\mu\zeta}{16L_1^2} + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.$$

1659

1660

1661

As $\alpha_p \leq \sqrt{\frac{\mu\zeta\hat{a}_p}{96L_1^2(H-1)a_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}}$, it follows that $\frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p} \leq \frac{\mu\zeta}{32L_1^2}$.

1662

1663

1664

$$\delta_{p+1} \leq \delta_p - \frac{\mu\zeta}{32L_1^2}.$$

1665

1666

1667

2. Suppose now that $\|\nabla f(\hat{x}_{t_p})\| \leq \frac{L_0}{L_1}$. For such p , we have $L_0 + L_1 \|\nabla f(\hat{x}_{t_p})\| = \hat{a}_p \leq 2L_0$. Hence,

1668

1669

$$\frac{\mu\zeta(f(\hat{x}_{t_p}) - f^*)}{4L_0} \leq f(\hat{x}_{t_p}) - f(\hat{x}_{t_{p+1}}) + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}.$$

1670

1671

Subtracting f^* on both sides and introducing $\delta_p \stackrel{\text{def}}{=} f(\hat{x}_{t_p}) - f^*$, we obtain

1672

1673

$$\delta_{p+1} \leq \delta_p \rho + \frac{3(H-1)a_p^3\alpha_p^2(f(\hat{x}_{t_p}) - f^* + \Delta^*)}{\hat{a}_p}, \quad \text{where } \rho \stackrel{\text{def}}{=} 1 - \frac{\mu\zeta}{4L_0}.$$

Let $\alpha_p \stackrel{\text{def}}{=} \alpha \hat{\alpha}_p$ and $\hat{\alpha}_p \leq \sqrt{\frac{A \hat{a}_p}{3(H-1)\alpha_p^3(f(\hat{x}_{t_p}) - f^* + \Delta^*)}}$ for some constant $A > 0$. Then,

$$\delta_{p+1} \leq \rho \delta_p + A\alpha^2.$$

Unrolling the recursion, we derive

$$\begin{aligned} \delta_P &\leq \rho^{P-\tilde{P}} \delta_0 + A\alpha^2 \sum_{i=0}^{\infty} \rho^i - \frac{\mu\zeta}{32L_1^2} \sum_{i=0}^{N-1} \rho^i \\ &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho} - \frac{1-\rho^{\tilde{P}}}{1-\rho} \frac{\mu\zeta}{32L_1^2}. \end{aligned}$$

Notice that $\delta_{p+1} \leq \delta_p + A\alpha^2$, which implies

$$\delta_P \leq \delta_0 + (P - \tilde{P}) A\alpha^2 - \tilde{P} \frac{\mu\zeta}{32L_1^2}.$$

Since $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$, we conclude that

$$0 \leq \delta_P \leq 2\delta_0 - \tilde{P} \frac{\mu\zeta}{32L_1^2}, \quad \Rightarrow \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}.$$

Therefore, for $P > \frac{64\delta_0 L_1^2}{\mu\zeta}$ we can guarantee that $P - \tilde{P} > 0$ and

$$\begin{aligned} \delta_P &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho} - \tilde{P} \rho^{\tilde{P}} \frac{\mu\zeta}{32L_1^2} \\ &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho}. \end{aligned}$$

□

Corollary 5. Fix $\varepsilon > 0$. Choose $\alpha \leq \min \left\{ \sqrt{\frac{\delta_0}{AP}}, L_1 \sqrt{\frac{8\delta_0 \varepsilon}{L_0 AP}} \right\}$. Then, if $P \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we have $\delta_P \leq \varepsilon$.

Proof of Corollary 5. Since $0 \leq \tilde{P} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$, $\alpha \leq \sqrt{\frac{\delta_0}{AP}}$, $\alpha \leq L_1 \sqrt{\frac{8\delta_0 \varepsilon}{L_0 AP}}$, due to the choice of $P \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we obtain that

$$\left(1 - \frac{\mu\zeta}{4L_0}\right)^{P-\tilde{P}} \delta_0 \leq e^{-\frac{\mu\zeta}{4L_0}(P-\tilde{P})} \delta_0 \leq \frac{\varepsilon}{2},$$

and that

$$\frac{4L_0 A}{\mu\zeta} \cdot \frac{\delta_0}{AP} \leq \frac{\varepsilon}{2}.$$

Therefore, $\delta_P \leq \varepsilon$. □

C RANDOM RESHUFFLING

There are several approaches, that fall under the category of permutation methods, and one of the most popular is **Random Reshuffling (RR)**. In each epoch t of the RR algorithm, we sample indices $\pi_t(1), \dots, \pi_t(M)$ without replacement from the set $\{1, 2, \dots, M\}$. In other words, $\pi_t(1), \dots, \pi_t(M)$ forms a random permutation of $\{1, 2, \dots, M\}$. We then perform M steps in the following manner:

$$x_t^m = x_t^{m-1} - \alpha_t \nabla f_{\pi_t(m)}(x_t^{m-1}), \quad (10)$$

where $f_{\pi_t(m)}$ is the m -th function after permutation π_t on epoch t , and α_t is a stepsize at t -th epoch. If we denote $x_t \equiv x_t^0$, we can rewrite this step as

$$x_t^m = x_t - \alpha_t \sum_{j=1}^m \nabla f_{\pi_t(j)}(x_t^{j-1}).$$

After each epoch we perform additional outer step with stepsize γ_t :

$$x_{t+1} = x_t - \gamma_t g_t, \quad g_t = \frac{1}{\alpha_t} M(x_t^M - x_t) = \frac{1}{M} \sum_{m=1}^M \nabla f_{\pi_t(m)}(x_t^{m-1}). \quad (11)$$

C.1 ASYMMETRIC GENERALIZED-SMOOTH NON-CONVEX FUNCTIONS

Theorem 3 (non-convex asymmetric generalized-smooth convergence analysis of Algorithm 2). *Let Assumptions 1 and 2 hold for functions f and $\{f_m\}_{m=1}^M$. Choose any $T \geq 1$. For all $0 \leq t \leq T-1$, denote*

$$\hat{a}_t = L_0 + L_1 \|\nabla f(x_t)\|, \quad \tilde{a}_t = L_0 + L_1 \max_{m,j} \|\nabla f_{m,j}(x_t)\|.$$

Put $\bar{\Delta}^* = f^* - \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M f_m^*$. Impose the following conditions on the client stepsizes α_t and global stepsizes γ_t :

$$\alpha_t \leq \min \left\{ \frac{\sqrt{2}}{\sqrt{3N(N-1)\tilde{a}_t}}, \frac{\sqrt{\hat{a}_t}}{c\tilde{a}_t^{3/2}} \right\}, \quad \zeta \leq \gamma_t \leq \frac{1}{4\hat{a}_t}, \quad 0 \leq t \leq T-1,$$

where $0 < \zeta \leq \frac{1}{4}$, $c \geq \sqrt{((N-1)(2N-1) + 2(N+1))T}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Then, the iterates $\{x_t\}_{t=0}^{T-1}$ of Algorithm 2 satisfy

$$\begin{aligned} \mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \\ \leq \frac{8 \left(1 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((N-1)(2N-1) + 2(N+1)) \right)^T}{T} \delta_0 + \frac{6\alpha_t^2 \tilde{a}_t^3}{\hat{a}_t} (N+1) \Delta^*. \end{aligned}$$

Lemma 6. Recall that $\tilde{a}_t = L_0 + L_1 \max_{m,j} \|\nabla f_{m,j}(x_t)\|$. Then

$$\frac{\gamma_t^2 \|g_t\|^2}{2} \leq \tilde{a}_t \gamma_t^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2 + \gamma_t^2 \|\nabla f(x_t)\|. \quad (12)$$

Proof.

$$\begin{aligned} \frac{\|g_t\|^2}{2} &= \frac{1}{2} \left\| \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \nabla f_{m,\pi_t(j)}(x_{t,j}^m) \right\|^2 \\ &= \left\| \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M (\nabla f_{m,\pi_t(j)}(x_{t,j}^m) - \nabla f_{m,\pi_t(j)}(x_t)) \right\|^2 + \|\nabla f(x_t)\|^2 \\ &\leq \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_{m,\pi_t(j)}(x_t)\|)^2 \|x_{t,j}^m - x_t\|^2 + \|\nabla f(x_t)\|^2 \\ &\leq (\tilde{a}_t)^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2 + \|\nabla f(x_t)\|^2 \\ &= (\tilde{a}_t)^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2 + \|\nabla f(x_t)\|^2. \end{aligned}$$

□

Lemma 7. Let Assumptions 1 and 2 hold for functions f and $\{f_m\}_{m=1}^M$. Then, if we choose $\alpha_t \leq \frac{\sqrt{2}}{\sqrt{3n(n-1)(\bar{a}_t)}}$, we get

$$\mathbb{E} \left[\frac{1}{M} \sum_{m=1}^M \|x_t^m - x_t\|^2 \middle| x_t \right] \leq 2\alpha_t^2 \bar{a}_t ((N-1)(2N-1) + 2(N+1)(f(x_t) - f^*)) + 4\alpha_t^2 \bar{a}_t (N+1) \bar{\Delta}^*. \quad (13)$$

Proof. From (10) we have

$$x_{t,j}^m = x_{t,j-1}^m - \alpha_t \nabla f_{m,\pi_t(j-1)}(x_{t,j-1}^m) = x_t - \sum_{k=0}^{j-1} \alpha_t \nabla f_{m,\pi_t(k)}(x_{t,k}^m).$$

Thus,

$$\begin{aligned} \|x_{t,j}^m - x_t\|^2 &= \left\| \sum_{k=0}^{j-1} \alpha_t \nabla f_{m,\pi_t(k)}(x_{t,k}^m) \right\|^2 \\ &\leq 2 \left\| \sum_{k=0}^{j-1} \alpha_t (\nabla f_{m,\pi_t(k)}(x_{t,k}^m) - \nabla f_{m,\pi_t(k)}(x_t)) \right\|^2 \\ &\quad + 2 \left\| \sum_{k=0}^{j-1} \alpha_t \nabla f_{m,\pi_t(k)}(x_t) \right\|^2 \\ &\leq 2j \sum_{k=0}^{j-1} (\alpha_t)^2 (L_0 + L_1 \|\nabla f_{m,\pi_t(k)}(x_t)\|)^2 \|x_{t,k}^m - x_t\|^2 \\ &\quad + 2 \left\| \sum_{k=0}^{j-1} \alpha_t \nabla f_{m,\pi_t(k)}(x_t) \right\|^2. \end{aligned}$$

Using last inequality, we get

$$\begin{aligned} \frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 &\leq \sum_{j=0}^{N-1} \frac{2j}{N} \sum_{k=0}^{j-1} (\alpha_t)^2 (L_0 + L_1 \|\nabla f_{m,\pi_t(k)}(x_t)\|)^2 \|x_{t,k}^m - x_t\|^2 \\ &\quad + \frac{2}{N} \sum_{j=0}^{N-1} \left\| \sum_{k=0}^{j-1} \alpha_t \nabla f_{m,\pi_t(k)}(x_t) \right\|^2 \\ &\leq (\alpha_t)^2 (\bar{a}_t)^2 \sum_{j=0}^{N-1} \frac{2j}{N} \sum_{k=0}^{j-1} \|x_{t,k}^m - x_t\|^2 + \frac{2\alpha_t^2}{N} \sum_{j=0}^{N-1} \left\| \sum_{k=0}^{j-1} \nabla f_{m,\pi_t(k)}(x_t) \right\|^2. \end{aligned}$$

