

---

# Anchor Sampling for Federated Learning with Partial Client Participation

---

Feijie Wu<sup>1</sup> Song Guo<sup>2</sup> Zhihao Qu<sup>3</sup> Shiqi He<sup>4</sup> Ziming Liu<sup>2</sup> Jing Gao<sup>1</sup>

## Abstract

Compared with full client participation, partial client participation is a more practical scenario in federated learning, but it may amplify some challenges in federated learning, such as data heterogeneity. The lack of inactive clients' updates in partial client participation makes it more likely for the model aggregation to deviate from the aggregation based on full client participation. Training with large batches on individual clients is proposed to address data heterogeneity in general, but their effectiveness under partial client participation is not clear. Motivated by these challenges, we propose to develop a novel federated learning framework, referred to as FedAMD, for partial client participation. The core idea is anchor sampling, which separates partial participants into anchor and miner groups. Each client in the anchor group aims at the local bullseye with the gradient computation using a large batch. Guided by the bullseyes, clients in the miner group steer multiple near-optimal local updates using small batches and update the global model. By integrating the results of the two groups, FedAMD is able to accelerate the training process and improve the model performance. Measured by  $\epsilon$ -approximation and compared to the state-of-the-art methods, FedAMD achieves the convergence by up to  $O(1/\epsilon)$  fewer communication rounds under non-convex objectives. Empirical studies on real-world datasets validate the effectiveness of FedAMD and demonstrate the superiority of the proposed algorithm: Not only does it considerably save computation and communication costs, but also the test accuracy significantly improves.

---

<sup>1</sup>Purdue University, West Lafayette, IN, USA <sup>2</sup>The Hong Kong Polytechnic University, Hong Kong <sup>3</sup>Hohai University, Nanjing, China <sup>4</sup>The University of British Columbia, Vancouver, Canada. Correspondence to: Feijie Wu <wu1977@purdue.edu>, Song Guo <song.guo@polyu.edu.hk>, Zhihao Qu <qzhihao@hhu.edu.cn>, Jing Gao <jinggao@purdue.edu>.

## 1. Introduction

Federated learning (FL) (Konečný et al., 2015; 2016; McMahan et al., 2017) has attained an increasing interest over the past few years. As a distributed training paradigm, it enables a group of clients to collaboratively train a global model from decentralized data under the orchestration of a central server. By this means, sensitive privacy is basically protected because the raw data are not shared across the clients. Due to the unreliable network connection and the rapid proliferation of FL clients, it is infeasible to require all clients to be simultaneously involved in the training. To address the issue, recent works (Li et al., 2019; Philippenko & Dieuleveut, 2020; Gorbunov et al., 2021; Karimireddy et al., 2020b; Yang et al., 2020; Li et al., 2020; Eichner et al., 2019; Yan et al., 2020; Ruan et al., 2021; Gu et al., 2021; Lai et al., 2021) introduce a practical setting where merely a portion of clients participates in the training. The partial-client scenario effectively avoids the network congestion at the FL server and significantly shortens the idle time as compared to traditional large-scale machine learning (Zinkevich et al., 2010; Bottou, 2010; Dean et al., 2012; Bottou et al., 2018).

However, the performance of a model trained with partial client participation is much worse than the one trained with full client participation (Yang et al., 2020). A well-known challenge in federated learning is data heterogeneity, i.e., different clients have non-i.i.d. data. When only a portion of the clients participate in the learning, there are inevitably updates missing from inactive clients. With data heterogeneity and partial client participation, the optimal model is subject to the local data distribution, and therefore, the local updates on the clients' models greatly deviate from the update towards optimal global parameters (Karimireddy et al., 2020b; Malinovskiy et al., 2020; Pathak & Wainwright, 2020; Wang et al., 2020; 2021; Mitra et al., 2021; Rothchild et al., 2020; Zhao et al., 2018; Wu et al., 2021). FedAvg (McMahan et al., 2017; Li et al., 2019; Yu et al., 2019a;b; Stich, 2018), for example, is less likely to follow a correct update towards the global minimizer because the model aggregation on the active clients critically deviates from the aggregation on the full clients, an expected direction towards global minimizer (Yang et al., 2020).

There are efforts toward addressing the data heterogeneity challenge in the general federated learning setting. One promising solution is to train the model using large batches.