Let $\alpha_t \leq \frac{\beta}{\bar{a}_t}$, where β is constant. Then, we take a conditional expectation of the last inequality and get the following

$$\begin{aligned} \mathbb{E} \left[\frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right] &\leq \mathbb{E} \left[\frac{2\beta^2}{N} \sum_{j=0}^{N-1} j \sum_{k=0}^{j-1} \|x_{t,k}^m - x_t\|^2 \middle| x_t \right] \\ &\quad + \frac{2\alpha_t^2}{N} \sum_{j=0}^{N-1} \mathbb{E} \left[\left\| \sum_{k=0}^{j-1} \nabla f_{m,\pi_t(k)}(x_t) \right\|^2 \middle| x_t \right]. \end{aligned}$$

Denote $\sigma_t^2 = \frac{1}{N} \sum_{j=0}^{N-1} \|\nabla f_{m,\pi_t(j)}(x_t) - f(x_t)\|^2$, and consider

$$\mathbb{E} \left[\left\| \sum_{k=0}^{j-1} \nabla f_{m,\pi_t(k)}(x_t) \right\|^2 \middle| x_t \right].$$

1836 From Malinovsky et al. (2022, Lemma 1) we get
1837

$$1838 \mathbb{E} \left[\left\| \sum_{k=0}^{j-1} \nabla f_{m, \pi_t(k)}(x_t) \right\|^2 \middle| x_t \right] \leq j^2 \|\nabla f(x_t)\|^2 + j^2 \mathbb{E} \left[\left\| \frac{1}{j} \sum_{k=0}^{j-1} (\nabla f_{m, \pi_t(k)}(x_t) - \nabla f(x_t)) \right\|^2 \middle| x_t \right]$$

$$1840 \leq j^2 \|\nabla f(x_t)\|^2 + \frac{j(N-j)}{N-1} \sigma_t^2.$$

1841 Thus,
1842

$$1843 \mathbb{E} \left[\frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right] \leq \mathbb{E} \left[\frac{2\beta^2}{N} \sum_{j=0}^{N-1} j \sum_{k=0}^{j-1} \|x_{t,k}^m - x_t\|^2 \middle| x_t \right]$$

$$1844 + \frac{2\alpha_t^2}{N} \sum_{j=0}^{N-1} \left(j^2 \|\nabla f(x_t)\|^2 + \frac{j(N-j)}{N-1} \sigma_t^2 \right)$$

$$1845 \leq \mathbb{E} \left[\frac{2\beta^2}{N} \cdot \frac{N(N-1)}{2} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right]$$

$$1846 + \frac{2\alpha_t^2}{N} \left(\frac{(N(N-1)(2N-1))}{6} \|\nabla f(x_t)\|^2 \right)$$

$$1847 + \frac{2\alpha_t^2}{N} \frac{N(N+1)}{3} \sigma_t^2.$$

1848 Further,
1849

$$1850 3 \cdot \mathbb{E} \left[\frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right] \leq 3\beta^2 N(N-1) \mathbb{E} \left[\frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right]$$

$$1851 + 2\alpha_t^2 \left(\frac{(N-1)(2N-1)}{2} \|\nabla f(x_t)\|^2 + (N+1)\sigma_t^2 \right).$$

1852 Thus, if we choose $\beta \leq \sqrt{\frac{2}{3N(N-1)}}$, we get
1853

$$1854 \mathbb{E} \left[\frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right] \leq (3 - 3\beta^2 N(N-1)) \mathbb{E} \left[\frac{1}{N} \sum_{j=0}^{N-1} \|x_{t,j}^m - x_t\|^2 \middle| x_t \right]$$

$$1855 \leq 2\alpha_t^2 \left(\frac{(N-1)(2N-1)}{2} \|\nabla f(x_t)\|^2 + (N+1)\sigma_t^2 \right)$$

$$1856 \leq 2\alpha_t^2 \left((N-1)(2N-1)(f(x_t) - f^*)(L_0 + L_1 \|\nabla f(x_t)\|) \right)$$

$$1857 + (N+1) \frac{1}{N} \sum_{j=0}^{N-1} \|\nabla f_{m, \pi_t(j)}(x_t)\|^2$$

$$1858 \stackrel{\text{Lemma 1}}{\leq} 2\alpha_t^2 \left((N-1)(2N-1)(f(x_t) - f^*)\tilde{a}_t \right)$$

$$1859 + 2(N+1)(\tilde{a}_t) \frac{1}{N} \sum_{j=0}^{N-1} (f_{m_j}(x_t) - f_{m_j}^*)$$

$$\leq 2\alpha_t^2 \left((N-1)(2N-1)(f(x_t) - f^*)\tilde{a}_t \right)$$

$$+ 2(N+1)(\tilde{a}_t) \frac{1}{N} \sum_{j=0}^{N-1} (f_{m_j}(x_t) - f_{m_j}^*)$$

Now, adding and removing f^* to the sum factor on the right-hand side, we get

$$\begin{aligned} \mathbb{E} \left[\frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2 \middle| x_t \right] &\leq 2\alpha_t^2 \tilde{a}_t (N-1)(2N-1)(f(x_t) - f^*) \\ &\quad + 4\alpha_t^2 \tilde{a}_t (N+1) \frac{1}{NM} \sum_{j=0}^{N-1} \sum_{m=1}^M (f_{mj}(x_t) - f_{mj}^*) \\ &= 2\alpha_t^2 \tilde{a}_t ((N-1)(2N-1) + 2(N+1)(f(x_t) - f^*)) \\ &\quad + 4\alpha_t^2 \tilde{a}_t (N+1) \overline{\Delta}^*. \end{aligned}$$

□

Proof of Theorem 3. From Lemma 1 and (11) we get

$$f(x_{t+1}) \leq f(x_t) - \gamma_t \langle \nabla f(x_t), g_t \rangle + (L_0 + L_1 \|\nabla f(x_t)\|) \frac{\gamma_t^2 \|g_t\|^2}{2}.$$

Additionally, from the fact that $2 \langle a, b \rangle = -\|a - b\|^2 + \|a\|^2 + \|b\|^2$ we can get

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) - \gamma_t \langle \nabla f(x_t), g_t \rangle + (L_0 + L_1 \|\nabla f(x_t)\|) \frac{\gamma_t^2 \|g_t\|^2}{2} \\ &\leq f(x_t) - \frac{\gamma_t}{2} (-\|\nabla f(x_t) - g_t\|^2 + \|\nabla f(x_t)\|^2 + \|g_t\|^2) \\ &\quad + (L_0 + L_1 \|\nabla f(x_t)\|) \frac{\gamma_t^2 \|g_t\|^2}{2} \\ &\leq f(x_t) - \frac{\gamma_t}{2} \|\nabla f(x_t)\|^2 + \frac{\gamma_t}{2} \|\nabla f(x_t) - g_t\|^2 + (L_0 + L_1 \|\nabla f(x_t)\|) \frac{\gamma_t^2 \|g_t\|^2}{2}. \end{aligned}$$

Consider $\frac{\gamma_t}{2} \|\nabla f(x_t) - g_t\|^2$ and denote $\hat{a}_t = (L_0 + L_1 \|\nabla f(x_t)\|)$ and $a_t = (L_0 + L_1 \max_m \|\nabla f_m(x_t)\|)$, then:

$$\begin{aligned} \frac{\gamma_t}{2} \|\nabla f(x_t) - g_t\|^2 &= \frac{\gamma_t}{2} \left\| \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \nabla f_{m,\pi_t(j)}(x_t) - \nabla f_{m,\pi_t(j)}(x_{t,j}^m) \right\|^2 \\ &\leq \frac{\gamma_t}{2} \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_{m,\pi_t(j)}(x_t)\|)^2 \|x_t - x_{t,j}^m\|^2 \\ &= \frac{\gamma_t}{2} \tilde{a}_t^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_t - x_{t,j}^m\|^2. \end{aligned}$$

From the above inequality and Lemma 6 we get

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) - \frac{\gamma_t}{2} \|\nabla f(x_t)\|^2 + \frac{\gamma_t}{2} \tilde{a}_t^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_t - x_{t,j}^m\|^2 + \hat{a}_t \frac{\gamma_t^2 \|g_t\|^2}{2} \\ &\stackrel{(12)}{\leq} f(x_t) - \frac{\gamma_t}{2} \|\nabla f(x_t)\|^2 + \frac{\gamma_t}{2} \tilde{a}_t^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_t - x_{t,j}^m\|^2 \\ &\quad + \hat{a}_t \tilde{a}_t^2 \gamma_t^2 \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2 + \hat{a}_t \gamma_t^2 \|\nabla f(x_t)\|^2 \\ &\leq f(x_t) + \left(\hat{a}_t \gamma_t^2 - \frac{\gamma_t}{2} \right) \|\nabla f(x_t)\|^2 + \left(\hat{a}_t \tilde{a}_t^2 \gamma_t^2 + \frac{\gamma_t}{2} \tilde{a}_t^2 \right) \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2. \end{aligned}$$

Let $\gamma_t \leq \frac{1}{4\hat{a}_t}$, then

$$f(x_{t+1}) \leq f(x_t) - \frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 + \left(\hat{a}_t \tilde{a}_t^2 \gamma_t^2 + \frac{\gamma_t}{2} \tilde{a}_t^2 \right) \frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2.$$

Now, if we take conditional expectation of this and use Lemma 7, we get

$$\begin{aligned} \mathbb{E}[f(x_{t+1})|x_t] &\leq f(x_t) - \frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 \\ &\quad + \left(\hat{a}_t \tilde{a}_t^2 \gamma_t^2 + \frac{\gamma_t}{2} \tilde{a}_t^2 \right) \mathbb{E} \left[\frac{1}{MN} \sum_{j=0}^{N-1} \sum_{m=1}^M \|x_{t,j}^m - x_t\|^2 \middle| x_t \right] \\ &\leq f(x_t) - \frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 \\ &\quad + 2\alpha_t^2 \tilde{a}_t \left(\hat{a}_t \tilde{a}_t^2 \gamma_t^2 + \frac{\gamma_t}{2} \tilde{a}_t^2 \right) \\ &\quad \times ((f(x_t) - f^*)((N-1)(2N-1) + 2(N+1)(f(x_t) - f^*) + 2(N+1)\bar{\Delta}^*)). \end{aligned}$$

Since $\gamma_t \leq \frac{1}{4\hat{a}_t}$, then

$$\begin{aligned} \frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 &\leq f(x_t) - \mathbb{E}[f(x_{t+1})|x_t] + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \\ &\quad \times ((f(x_t) - f^*)((N-1)(2N-1) + 2(N+1)\delta_t + 2(N+1)\bar{\Delta}^*)). \end{aligned} \quad (14)$$

Consider the left-hand side of (14). Due to the bounds $\frac{1}{4\hat{a}_t} \geq \gamma_t \geq \frac{\zeta}{\hat{a}_t}$ on γ_t , we have

$$\frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 \geq \frac{\zeta \|\nabla f(x_t)\|^2}{4\hat{a}_t}.$$

Then, we get

$$\begin{aligned} \frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 &\geq \begin{cases} \frac{\zeta \|\nabla f(x_t)\|^2}{8L_0}, & \|\nabla f(x_t)\| \leq \frac{L_0}{L_1}, \\ \frac{\zeta \|\nabla f(x_t)\|}{8L_1}, & \|\nabla f(x_t)\| > \frac{L_0}{L_1} \end{cases} \\ &= \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \end{aligned} \quad (15)$$

Denote $\delta_t \equiv f(x_t) - f^*$, then from (14) and (15) we get

$$\begin{aligned} &f(x_t) - f(x_{t+1}) \\ &\quad + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((f(x_t) - f^*)((N-1)(2N-1) + 2(N+1)\delta_t + 2(N+1)\bar{\Delta}^*)) \\ &= \delta_t - \delta_{t+1} + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((f(x_t) - f^*)((N-1)(2N-1) + 2(N+1)\delta_t + 2(N+1)\bar{\Delta}^*)) \\ &\geq \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \end{aligned}$$

Let $\alpha_t \leq \frac{1}{c\hat{a}_t} \cdot \sqrt{\frac{\hat{a}_t}{c}}$, where c is a constant such that $\sqrt{((N-1)(2N-1) + 2(N+1))T} \leq c$. Now we take full expectation and use from Mishchenko et al. (2020, Lemma 6):

$$\begin{aligned} &\mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \\ &\leq \frac{\left(1 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((N-1)(2N-1) + 2(N+1)) \right)^T}{T} \delta_0 + \frac{3\alpha_t^2 \tilde{a}_t^3}{4\hat{a}_t} (N+1)\bar{\Delta}^*. \end{aligned}$$

□

Corollary 3. Fix $\varepsilon > 0$. Choose $c = \sqrt{((N-1)(2N-1) + 2(N+1))T}$. Let $\alpha_t \leq 8\sqrt{\frac{\hat{a}_t}{3\tilde{a}_t^3 T(N+1)\Delta^*}}$. Then, if $T \geq \frac{256\delta_0}{\zeta\varepsilon}$, we have

$$\mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \leq \varepsilon.$$

Proof of Corollary 3. Since $c = \sqrt{((N-1)(2N-1) + 2(N+1))T}$ and $\alpha_t \leq \frac{1}{ca_t} \sqrt{\frac{\hat{a}_t}{a_t}}$, $\alpha_t \leq 8\sqrt{\frac{\hat{a}_t}{3\tilde{a}_t^3 T(N+1)\Delta^*}}$, due to the choice of $T \geq \frac{256\delta_0}{\zeta\varepsilon}$, we obtain that

$$\frac{\left(1 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((N-1)(2N-1) + 2(N+1))\right)^T}{T} \delta_0 \leq \frac{e^{\frac{3}{8}} \delta_0}{T} \leq \frac{2\delta_0}{T} \leq \frac{\zeta\varepsilon}{16},$$

and that

$$\frac{3\alpha_t^2 \tilde{a}_t^3}{4\hat{a}_t} (N+1)\Delta^* \leq \frac{\zeta\varepsilon}{16}.$$

Therefore, $\mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \leq \varepsilon$. \square

C.2 ASYMMETRIC GENERALIZED-SMOOTH FUNCTIONS UNDER PŁ-CONDITION

Theorem 7. Let Assumptions 1 and 2 hold for functions f , $\{f_m\}_{m=1}^M$ and $\{f_{m,j}\}_{m=1, j=0}^{M, N-1}$. Let Assumption 4 hold. Choose $0 < \zeta \leq \frac{1}{4}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Choose any integer $T > \frac{64\delta_0 L_1^2}{\mu\zeta}$. For all $0 \leq t \leq T-1$, denote

$$\hat{a}_t = L_0 + L_1 \|\nabla f(x_t)\|, \quad a_t = L_0 + L_1 \max_m \|\nabla f_m(x_t)\|.$$

Put $\bar{\Delta}^* = f^* - \frac{1}{MN} \sum_{m=1}^M \sum_{j=0}^{N-1} f_{m,j}^*$. Impose the following conditions on the client stepsizes α_t and global stepsizes γ_t :

$$\alpha_t \leq \min \left\{ \frac{\sqrt{2}}{\sqrt{3M(M-1)\tilde{a}_t}}, \frac{\sqrt{\hat{a}_t}}{c\tilde{a}_t^{3/2}}, \sqrt{\frac{\hat{a}_t \mu \zeta}{12L_1^2 \tilde{a}_t^3 \left(\delta_t ((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right)}}, \sqrt{\frac{8\hat{a}_t \delta_0}{3T \tilde{a}_t^3 \left(\delta_t ((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right)}} \right\},$$

$$\frac{\zeta}{\hat{a}_t} \leq \gamma_t \leq \frac{1}{4\hat{a}_t}, \quad 0 \leq t \leq T-1,$$

where $c \geq \sqrt{((N-1)(2N-1) + 2(N+1))T}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Let \tilde{T} be an integer such that $0 \leq \tilde{T} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$ be a constant, $\alpha \leq \sqrt{\frac{\delta_0}{AT}}$. Then, the iterates $\{x_t\}_{t=0}^{T-1}$ of Algorithm 2 satisfy

$$\delta_T \leq \left(1 - \frac{\mu\zeta}{4L_0}\right)^{T-\tilde{T}} \delta_0 + \frac{4L_0 A \alpha^2}{\mu\zeta},$$

where $\delta_T \stackrel{\text{def}}{=} f(x_T) - f^*$.

Proof of Theorem 7. Let us follow the first steps of the proof of Theorem 3. Consider (14):

$$\begin{aligned} \frac{\gamma_t}{4} \|\nabla f(x_t)\|^2 &\leq f(x_t) - f(x_{t+1}) \\ &+ \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} ((f(x_t) - f^*)((M-1)(2M-1) + 2(M+1)) + 2(M+1)\Delta^*). \end{aligned}$$

2052 Since $\gamma_p \geq \frac{\zeta}{\hat{a}_p}$, and f satisfies Polyak–Łojasiewicz Assumption 4, we obtain that
 2053

$$2054 \frac{\mu\zeta (f(\hat{x}_{t_p}) - f^*)}{2\hat{a}_p} \leq f(x_t) - f(x_{t+1})$$

$$2055 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right).$$

2059 **1.** Let \tilde{T} be the number of steps t , so that $\|\nabla f(\hat{x}_t)\| \geq \frac{L_0}{L_1}$. For such t , we have $L_0 + L_1 \|\nabla f(x_t)\| =$
 2060 $\hat{a}_t \leq 2L_1 \|\nabla f(x_t)\|$. Therefore, we get
 2061

$$2062 \frac{\mu\zeta (f(x_t) - f^*)}{4L_1 \|\nabla f(x_t)\|} \leq f(x_t) - f(x_{t+1})$$

$$2063 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right).$$

2066 Notice that the relation $\hat{a}_t \leq 2L_1 \|\nabla f(x_t)\|$ and Lemma 1 together imply
 2067

$$2068 \frac{\|\nabla f(x_t)\|}{4L_1} \leq \frac{\|\nabla f(x_t)\|^2}{2\hat{a}_t} \leq f(x_t) - f^*.$$

2070 Hence, we have
 2071

$$2072 \frac{\mu\zeta}{16L_1^2} \leq f(x_t) - f(x_{t+1})$$

$$2073 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right).$$

2077 Subtracting f^* on both sides and introducing $\delta_t \stackrel{\text{def}}{=} f(x_t) - f^*$, we obtain
 2078

$$2079 \delta_{t+1} \leq \delta_t - \frac{\mu\zeta}{16L_1^2}$$

$$2080 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right).$$

2083 As $\alpha_t \leq \sqrt{\frac{\hat{a}_t \mu \zeta}{12L_1^2 \tilde{a}_t^3 (\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^*)}}$, it follows that
 2084

$$2085 \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right) \leq \frac{\mu\zeta}{32L_1^2}.$$

2088 Therefore, we get
 2089

$$2090 \delta_{t+1} \leq \delta_t - \frac{\mu\zeta}{32L_1^2}.$$

2092 **2.** Suppose now that $\|\nabla f(x_t)\| \leq \frac{L_0}{L_1}$. For such t , we have $L_0 + L_1 \|\nabla f(x_t)\| = \hat{a}_p \leq 2L_0$. Hence,
 2093

$$2094 \frac{\mu\zeta (f(x_t) - f^*)}{4L_0} \leq f(x_t) - f(x_{t+1})$$

$$2095 + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right).$$

2099 Subtracting f^* on both sides and introducing $\delta_t \stackrel{\text{def}}{=} f(x_t) - f^*$, we obtain
 2100

$$2101 \delta_{t+1} \leq \delta_t \rho + \frac{3\alpha_t^2 \tilde{a}_t^3}{8\hat{a}_t} \left(\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^* \right).$$

2103 where $\rho \stackrel{\text{def}}{=} 1 - \frac{\mu\zeta}{4L_0}$. Let $\alpha_t \stackrel{\text{def}}{=} \alpha \hat{\alpha}_t$ with $\hat{\alpha}_t \leq \sqrt{\frac{8\hat{a}_t A}{3\tilde{a}_t^3 (\delta_t((N-1)(2N-1) + 2(N+1)) + 2(N+1)\bar{\Delta}^*)}}$ for
 2104 some constant $A > 0$. Then,
 2105

$$\delta_{t+1} \leq \rho \delta_t + A \alpha^2.$$

Unrolling the recursion, we derive

$$\begin{aligned}\delta_T &\leq \rho^{T-\tilde{T}}\delta_0 + A\alpha^2 \sum_{i=0}^{\infty} \rho^i - \frac{\mu\zeta}{32L_1^2} \sum_{i=0}^{N-1} \rho^i \\ &\leq \rho^{P-\tilde{P}}\delta_0 + \frac{A\alpha^2}{1-\rho} - \frac{1-\rho^{\tilde{P}}}{1-\rho} \frac{\mu\zeta}{32L_1^2}.\end{aligned}$$

Notice that $\delta_{t+1} \leq \delta_t + A\alpha^2$, which implies

$$\delta_T \leq \delta_0 + (T - \tilde{T})A\alpha^2 - \tilde{T} \frac{\mu\zeta}{32L_1^2}.$$

Since $\alpha \leq \sqrt{\frac{\delta_0}{AT}}$, we conclude that

$$0 \leq \delta_T \leq 2\delta_0 - \tilde{T} \frac{\mu\zeta}{32L_1^2}, \quad \Rightarrow \tilde{T} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}.$$

Therefore, for $T > \frac{64\delta_0 L_1^2}{\mu\zeta}$ we can guarantee that $T - \tilde{T} > 0$ and

$$\begin{aligned}\delta_T &\leq \rho^{T-\tilde{T}}\delta_0 + \frac{A\alpha^2}{1-\rho} - \tilde{T} \frac{\mu\zeta}{32L_1^2} \\ &\leq \rho^{T-\tilde{T}}\delta_0 + \frac{A\alpha^2}{1-\rho}.\end{aligned}$$

□

Corollary 6. Fix $\varepsilon > 0$. Choose $\alpha \leq \min \left\{ \sqrt{\frac{\delta_0}{AT}}, L_1 \sqrt{\frac{8\delta_0\varepsilon}{L_0AT}} \right\}$. Then, if $T \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we have $\delta_T \leq \varepsilon$.

Proof of Corollary 6. Since $0 \leq \tilde{T} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$, $\alpha \leq \sqrt{\frac{\delta_0}{AT}}$, $\alpha \leq L_1 \sqrt{\frac{8\delta_0\varepsilon}{L_0AT}}$, due to the choice of $T \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we obtain that

$$\left(1 - \frac{\mu\zeta}{4L_0}\right)^{T-\tilde{T}} \delta_0 \leq e^{-\frac{\mu\zeta}{4L_0}(T-\tilde{T})} \delta_0 \leq \frac{\varepsilon}{2},$$

and that

$$\frac{4L_0 A}{\mu\zeta} \cdot \frac{\delta_0}{AT} \leq \frac{\varepsilon}{2}.$$

Therefore, $\delta_T \leq \varepsilon$. □

D PARTIAL PARTICIPATION

D.1 ASYMMETRIC GENERALIZED-SMOOTH NON-CONVEX FUNCTIONS

Theorem 4 Let Assumptions 1 and 2 hold for functions f , $\{f_m\}_{m=1}^M$ and $\{f_{m,j}\}_{m=1,j=1}^{M,N}$. Choose any $T \geq 1$. For all $0 \leq t \leq T-1$, denote

$$\hat{a}_t = L_0 + L_1 \|\nabla f(x_t)\|, \quad a_t = L_0 + L_1 \max_m \|\nabla f_m(x_t)\|, \quad \tilde{a}_t = L_0 + L_1 \max_{m,j} \|\nabla f_m^{\pi_j}(x_t)\|.$$

Put $\Delta^* = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$ and $\bar{\Delta}^* = f^* - \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} f_{m,j}^*$. Impose the following conditions on the local stepsizes γ_t , server stepsizes η_t , global stepsizes θ_t :

$$\gamma_t N R \leq \eta_t R \leq \min \left\{ \frac{1}{16\hat{a}_t}, \frac{2\hat{a}_t}{c} \sqrt{\frac{1}{a_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}} \right\}, \quad \gamma_t \leq \frac{2\hat{a}_t}{cRN} \sqrt{\frac{1}{\tilde{a}_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}},$$

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

$$\frac{\zeta}{\hat{a}_t} \leq \theta_t \leq \frac{1}{4\hat{a}_t}, \quad 0 \leq t \leq T-1,$$

where $c \geq \sqrt{T}$, $0 < \zeta \leq \frac{1}{4}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Then, the iterates $\{x_t\}_{t=0}^{T-1}$ of Algorithm 3 satisfy

$$\begin{aligned} \mathbb{E} \left[\min_{0 \leq t \leq T-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \\ \leq \frac{\left(1 + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t)\right)^T}{T} \delta_0 \\ + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \left(\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^* \right). \end{aligned}$$

We need to use the following relations to establish convergence guarantees:

$$\begin{aligned} x_t^R &= x_t - \eta_t \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j} (x_{m,t}^{r,j}), \\ x_{m,t}^{r,j} &= x_t - \eta_t \sum_{k=0}^{r-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^k}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j} (x_{m,t}^{k,j}) - \gamma_t \sum_{l=0}^{j-1} \nabla f_m^{\pi^l} (x_{m,t}^{r,l}), \\ x_{t+1} &= x_t - \frac{\theta_t}{\eta_t R} (x_t - x_t^R). \end{aligned}$$

We assume that the whole sum is zero when the upper summation index is smaller than the lower index. We can derive the following recursion from the above relations:

$$\begin{aligned} x_t - x_{t+1} &= \frac{\theta_t}{\eta_t R} (x_t - x_t^R) \\ &= \frac{\theta_t}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j} (x_{m,t}^{r,j}). \end{aligned}$$

Further, the first statement of Lemma 1 yields the following inequality:

$$f(x_{t+1}) \leq f(x_t) - \langle \nabla f(x_t), x_t - x_{t+1} \rangle + (L_0 + L_1 \|\nabla f(x_t)\|) \frac{\|x_t - x_{t+1}\|^2}{2}.$$

We deal with the last term, using the second statement of Lemma 1:

$$\begin{aligned} \|x_t - x_{t+1}\|^2 &= \theta_t^2 \left\| \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j} (x_{m,t}^{r,j}) \right\|^2 \\ &\leq 2\theta_t^2 \left\| \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} (\nabla f_m^{\pi^j} (x_{m,t}^{r,j}) - \nabla f(x_t)) \right\|^2 + 2\theta_t^2 \|\nabla f(x_t)\|^2 \\ &\leq \frac{2\theta_t^2}{RCN} (L_0 + L_1 \|\nabla f(x_t)\|)^2 \sum_{r=0}^{R-1} \sum_{m \in S_t^{\lambda^r}} \sum_{j=0}^{N-1} \|x_t - x_{m,t}^{r,j}\|^2 \\ &\quad + 4\theta_t^2 (L_0 + L_1 \|\nabla f(x_t)\|) (f(x_t) - f^*). \end{aligned}$$

We use the following notation: $\hat{a}_t = L_0 + L_1 \|\nabla f(x_t)\|$, $a_t = L_0 + L_1 \max_m \|\nabla f_m(x_t)\|$, $\tilde{a}_t = L_0 + L_1 \max_{m,j} \|\nabla f_m^{\pi^j}(x_t)\|$. Next, we have that

$$\begin{aligned} \|x_t - x_{m,t}^{r,j}\|^2 &= \left\| \eta_t \sum_{k=0}^{r-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^k}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j} (x_{m,t}^{k,j}) + \gamma_t \sum_{l=0}^{j-1} \nabla f_m^{\pi^l} (x_{m,t}^{r,l}) \right\|^2 \\ &\leq 2 \left\| \eta_t \sum_{k=0}^{r-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^k}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j} (x_{m,t}^{k,j}) \right\|^2 + 2 \left\| \gamma_t \sum_{l=0}^{j-1} \nabla f_m^{\pi^l} (x_{m,t}^{r,l}) \right\|^2. \end{aligned}$$

Using Young's inequality, we obtain

$$\begin{aligned}
\|x_t - x_{m,t}^{r,j}\|^2 &\leq 4\eta_t^2 \left\| \sum_{k=0}^{r-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^k}} \frac{1}{N} \sum_{j=0}^{N-1} \left(\nabla f_m^{\pi^j}(x_{m,t}^{k,j}) - \nabla f_m^{\pi^j}(x_t) \right) \right\|^2 \\
&\quad + 4\eta_t^2 \left\| \sum_{k=0}^{r-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^k}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi^j}(x_t) \right\|^2 \\
&\quad + 4\gamma_t^2 \left\| \sum_{l=0}^{j-1} \left(\nabla f_m^{\pi^l}(x_{m,t}^{r,l}) - \nabla f_m^{\pi^l}(x_t) \right) \right\|^2 \\
&\quad + 4\gamma_t^2 \left\| \sum_{l=0}^{j-1} \nabla f_m^{\pi^l}(x_t) \right\|^2.
\end{aligned}$$

Using Malinovsky et al. (2022, Lemma 1), we derive the following upper bound on $\|x_t - x_{m,t}^{r,j}\|^2$:

$$\begin{aligned}
\|x_t - x_{m,t}^{r,j}\|^2 &\leq 4\eta_t^2 r^2 (\hat{\alpha}_t)^2 \frac{1}{rCN} \sum_{k=0}^{r-1} \sum_{m \in S_t^{\lambda^k}} \sum_{j=0}^{N-1} \|x_t - x_{m,t}^{k,j}\|^2 \\
&\quad + 4\eta_t^2 \frac{1}{N^2 C^2} \left(N^2 C^2 r^2 \|\nabla f(x_t)\|^2 + \frac{Cr(M-Cr)}{M-1} \sigma_t^2 \right) \\
&\quad + 4\gamma_t^2 j (\hat{\alpha}_t)^2 \sum_{l=0}^{j-1} \|x_t - x_{m,t}^{r,l}\|^2 \\
&\quad + 4\gamma_t^2 \left(j^2 \|\nabla f_m(x_t)\|^2 + \frac{j(N-j)}{N-1} \sigma_{m,t}^2 \right),
\end{aligned}$$

where

$$\begin{aligned}
\sigma_t^2 &= \frac{1}{MN} \sum_{m=1}^M \sum_{j=0}^{N-1} \|\nabla f_m^{\pi^j}(x_t) - \nabla f_m(x_t)\|^2, \\
\sigma_{m,t}^2 &= \frac{1}{N} \sum_{j=0}^{N-1} \|\nabla f_m^{\pi^j}(x_t) - \nabla f_m(x_t)\|^2.
\end{aligned}$$

Using this bound on $\|x_t - x_{m,t}^{r,j}\|^2$, for $V_t \stackrel{\text{def}}{=} \frac{1}{CRN} \sum_{r=0}^{R-1} \sum_{m \in S_t^{\lambda^r}} \sum_{j=0}^{N-1} \|x_{m,t}^{r,j} - x_t\|^2$, we obtain

$$\begin{aligned}
\mathbb{E}[V_t] &= \frac{1}{CRN} \sum_{r=0}^{R-1} \sum_{m \in S_t^{\lambda^r}} \sum_{j=0}^{N-1} \mathbb{E} \|x_{m,t}^{r,j} - x_t\|^2 \\
&\leq \frac{(\hat{\alpha}_t)^2}{CRN} \\
&\quad \times \sum_{r=0}^{R-1} \sum_{m \in S_t^{\lambda^r}} \sum_{j=0}^{N-1} \left(4\eta_t^2 r^2 \frac{1}{rCN} \sum_{k=0}^{r-1} \sum_{m \in S_t^{\lambda^k}} \sum_{j=0}^{N-1} \|x_t - x_{m,t}^{k,j}\|^2 + 4\gamma_t^2 j \sum_{l=0}^{j-1} \|x_{m,t}^{r,l} - x_t\|^2 \right) \\
&\quad + \frac{1}{CRN} \sum_{r=0}^{R-1} \sum_{m \in S_t^{\lambda^r}} \sum_{j=0}^{N-1} \left(4\gamma_t^2 \left(j^2 \|\nabla f_m(x_t)\|^2 + \frac{j(N-j)}{N-1} \sigma_{m,t}^2 \right) \right) \\
&\quad + \frac{1}{CRN} \sum_{r=0}^{R-1} \sum_{m \in S_t^{\lambda^r}} \sum_{j=0}^{N-1} \left(4\eta_t^2 \frac{1}{N^2 C^2} \left(N^2 C^2 r^2 \|\nabla f(x_t)\|^2 + \frac{Cr(M-Cr)}{M-1} \sigma_t^2 \right) \right).
\end{aligned}$$

2268 Recall that $\gamma_t NR \leq \eta_t R \leq \frac{1}{16\hat{a}_t}$. Summing over indices, we arrive at

$$\begin{aligned}
2270 \mathbb{E}[V_t] &\leq \frac{R(R-1)}{2} 4\eta_t^2 (\hat{a}_t)^2 \mathbb{E}[V_t] + \frac{M(M-1)}{2} 4\gamma_t^2 (\hat{a}_t)^2 \mathbb{E}[V_t] \\
2272 &+ \frac{2}{3}\gamma_t^2 \frac{1}{M} \sum_{m=1}^M \|\nabla f_m(x_t)\|^2 (N-1)(2M-1) + \frac{2}{3}\gamma_t^2 (N+1) \frac{1}{M} \sum_{m=1}^M \sigma_{m,t}^2 \\
2273 &+ \frac{2}{3}\eta_t^2 \|\nabla f(x_t)\|^2 (R-1)(2R-1) + \frac{2}{3} \frac{M-C}{(M-1)C} \eta_t^2 \frac{R+1}{N^2} \sigma_t^2 \\
2274 & \\
2275 &\leq 2\eta_t^2 (\hat{a}_t)^2 (1+R^2) \mathbb{E}[V_t] + \frac{2}{3}\gamma_t^2 \frac{1}{M} \sum_{m=1}^M \|\nabla f_m(x_t)\|^2 (N-1)(2M-1) \\
2276 & \\
2277 &+ \frac{2}{3}\eta_t^2 \|\nabla f(x_t)\|^2 (R-1)(2R-1) + \frac{2}{3}\gamma_t^2 (N+1) \frac{1}{M} \sum_{m=1}^M \sigma_{m,t}^2 \\
2278 & \\
2279 &+ \frac{2}{3}\eta_t^2 \frac{R+1}{N^2} \frac{M-C}{(M-1)C} \sigma_t^2.
\end{aligned}$$

2285 To derive the bound on $\mathbb{E}[V_t]$ we need to require that $\gamma_t NR \leq \eta_t R \leq \frac{1}{16\hat{a}_t}$ to have $1 - 2\eta_t^2 (\hat{a}_t)^2 (1 + R^2) > 0$. Using Lemma 1, we have

$$\begin{aligned}
2288 \mathbb{E}[V_t] &\leq 2\gamma_t^2 N^2 \frac{1}{M} \sum_{m=1}^M \|\nabla f_m(x_t)\|^2 + 2\eta_t^2 R^2 \|\nabla f(x_t)\|^2 \\
2289 &+ 2\gamma_t^2 N \frac{1}{M} \sum_{m=1}^M \sigma_{m,t}^2 + 2\eta_t^2 \frac{R}{N^2} \frac{M-C}{(M-1)C} \sigma_t^2 \\
2290 & \\
2291 &\leq 4\gamma_t^2 N^2 \frac{1}{M} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(x_t)\|) (f_m(x_t) - f_m^*) \\
2292 & \\
2293 &+ 4\eta_t^2 R^2 \hat{a}_t (f(x_t) - f(x_*)) + 2\gamma_t^2 N \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} \|\nabla f_m^{\pi_j}(x_t)\|^2 \\
2294 & \\
2295 &+ 2\eta_t^2 R \frac{M-C}{(M-1)C} \frac{1}{M} \sum_{m=1}^M \|\nabla f_m(x_t)\|^2 \\
2296 & \\
2297 &\leq 4\eta_t^2 \frac{1}{M} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(x_t)\|) (f_m(x_t) - f_m^*) \\
2298 & \\
2299 &+ 4\eta_t^2 R^2 \hat{a}_t (f(x_t) - f(x_*)) \\
2300 & \\
2301 &+ 4\gamma_t^2 N \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} (L_0 + L_1 \|\nabla f_m^{\pi_j}(x_t)\|) (f_m^{\pi_j}(x_t) - f_m^{\pi_j,*}) \\
2302 & \\
2303 &+ 4\eta_t^2 R \frac{M-C}{(M-1)C} \frac{1}{M} \sum_{m=1}^M (L_0 + L_1 \|\nabla f_m(x_t)\|) (f_m(x_t) - f_m^*).
\end{aligned}$$

2312 The bound for $\mathbb{E}[V_t]$ is given by the following:

$$\begin{aligned}
2313 \mathbb{E}[V_t] &\leq 4\eta_t^2 a_t \left(f(x_t) - f^* + \left(f^* - \frac{1}{M} \sum_{m=1}^M f_m^* \right) \right) + 4\eta_t^2 R^2 \hat{a}_t (f(x_t) - f^*) \\
2314 & \\
2315 &+ 4\gamma_t^2 N \tilde{a}_t \left(f(x_t) - f^* + \left(f^* - \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} f_m^{\pi_j,*} \right) \right) \\
2316 & \\
2317 &+ 4\eta_t^2 R a_t \frac{M-C}{(M-1)C} \left(f(x_t) - f^* + \left(f^* - \frac{1}{M} \sum_{m=1}^M f_m^* \right) \right).
\end{aligned}$$

2320

2321

2322 Recall that $\Delta^* = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$, $\bar{\Delta}^* = f^* - \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} f_m^{\pi_j, *}$. Therefore,

$$2323 \mathbb{E}[V_t] \leq 4\eta_t^2 a_t (f(x_t) - f^* + \Delta^*) + 4\eta_t^2 R^2 \hat{a}_t (f(x_t) - f^*)$$

$$2324 + 4\gamma_t^2 N \tilde{a}_t (f(x_t) - f^* + \bar{\Delta}^*) + 4\eta_t^2 R a_t \frac{M-C}{(M-1)C} (f(x_t) - f^* + \Delta^*).$$

2327 Rewriting, we obtain

$$2329 \mathbb{E}[V_t] \leq 4(f(x_t) - f^*) \left(\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t \frac{M-C}{(M-1)C} \right)$$

$$2330 + 4\eta_t^2 a_t \Delta^* + 4\gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + 4\eta_t^2 R a_t \frac{M-C}{(M-1)C} \Delta^*$$

$$2331 \leq 4(f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t)$$

$$2332 + 4\eta_t^2 a_t \Delta^* + 4\gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + 4\eta_t^2 R a_t \Delta^*.$$

2336 Following this, we need to establish a bound for the scalar product

$$2337 -\langle \nabla f(x_t), x_t - x_{t+1} \rangle = \theta_t \left\langle \nabla f(x_t), -\frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi_j} (x_{m,t}^{r,j}) \right\rangle.$$

2341 Using the identity $2\langle a, b \rangle = \|a + b\|^2 - \|a\|^2 - \|b\|^2$, we obtain

$$2342 -\langle \nabla f(x_t), x_t - x_{t+1} \rangle = - \left(\frac{\theta_t}{2} \|\nabla f(x_t)\|^2 + \frac{\theta_t}{2} \left\| \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi_j} (x_{m,t}^{r,j}) \right\|^2 \right)$$

$$2343 + \frac{\theta_t}{2} \left\| \nabla f(x_t) - \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi_j} (x_{m,t}^{r,j}) \right\|^2$$

$$2344 = - \left(\frac{\theta_t}{2} \|\nabla f(x_t)\|^2 + \frac{\theta_t}{2} \left\| \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \nabla f_m^{\pi_j} (x_{m,t}^{r,j}) \right\|^2 \right)$$

$$2345 + \frac{\theta_t}{2} \left\| \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} (\nabla f_m^{\pi_j} (x_{m,t}^{r,j}) - \nabla f_m^{\pi_j} (x_t)) \right\|^2.$$

2357 Using Lemma 1 and omitting one of the terms, we get

$$2358 -\langle \nabla f(x_t), x_t - x_{t+1} \rangle \leq -\frac{\theta_t}{2} \|\nabla f(x_t)\|^2 + \frac{\theta_t}{2} (\tilde{a}_t)^2 \frac{1}{R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \|x_{m,t}^{r,j} - x_t\|^2.$$

2362 Taking the expectation with respect to the randomness of the algorithm, we have

$$2363 \mathbb{E}[f(x_{t+1})] \leq f(x_t) - \frac{\theta_t}{2} \|\nabla f(x_t)\|^2$$

$$2364 + \frac{\theta_t \tilde{a}_t^2}{2R} \sum_{r=0}^{R-1} \frac{1}{C} \sum_{m \in S_t^{\lambda^r}} \frac{1}{N} \sum_{j=0}^{N-1} \|x_{m,t}^{r,j} - x_t\|^2$$

$$2365 + \frac{\hat{a}_t}{2} \|x_t - x_{t+1}\|^2.$$

2371 Recalling the definition of $\mathbb{E}[V_t]$ and taking the conditional expectation, we obtain

$$2372 \mathbb{E}[f(x_{t+1}) | x_t] \leq f(x_t) - \frac{\theta_t}{2} \|\nabla f(x_t)\|^2 + \frac{\theta_t \tilde{a}_t^2}{2} \mathbb{E}[V_t] + \theta_t^2 \hat{a}_t \|\nabla f(x_t)\|^2 + \theta_t^2 \hat{a}_t^3 \mathbb{E}[V_t]$$

$$2373 = f(x_t) - \frac{\theta_t}{2} (1 - 2\theta_t \hat{a}_t) \|\nabla f(x_t)\|^2 + \frac{\theta_t \tilde{a}_t^2}{2} \mathbb{E}[V_t] + \theta_t^2 \hat{a}_t^3 \mathbb{E}[V_t].$$

Using the fact that $\theta_t \leq \frac{1}{4\hat{a}_t}$, we arrive at

$$\mathbb{E}[f(x_{t+1}) | x_t] \leq f(x_t) - \frac{\theta_t}{4} \|\nabla f(x_t)\|^2 + \frac{\tilde{a}_t^2}{8\hat{a}_t} \mathbb{E}[V_t] + \frac{\hat{a}_t^3}{16\hat{a}_t^2} \mathbb{E}[V_t].$$

Recalling the bound on $\mathbb{E}[V_t]$, we obtain

$$\begin{aligned} \mathbb{E}[f(x_{t+1}) | x_t] &\leq f(x_t) - \frac{\theta_t}{4} \|\nabla f(x_t)\|^2 + \frac{\tilde{a}_t^2}{8\hat{a}_t} \mathbb{E}[V_t] + \frac{\hat{a}_t^3}{16\hat{a}_t^2} \mathbb{E}[V_t] \\ &\leq f(x_t) - \frac{\theta_t}{4} \|\nabla f(x_t)\|^2 \\ &\quad + \left(\frac{\tilde{a}_t^2}{2\hat{a}_t} + \frac{\hat{a}_t^3}{4\hat{a}_t^2} \right) (f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \left(\frac{\tilde{a}_t^2}{2\hat{a}_t} + \frac{\hat{a}_t^3}{4\hat{a}_t^2} \right) (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned} \quad (16)$$

Using the fact that $\theta_t \geq \frac{\zeta}{\hat{a}_t}$, we get that

$$\frac{\theta_t \|\nabla f(x_t)\|^2}{4} \geq \frac{\zeta \|\nabla f(x_t)\|^2}{4\hat{a}_t}.$$

Therefore,

$$\frac{\theta_t \|\nabla f(x_t)\|^2}{4} \geq \begin{cases} \frac{\zeta \|\nabla f(x_t)\|^2}{8L_0}, & \|\nabla f(x_t)\| \leq \frac{L_0}{L_1}, \\ \frac{\zeta \|\nabla f(x_t)\|}{8L_1}, & \|\nabla f(x_t)\| > \frac{L_0}{L_1}, \end{cases} = \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\}.$$

Denote $\delta_t \stackrel{\text{def}}{=} f(x_t) - f^*$. Then we have

$$\begin{aligned} \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} &\leq -\delta_{t+1} + \delta_t \\ &\quad + \frac{2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \delta_t \\ &\quad + \frac{2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned}$$

Recall that $\eta_t \leq \frac{2\hat{a}_t}{cR} \sqrt{\frac{1}{a_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3)}}$, $\gamma_t \leq \frac{2\hat{a}_t}{cRN} \sqrt{\frac{1}{\tilde{a}_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3)}}$, $c \geq \sqrt{T}$. Using Mishchenko et al. (2020, Lemma 6), we appear at

$$\begin{aligned} \min_{t=0,1,\dots,T-1} \left\{ \frac{\zeta}{8} \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \\ \leq \frac{\left(1 + \frac{2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \right)^T}{T} \delta_0 \\ + \frac{2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned}$$

Corollary 4. Fix $\varepsilon > 0$. Choose $c = 2\sqrt{T}$. Let $\eta_t \leq 2\hat{a}_t \sqrt{\frac{3\delta_0}{2a_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3) \Delta^* RT}}$, $\gamma_t \leq \frac{2\hat{a}_t}{N} \sqrt{\frac{3\delta_0}{\tilde{a}_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3) \Delta^* RT}}$. Then, if $T \geq \frac{72\delta_0}{\varepsilon^2}$, we have

$$\mathbb{E} \left[\min_{t=0,\dots,T-1} \left\{ \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \leq \varepsilon.$$

Proof of Corollary 4. Since $c = 2\sqrt{T}$ and $\eta_t \leq \frac{2\hat{a}_t}{cR} \sqrt{\frac{1}{a_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3)}}$, $\gamma_t \leq \frac{2\hat{a}_t}{cRN} \sqrt{\frac{1}{\tilde{a}_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3)}}$, and $\eta_t \leq 2\hat{a}_t \sqrt{\frac{3\delta_0}{2a_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3) \Delta^* RT}}$, $\gamma_t \leq \frac{2\hat{a}_t}{N} \sqrt{\frac{3\delta_0}{\tilde{a}_t(2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3) \Delta^* RT}}$ due to the choice of $T \geq$

2430 $\max \left\{ \frac{72\delta_0}{\zeta\varepsilon}, \frac{12\Delta^*}{\zeta\varepsilon}, \frac{6\bar{\Delta}^*}{\zeta\varepsilon} \right\}$, we obtain that

$$2432 \frac{\left(1 + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t)\right)^T}{T} \delta_0 \leq \frac{e\delta_0}{T} \leq \frac{3\delta_0}{T} \leq \frac{\zeta\varepsilon}{24},$$

$$2434 \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \eta_t^2 R a_t \Delta^*) \leq \frac{\zeta\varepsilon}{24},$$

2437 and that

$$2438 \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* \leq \frac{\zeta\varepsilon}{24}.$$

2440 Therefore, $\mathbb{E} \left[\min_{t=0, \dots, T-1} \left\{ \min \left\{ \frac{\|\nabla f(x_t)\|^2}{L_0}, \frac{\|\nabla f(x_t)\|}{L_1} \right\} \right\} \right] \leq \varepsilon.$ \square

2442 D.2 ASYMMETRIC GENERALIZED-SMOOTH FUNCTIONS UNDER PL-COMDITION

2444 **Theorem 8.** *Let Assumptions 1 and 2 hold for functions f , $\{f_m\}_{m=1}^M$ and $\{f_{mj}\}_{m=1, j=1}^{M, N}$. Let*
 2445 *Assumption 4 hold. Choose $0 < \zeta \leq \frac{1}{4}$. Let $\delta_0 \stackrel{\text{def}}{=} f(x_0) - f^*$. Choose any integer $T > \frac{64\delta_0 L_1^2}{\mu\zeta}$. For*
 2446 *all $0 \leq t \leq T-1$, denote*

$$2448 \hat{a}_t = L_0 + L_1 \|\nabla f(x_t)\|, \quad a_t = L_0 + L_1 \max_m \|\nabla f_m(x_t)\|, \quad \tilde{a}_t = L_0 + L_1 \max_{m,j} \|\nabla f_m^{\pi_j}(x_t)\|.$$

2450 Put $\Delta^* = f^* - \frac{1}{M} \sum_{m=1}^M f_m^*$ and $\bar{\Delta}^* = f^* - \frac{1}{M} \sum_{m=1}^M \frac{1}{N} \sum_{j=0}^{N-1} f_{mj}^*$. Impose the following
 2451 conditions on the local stepsizes γ_t , server stepsizes η_t , global stepsizes θ_t :

$$2453 \gamma_t N R \leq \eta_t R \leq \min \left\{ \frac{1}{16\hat{a}_t}, \frac{2\hat{a}_t}{c} \sqrt{\frac{1}{a_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}}, \sqrt{\frac{\hat{a}_t^2 \mu \zeta}{32L_1^2(\delta_t + \Delta^*) a_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}}, \right.$$

$$2456 \left. \sqrt{\frac{\hat{a}_t^2 \delta_0}{TL_1^2(\delta_t + \Delta^*) a_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}} \right\},$$

$$2459 \gamma_t N R \leq \min \left\{ \frac{2\hat{a}_t}{c} \sqrt{\frac{1}{\tilde{a}_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}}, \sqrt{\frac{\hat{a}_t^2 \delta_0}{TL_1^2(\delta_t + \bar{\Delta}^*) \tilde{a}_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}}, \right.$$

$$2463 \left. \sqrt{\frac{\hat{a}_t^2 \mu \zeta}{32L_1^2(\delta_t + \bar{\Delta}^*) \tilde{a}_t(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}} \right\}.$$

$$2467 \frac{\zeta}{\hat{a}_t} \leq \theta_t \leq \frac{1}{4\hat{a}_t}, \quad 0 \leq t \leq T-1,$$

2469 where $c \geq \sqrt{T}$. Let \tilde{T} be an integer such that $0 \leq \tilde{T} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$ be a constant, $\alpha \leq \sqrt{\frac{\delta_0}{AT}}$.

2470 Then, the iterates $\{x_t\}_{t=0}^{T-1}$ of Algorithm 3 satisfy

$$2472 \delta_T \leq \left(1 - \frac{\mu\zeta}{4L_0}\right)^{T-\tilde{T}} \delta_0 + \frac{4L_0 A \alpha^2}{\mu\zeta},$$

2475 where $\delta_T \stackrel{\text{def}}{=} f(x_T) - f^*$.

2476 *Proof of Theorem 8.* Let us follow the first steps of the proof of Theorem 4. Consider (16):

$$2478 \frac{\theta_t}{4} \|\nabla f(x_t)\|^2 \leq f(x_t) - f(x_{t+1})$$

$$2480 + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t)$$

$$2482 + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*).$$

Since $\theta_t \geq \frac{\zeta}{\hat{a}_t}$, and f satisfies Polyak–Łojasiewicz Assumption 4, we obtain that

$$\begin{aligned} \frac{\mu\zeta(f(x_t) - f^*)}{2\hat{a}_t} &\leq f(x_t) - f(x_{t+1}) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned}$$

1. Let \tilde{T} be the number of steps t , so that $\|\nabla f(\hat{x}_t)\| \geq \frac{L_0}{L_1}$. For such t , we have $L_0 + L_1 \|\nabla f(x_t)\| = \hat{a}_t \leq 2L_1 \|\nabla f(x_t)\|$. Therefore, we get

$$\begin{aligned} \frac{\mu\zeta(f(x_t) - f^*)}{4L_1 \|\nabla f(x_t)\|} &\leq f(x_t) - f(x_{t+1}) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned}$$

Notice that the relation $\hat{a}_t \leq 2L_1 \|\nabla f(x_t)\|$ and Lemma 1 together imply

$$\frac{\|\nabla f(x_t)\|}{4L_1} \leq \frac{\|\nabla f(x_t)\|^2}{2\hat{a}_t} \leq f(x_t) - f^*.$$

Hence, we have

$$\begin{aligned} \frac{\mu\zeta}{16L_1^2} &\leq f(x_t) - f(x_{t+1}) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned}$$

Subtracting f^* on both sides and introducing $\delta_t \stackrel{\text{def}}{=} f(x_t) - f^*$, we obtain

$$\begin{aligned} \delta_{t+1} &\leq \delta_t - \frac{\mu\zeta}{16L_1^2} \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \delta_t (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*). \end{aligned}$$

As $\gamma_t \leq \sqrt{\frac{4\hat{a}_t^2 \mu\zeta}{128L_1^2(\delta_t + \bar{\Delta}^*)_{a_t} R^2 N^2 (2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3)}}$ and $\eta_t \leq \sqrt{\frac{4\hat{a}_t^2 \mu\zeta}{128L_1^2(\delta_t + \bar{\Delta}^*)_{a_t} R^2 (2\hat{a}_t \tilde{a}_t^2 + \hat{a}_t^3)}}$, it follows that

$$\begin{aligned} \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \delta_t (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) + \\ + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^*) \leq \frac{\mu\zeta}{32L_1^2}. \end{aligned}$$

Therefore, we get

$$\delta_{t+1} \leq \delta_t - \frac{\mu\zeta}{32L_1^2}.$$

2. Suppose now that $\|\nabla f(x_t)\| \leq \frac{L_0}{L_1}$. For such t , we have $L_0 + L_1 \|\nabla f(x_t)\| = \hat{a}_p \leq 2L_0$. Hence,

$$\begin{aligned} \frac{\mu\zeta(f(x_t) - f^*)}{4L_0} &\leq f(x_t) - f(x_{t+1}) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} (f(x_t) - f^*) (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \left(\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^* \right). \end{aligned}$$

Subtracting f^* on both sides and introducing $\delta_t \stackrel{\text{def}}{=} f(x_t) - f^*$, we obtain

$$\begin{aligned} \delta_{t+1} &\leq \delta_t \rho + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \delta_t (\eta_t^2 a_t + \eta_t^2 R^2 \hat{a}_t + \gamma_t^2 N \tilde{a}_t + \eta_t^2 R a_t) \\ &\quad + \frac{2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3}{4\hat{a}_t^2} \left(\eta_t^2 a_t \Delta^* + \gamma_t^2 N \tilde{a}_t \bar{\Delta}^* + \eta_t^2 R a_t \Delta^* \right). \end{aligned}$$

where $\rho \stackrel{\text{def}}{=} 1 - \frac{\mu\zeta}{4L_0}$. Let $\gamma_t \stackrel{\text{def}}{=} \alpha \hat{\gamma}_t$ and $\eta_t \stackrel{\text{def}}{=} \alpha \hat{\eta}_t$ with $\hat{\gamma}_t \leq \sqrt{\frac{4\hat{a}_t^2 A}{4L_1^2(\delta_t + \Delta^*)\tilde{a}_t R^2 N^2(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}}$ and

$\hat{\eta}_t \leq \sqrt{\frac{4\hat{a}_t^2 A}{4L_1^2(\delta_t + \Delta^*)a_t R^2(2\hat{a}_t\tilde{a}_t^2 + \hat{a}_t^3)}}$, for some constant $A > 0$. Then,

$$\delta_{t+1} \leq \rho \delta_t + A\alpha^2.$$

Unrolling the recursion, we derive

$$\begin{aligned} \delta_T &\leq \rho^{T-\tilde{T}} \delta_0 + A\alpha^2 \sum_{i=0}^{\infty} \rho^i - \frac{\mu\zeta}{32L_1^2} \sum_{i=0}^{N-1} \rho^i \\ &\leq \rho^{P-\tilde{P}} \delta_0 + \frac{A\alpha^2}{1-\rho} - \frac{1-\rho^{\tilde{P}}}{1-\rho} \frac{\mu\zeta}{32L_1^2}. \end{aligned}$$

Notice that $\delta_{t+1} \leq \delta_t + A\alpha^2$, which implies

$$\delta_T \leq \delta_0 + (T - \tilde{T}) A\alpha^2 - \tilde{T} \frac{\mu\zeta}{32L_1^2}.$$

Since $\alpha \leq \sqrt{\frac{\delta_0}{AT}}$, we conclude that

$$0 \leq \delta_T \leq 2\delta_0 - \tilde{T} \frac{\mu\zeta}{32L_1^2}, \quad \Rightarrow \tilde{T} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}.$$

Therefore, for $T > \frac{64\delta_0 L_1^2}{\mu\zeta}$ we can guarantee that $T - \tilde{T} > 0$ and

$$\begin{aligned} \delta_T &\leq \rho^{T-\tilde{T}} \delta_0 + \frac{A\alpha^2}{1-\rho} - \tilde{T} \rho^{\tilde{T}} \frac{\mu\zeta}{32L_1^2} \\ &\leq \rho^{T-\tilde{T}} \delta_0 + \frac{A\alpha^2}{1-\rho}. \end{aligned}$$

□

Corollary 7. Fix $\varepsilon > 0$. Choose $\alpha \leq \min \left\{ \sqrt{\frac{\delta_0}{AT}}, L_1 \sqrt{\frac{8\delta_0 \varepsilon}{L_0 AT}} \right\}$. Then, if $T \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we have $\delta_T \leq \varepsilon$.

Proof of Corollary 7. Since $0 \leq \tilde{T} \leq \frac{64\delta_0 L_1^2}{\mu\zeta}$, $A > 0$, $\alpha \leq \sqrt{\frac{\delta_0}{AT}}$, $\alpha \leq L_1 \sqrt{\frac{8\delta_0 \varepsilon}{L_0 AT}}$, due to the choice of $T \geq \frac{64\delta_0 L_1^2}{\mu\zeta} + \frac{4L_0}{\mu\zeta} \ln \frac{2\delta_0}{\varepsilon}$, we obtain that

$$\left(1 - \frac{\mu\zeta}{4L_0}\right)^{T-\tilde{T}} \delta_0 \leq e^{-\frac{\mu\zeta}{4L_0}(T-\tilde{T})} \delta_0 \leq \frac{\varepsilon}{2},$$

and that

$$\frac{4L_0 A}{\mu\zeta} \cdot \frac{\delta_0}{AT} \leq \frac{\varepsilon}{2}.$$

Therefore, $\delta_T \leq \varepsilon$. □

E ADDITIONAL EXPERIMENTAL DETAILS FOR MAIN PART

In this section, we provide additional experimental details: parameters search grids and some technical details that did not fit in the main text. For all the plots we provide in the legend all the best parameters found by the grid search. The parameter grids are provided as table for every method. All the code can be seen at https://anonymous.4open.science/r/local_steps_rr-BA8E/.

It can be seen from pseudocode of Algorithms 1, 2, 3, that global stepsize depends on the full gradient. However, our numerical tests showed that use of gradient approximations g_p for Algorithm 1 and g_t for Algorithms 2, 3 gives better numerical results while being less computationally expensive. Thus, in our practical experiments we decided to use this approximation in calculation of global stepsize. We want to point out, that the theoretical analysis for this “practical” version of the algorithm can be done by considering very small inner stepsizes. Although, we decided not to include it in the current version to keep the presentation more concise and avoid additional complexities.

E.1 METHODS WITH RANDOM RESHUFFLING

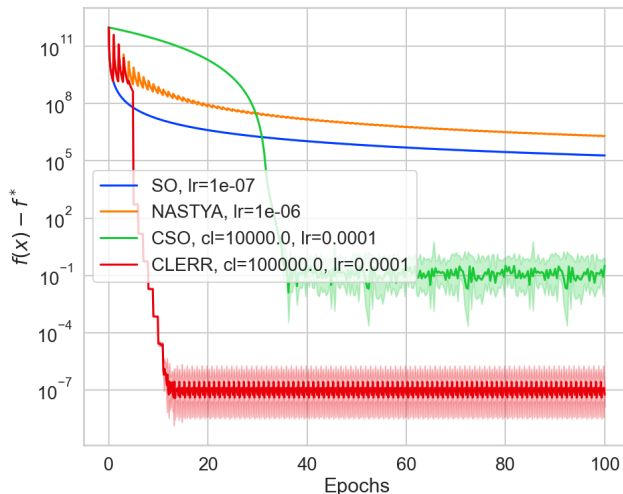


Figure 5: Function residual for (4), $\alpha_t = 10^{-7}$. The best parameters are provided in the legend.

In these experiments we compare methods with random reshuffling, that shuffle data once at the start of training process. The main idea is to show the positive impact of random reshuffling and clipping on algorithm performance. We incorporate these two techniques inside our CLERR method (Algorithm 2).

Firstly, consider (4). For these experiments we take $d = 1$ and randomly sample 1000 shifts $x_i \in [-10, 10]$. We run all the methods for 10 different seeds on a logarithmic hyperparameter grid. Then we choose the best hyperparameters according to the best mean loss values on the second half of epochs. The parameter grid is provided in Table 1. To find f^* , we run the Newton method for couple iterations until convergence.

Since both Nastya and Algorithm 2 have jumping at the end of every epoch, if we tuned the inner stepsize along with other parameters, the inner stepsize would go to zero and the outer stepsize would be selected such as these methods solve the problem in 1 step. This would be unfair because other baselines do not use a jumping technique, so they would not be able to achieve such performance. Thus, we decided to fix the inner stepsize for Algorithm 2 and Nastya equal to the best stepsize, chosen for SO, and tune the clipping level and outer stepsize with the outer stepsize not exceeding the values supported in theory. Here and later, for simplicity, we speak about Algorithm 2 in terms of stepsize and clipping level, that we can obtain from c_0 and c_1 from (3). The best stepsize for SO is 10^{-7} , so we choose inner stepsize for Nastya and Algorithm 2 the same. Nastya chooses outer stepsize equal 10^{-7} , while CSO and CLERR (Algorithm 2) choose it equal to 10^{-4} . CSO clips gradients at the level 10^4 , while CLERR – at the level 10^5 .

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

Method	Stepsize	Clipping Level	Inner Stepsize
SO	$[10^{-8}, 10^{-2}]$	-	-
NASTYA	$[10^{-8}, 10^{-2}]$	-	10^{-7}
CSO	$[10^{-8}, 10^{-2}]$	$[10^0, 10^5]$	-
Algorithm 2	$[10^{-8}, 10^{-2}]$	$[10^0, 10^5]$	10^{-7}

Table 1: Parameter grids for experiments on methods with random reshuffling on (4).

E.1.1 RESNET-18 ON CIFAR-10

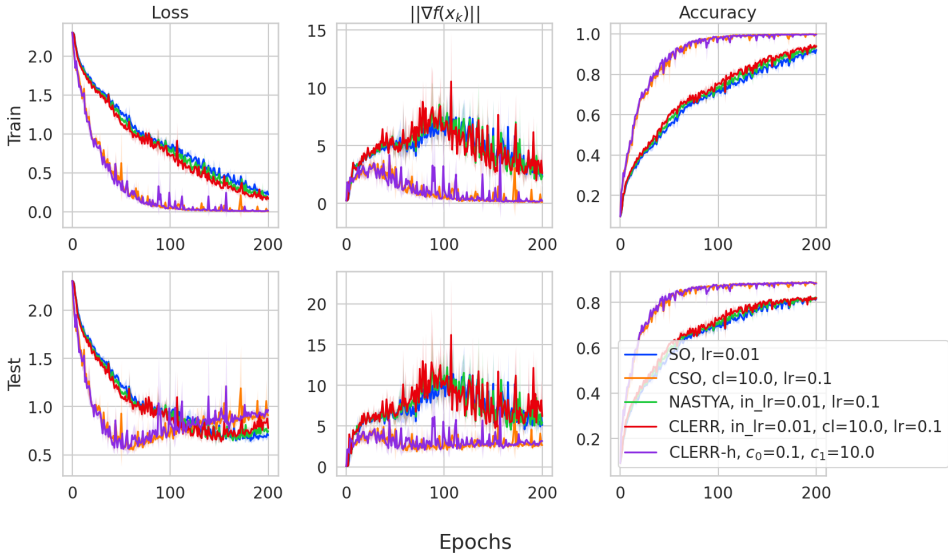


Figure 6: Loss, gradient norm and accuracy on train and test dataset for ResNet-18 on CIFAR-10. The best parameters are provided in the legend.

In this experiment we consider image classification task. We train ResNet-18 He et al. (2016) on the CIFAR-10 Krizhevsky et al. (2009) dataset. The implementation of ResNet-18 was taken from <https://github.com/kuangliu/pytorch-cifar>. All the methods are run on 3 different random seeds on logarithmic hyperparameter grid. Then we choose the best hyperparameters according to the best mean test accuracy on the last 25% of epochs.

In this experiment, we do not fix the inner stepsize for Nastya and CLERR, since methods do not try to make it as small as possible, as it was in the previous experiment. However, both SO, Nastya, and CLERR choose the same inner stepsize 10^{-2} as the best. Then, both Nastya and CLERR choose bigger outer step size 10^{-1} , and CLERR also chooses clipping level on outer step size as 10. Despite the fact that both Nastya and CLERR choose bigger outer stepsizes compared to inner stepsize, jumping does not have any impact on this problem. CLERR clips outer gradients at the level of 10, so this also does not help method to converge to a better area.

Moreover, we provide results of heuristically modified Algorithm 2, where we fix clipping level and inner stepsize of Algorithm 2 equal to the best clipping level and the best stepsize from CSO correspondingly. The tunable parameters are only c_0 and c_1 for outer stepsize. We call this method CLERR-h. CLERR-h chooses an outer stepsize equal to 5, while the clipping level is very tiny and equal to 10^{-2} . All the parameter grids are provided in Table 2.

Method	Stepsize	Clipping Level	Inner Stepsize	c_0	c_1
RR	$[10^{-3}, 10^{-1}]$	-	-	-	-
NASTYA	$[10^{-3}, 10^{-1}]$	-	$[10^{-4}, 10^0]$	-	-
CRR	$[10^{-3}, 10^{-1}]$	$[10^0, 10^3]$	-	-	-
CLERR	$[10^{-3}, 10^{-1}]$	$[10^0, 10^3]$	$[10^{-4}, 10^0]$	-	-
CLERR-h	-	10^1	10^{-1}	$[10^{-2}, 10^1]$	$[10^{-2}, 10^1]$

Table 2: Parameter grids for experiments on methods with random reshuffling on ResNet-18 on CIFAR-10.

E.2 METHODS WITH LOCAL STEPS

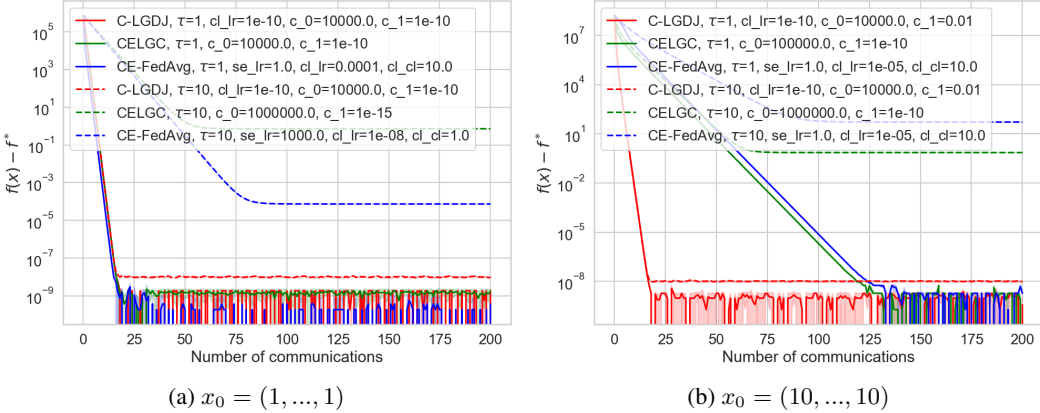


Figure 7: Function residual for (4), starting from different x_0 for different number of local steps on the client device τ . The best parameters are provided in the legend.

In these experiments we compare methods with local steps: Algorithm 1 (C-LGDJ) with Communication Efficient Local Gradient Clipping (CELGC) (Liu et al., 2022) and Clipping-Enabled-FedAvg (CE-FedAvg) Zhang et al. (2022). For comparison we take problem (4) for $d = 100$, where we randomly sample 1000 shifts $x_i \in [-10, 10]^d$. To make the distributions of data on each client more distinct between each other, we sort the whole dataset at the beginning of the experiment by $\|x_i\|$. Each method has 10 clients, where each client has equal number of data. We provide results for two starting points: $x_0 = (1, \dots, 1)$ and $x_0 = (10, \dots, 10)$. All the methods are run for 10 different random seeds on logarithmic hyperparameter grid. The best hyperparameters are chosen according to the best mean loss on the last 25% of epochs.

Each client performs $\tau = 1$ or $\tau = 10$ local steps, and each local step is performed on the whole local data. For ease of implementation and due to computational limitations we iterate over all the clients sequentially.

We reformulate constants c_0 and c_1 as server stepsize and clipping level from (3) to better interpret the experimental results. We start by paying attention to results with a single local step. Firstly, consider C-LGDJ (Algorithm 1). It chooses tiny client stepsizes 10^{-10} and small server stepsizes $5 \cdot 10^{-5}$ for both starting points. For Figure 7a it also takes very big clipping level for server 10^{14} , compared to Figure 7b, where it clips on level 10^6 , which is obvious because on the second picture methods start farther from the minimum and have bigger gradients. Secondly, consider CELGC. In both cases, it takes very small client stepsizes: $5 \cdot 10^{-5}$ and $5 \cdot 10^{-6}$ respectively, and very big clipping levels: 10^{14} and 10^{15} respectively. Finally, CE-FedAvg also takes small client stepsizes: 10^{-4} and 10^{-5} , rather big server stepsizes, which are equal to 1, and average client clipping levels: 10 in both cases. For $\tau = 10$ we have the same parameters for C-LGDJ, CELGC tries to make even smaller steps with high clipping levels, while CE-FedAvg uses a much bigger server stepsize and much smaller client stepsize, for the case from Figure 7a.

The grids of hyperparameters for $x_0 = (1, \dots, 1)$ are provided in Table 3, and for $x_0 = (10, \dots, 10)$ – in Table 4.

Method	Cl. Stepsize	Se. Stepsize	Cl. Clip Level	c_0	c_1
Clipped-L-SGD-J	$[10^{-10}, 10^0]$	-	-	$[10^{-10}, 10^6]$	$[10^{-10}, 10^6]$
CELGC	-	-	-	$[10^{-15}, 10^{10}]$	$[10^{-15}, 10^{10}]$
CE-FedAvg	$[10^{-10}, 10^0]$	$[10^{-10}, 10^3]$	$[10^0, 10^4]$	-	-

Table 3: Parameter grids for experiments on methods with local steps on (4) for $x_0 = (1, \dots, 1)$. Here "cl." means "Client", and "se." – "server".

Method	Cl. Stepsize	Se. Stepsize	Cl. Clip Level	c_0	c_1
Clipped-L-SGD-J	$[10^{-10}, 10^0]$	-	-	$[10^{-10}, 10^6]$	$[10^{-10}, 10^6]$
CELGC	-	-	-	$[10^{-10}, 10^{10}]$	$[10^{-10}, 10^{10}]$
CE-FedAvg	$[10^{-10}, 10^0]$	$[10^{-10}, 10^0]$	$[10^0, 10^4]$	-	-

Table 4: Parameter grids for experiments on methods with local steps on (4) for $x_0 = (10, \dots, 10)$. Here "cl." means "client", and "se." – "server".

E.3 METHODS WITH LOCAL STEPS, RANDOM RESHUFFLING AND PARTIAL PARTICIPATION

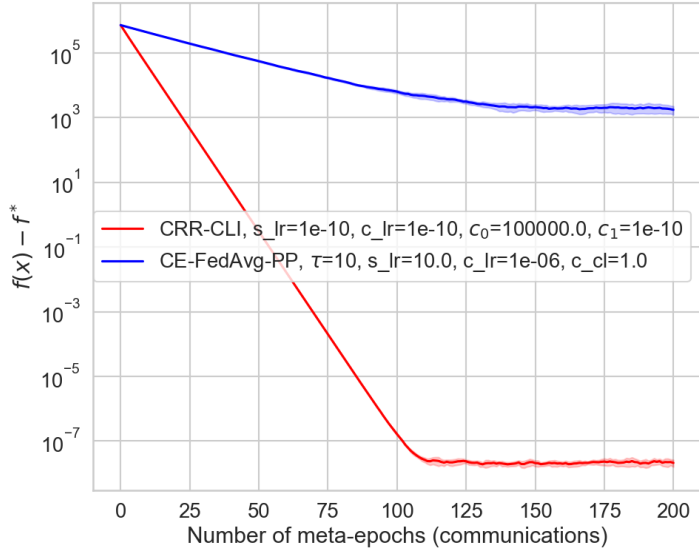


Figure 8: Function residual for (4), starting from $x_0 = (1, \dots, 1)$ with batch size 16. The best parameters are provided in the legend.

In these experiments we compare methods with clipping, random reshuffling, local steps and partial participation: Algorithm 3 (CRR-CLI) and with CE-FedAvg Zhang et al. (2022) with partial participation (CE-FedAvg-PP). For comparison we take problem (4) for $d = 100$, where we randomly sample 1000 shifts $x_i \in [-10, 10]^d$. Again, to make the distributions of data on each client more distinct between each other, we sort the whole dataset at the beginning of the experiment by $\|x_i\|$. All the methods are run for 10 different random seeds on logarithmic hyperparameter grid. The best hyperparameters are chosen according to the best mean loss on the last 25% of epochs.

Each method has 10 clients, where each client has the same amount of data. The size of the cohort is chosen to be 2. The method performs local steps on each client from the cohort, after which it performs communication and goes to the next cohort. In the Algorithm 3 the clients to the cohort are chosen sequentially with sliding window after Client-Reshuffling. In CE-FedAvg-PP clients to the cohort are always chosen randomly. The starting point is chosen $x_0 = (1, \dots, 1)$. All the methods are run for 10 different random seeds. The best hyperparameters are chosen according to the best mean loss on the last 25% of epochs.

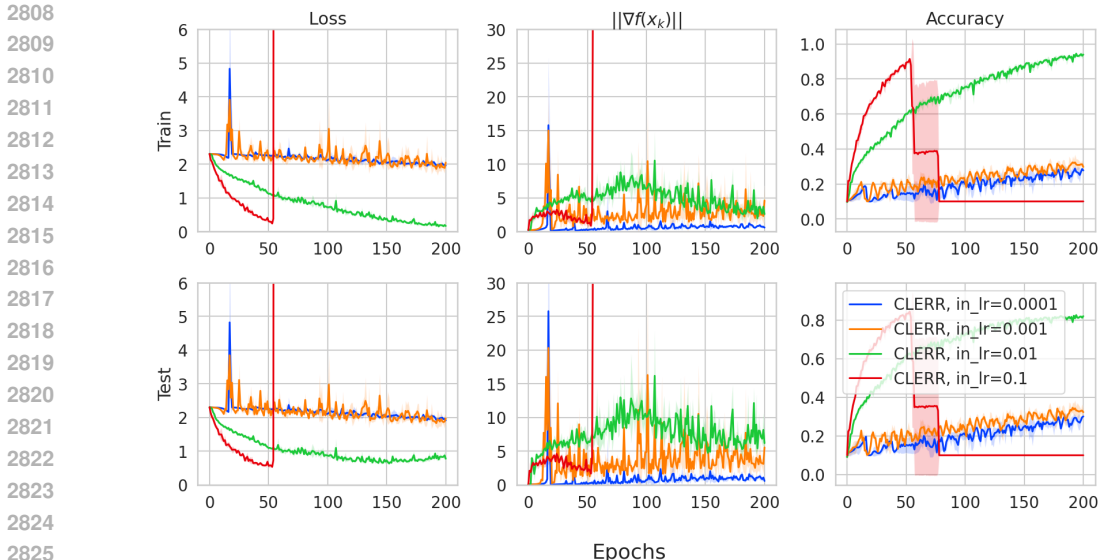


Figure 9: Algorithm 2 with different step sizes on ResNet-18 on CIFAR-10.

For local steps we chose batch size equal to 16. In Algorithm 3 every client goes sequentially over the whole shuffled local dataset with batch size window. In CE-FedAvg-PP we fix number of local steps to 10, and each client samples batch on every local step.

Just like in previous experiment in Section 5.2, all the methods try to reduce the influence of local steps by making inner stepsizes very small. Algorithm 3 chooses both client and server stepsizes equal 10^{-10} , and CE-FedAvg-PP chooses client stepsize equal 10^{-6} and client clipping level equals 1. Speaking of outer steps, Algorithm 3 chooses global stepsize equal to $5 \cdot 10^{-7}$ with clipping level 10^{16} . And CE-FedAvg-PP has server stepsize equal to 10. The grids of hyperparameters are provided in Table 5.

Method	Cl. Step size	Se. Step size	Cl. Clip Level	c_0	c_1
CRR-CLI	$[10^{-10}, 10^6]$	$[10^{-10}, 10^6]$	-	$[10^{-10}, 10^5]$	$[10^{-10}, 10^5]$
CE-FedAvg-PP	$[10^{-10}, 10^3]$	$[10^{-10}, 10^3]$	$[10^0, 10^4]$	-	-

Table 5: Parameter grids for experiments on methods with clipping, random reshuffling, local steps and partial participation. Here "cl." means "client", and "se." – "server".

F ADDITIONAL EXPERIMENTS

In this section, we provide additional numerical experiments, that did not fit in the main paper: in Section F.1 we investigate the influence of that inner step size on the behavior of Algorithm 2, and in Section F.2 we provide additional experiments on logistic regression, where we compare Algorithm 2 with clipped SGD.

F.1 HOW THE INNER STEP SIZE AFFECTS CONVERGENCE OF THE METHOD

In this experiment, we investigate the influence of the inner step size on the behavior of Algorithm 2 on ResNet-18 on CIFAR-10. To do this, we take the same hyperparameters for Algorithm 2 as in Sections 5.1.1, E.1.1 and only change the inner step size. The results are provided in Figure 9.

On the one hand, if we take the inner step size too small (blue and orange lines), it converges very slowly. This is obvious since Algorithm 2 becomes regular Clipped-GD, which can be seen from pseudocode. Because Clipped-GD performs a single step per epoch, it has slow convergence. On the other hand, if we take the inner step size too big (red line), the method diverges. It does not have clipping on the inner step, so such behavior is expected. To summarize, it is important to take the inner step size small, but not too small, because it may slow down the convergence.

2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915

F.2 LOGISTIC REGRESSION EXPERIMENTS

Since in the experiments on neural networks (Sections 5.1.1, E.1.1) regular CSO (SGD with clipping) showed very good results, we decided to conduct additional experiments on logistic regression, where we compare CSO with our Algorithm 2. We consider gisette and realsim datasets from libsvm library Chang & Lin (2011). All the methods are run for 3 different random seeds on logarithmic hyperparameter grid. The best hyperparameters are chosen according to the best mean loss on the last 25% of epochs. The results are presented in Figure 9.

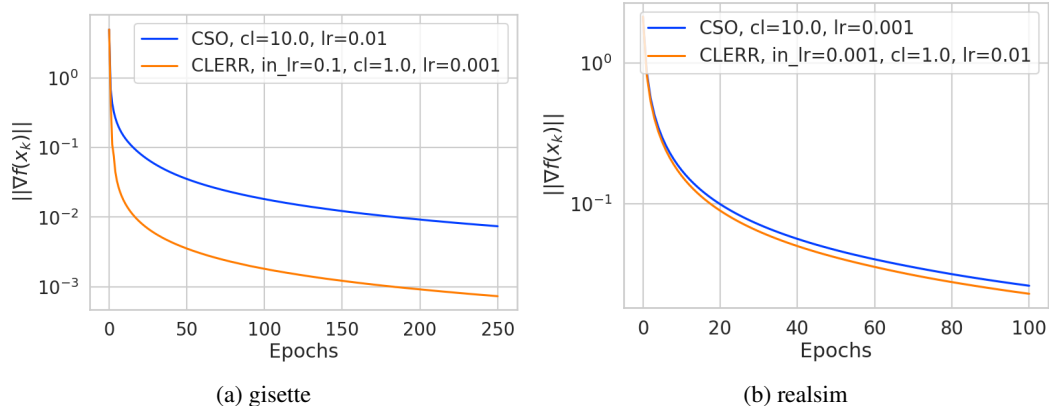


Figure 10: Gradient norm for logistic regression problem on gisette and realsim datasets. The best parameters are provided in the legend.

Since the inner stepsize for CLERR has the same meaning as stepsize for CSO, we decided to take the same parameter grids for these two parameters. The same goes for clipping levels in spite of the fact that CLERR clips the gradient approximation only in the end of the epoch. This experiment shows that CLERR either has the same performance as CSO or better. Since logistic regression is (L_0, L_1) -smooth, such result is expected, as Algorithm 2 is designed for such type of functions. Figure 10a shows us that CLERR chooses very small outer stepsize 10^{-3} , while inner step size is bigger than the one in CSO: 10^{-1} vs 10^{-2} . In the Figure 10b CLERR chooses parameters in the opposite way: inner step size is very small and equal to the one from CSO, while the outer stepsize is bigger. The parameter grids for gisette dataset is presented in Table 6, and for realsim – in Table 7.

2916
 2917
 2918
 2919
 2920
 2921
 2922
 2923
 2924
 2925
 2926
 2927
 2928
 2929
 2930
 2931
 2932
 2933
 2934
 2935
 2936
 2937
 2938
 2939
 2940
 2941
 2942
 2943
 2944
 2945
 2946
 2947
 2948
 2949
 2950
 2951
 2952
 2953
 2954
 2955
 2956
 2957
 2958
 2959
 2960
 2961
 2962
 2963
 2964
 2965
 2966
 2967
 2968
 2969

	Stepsize	Clipping Level	Inner Stepsize
CSO	$[10^{-3}, 10^{-1}]$	$[10^0, 10^2]$	-
CLERR	$[10^{-3}, 10^{-1}]$	$[10^0, 10^2]$	$[10^{-3}, 10^{-1}]$

Table 6: Parameter grids for logistic regression experiments on gisette dataset

	Stepsize	Clipping Level	Inner Stepsize
CSO	$[10^{-5}, 10^{-1}]$	$[10^0, 10^2]$	-
CLERR	$[10^{-3}, 10^{-1}]$	$[10^0, 10^2]$	$[10^{-5}, 10^{-1}]$

Table 7: Parameter grids for logistic regression experiments on realsim dataset