Convexity	Method	Partial Clients	Communication Rounds
Non-convex	Minibatch SGD (Wang & Srebro, 2019)	✗	$\frac{1}{MK\epsilon^2} + \frac{1}{\epsilon}$
	FedAvg (Yang et al., 2020)	✓	$\frac{K}{A\epsilon^2} + \frac{1}{\epsilon}$
	SCAFFOLD (Karimireddy et al., 2020b)	✓	$\frac{\sigma^2}{AK\epsilon^2} + \left(\frac{M}{A}\right)^{2/3} \frac{1}{\epsilon}$
	BVR-L-SGD (Murata & Suzuki, 2021)	✗	$\frac{1}{MK\epsilon^{3/2}} + \frac{1}{\epsilon}$
	VR-MARINA (Gorunov et al., 2021)	✗	$\frac{\sigma}{M\epsilon^{3/2}} + \frac{\sigma^2}{M\epsilon} + \frac{1}{\epsilon}$
	FedAMD (Sequential) (Corollary 4.7)	✓	$\frac{M}{A\epsilon}$
	FedAMD (Constant) (Corollary 4.9)	✓	$\frac{M}{A\epsilon}$
PL condition (or *strongly-convex)	Minibatch SGD* (Woodworth et al., 2020b)	✗	$\frac{\sigma^2}{\mu MK\epsilon} + \frac{1}{\mu} \log \frac{1}{\mu\epsilon}$
	FedAvg (Karimireddy et al., 2020a)	✓	$\frac{1+\sigma^2/K}{\mu A\epsilon} + \sqrt{\frac{1+\sigma^2/K}{\mu\sqrt{\epsilon}}} + \frac{1}{\mu} \log \frac{1}{\mu\epsilon}$
	SCAFFOLD* (Karimireddy et al., 2020b)	✓	$\frac{\sigma^2}{\mu AK\epsilon} + \left(\frac{M}{A} + \frac{1}{\mu}\right) \log \frac{M\mu}{A\epsilon}$
	VR-MARINA (Gorunov et al., 2021)	✗	$\left(\frac{\sigma^2}{\mu M\epsilon} + \frac{\sigma}{\mu^{3/2} M\sqrt{\epsilon}} + \frac{1}{\mu}\right) \log \frac{1}{\mu\epsilon}$
	FedAMD (Constant) (Corollary 4.12)	✓	$\left(\frac{1}{\mu} + \frac{M}{\mu^2 A} + \frac{M}{A}\right) \log \frac{1}{\mu\epsilon}$

Table 1. Number of communication rounds that achieve  $\mathbb{E} \|\nabla F(\tilde{\mathbf{x}}_{out})\|_2^2 \leq \epsilon$  for non-convex objectives (or  $\mathbb{E} F(\tilde{\mathbf{x}}_{out}) - F_* \leq \epsilon$  for PL condition or strongly-convex with the parameter of  $\mu$ ). We optimize an online scenario and set the small batch size to 1. The symbol ✓ or ✗ for "Partial Clients" is determined by whether full-client participation is required.

Specifically, each client fully utilizes the local training set or generates a large number of training samples, and thus could compute a gradient for an accurate estimation of the expected local update. Under full client participation, this approach follows the optimal update towards the global minimizer. Recent studies (Gorunov et al., 2021; Murata & Suzuki, 2021; Tyurin & Richtárik, 2022; Zhao et al., 2021) have demonstrated its great potential from a theoretical perspective. Measured by  $\epsilon$ -approximation, MARINA (Gorunov et al., 2021), for instance, realizes  $O(1/M\epsilon^{1/2})$  faster while using large batches, where  $M$  indicates the number of clients.

Despite the success, the large-batch update method has some drawbacks. Typically, a large batch update involves several gradient computations compared to a small batch update. This increases the burden of FL clients, especially on IoT devices like smartphones. On small devices, the hardware hardly accommodates all samples in a large batch simultaneously. Instead, the large batch has to be divided into several small batches to obtain the final gradient.

This large-batch approach is further challenged in the partial client participation scenario studied in this paper. Especially, the convergence property has not been analyzed. Some existing work, such as BVR-L-SGD (Murata & Suzuki, 2021) and FedPAGE (Zhao et al., 2021), claim that they can work under partial client participation, but they still require all clients' participation when the algorithms come to the synchronization using a large batch. In this paper, partial client participation refers to the case where only a portion of clients take part at every round during the entire training.

On the other hand, small-batch methods, such as mini-batch SGD and local SGD, are computational-friendly to perform multiple local updates, in contrast to the large batches when doing so. However, they suffer from data heterogeneity and partial client participation.

Motivated by the aforementioned observations, we propose a novel framework named FedAMD under federated learning with partial client participation. As far as we know, this is the first work that creatively integrates the complementary advantages of small-batch and large-batch approaches. The framework is based on anchor sampling, which separates the partial participants into two disjoint groups, namely, anchor and miner groups. In the anchor group, clients (a.k.a. anchors) compute the gradient using a large batch cached in the server to estimate the global orientation. In the miner group, clients (a.k.a. miners) perform multiple updates corrected according to the previous and the current local parameters and the last local update volume. Compared with using large-batch updates only, the benefit of involving miners is twofold. First, multiple local updates without severe deviation can effectively accelerate the training process. Second, FedAMD updates the global model using the local models from the miner group only. Compared with using small-batch updates only, we demonstrate a much better convergence property of the proposed method. Since anchor sampling could distinguish the clients with time-varying probability, we separately consider constant and sequential probability settings.

**Contributions.** We summarize our contributions below:

- **Algorithmically**, we propose a unified federated learn-

ing framework FedAMD based on anchor sampling, which identifies a participant as an anchor or a miner. Clients in the anchor group aim to obtain the bullseyes of their local data with a large batch, while the miners target to accelerate the training with multiple local updates using small batches.

- **Theoretically**, we establish the convergence rate for FedAMD under non-convex objectives in both constant and sequential probability settings. Incorporating the large batches during the training, we prove that the convergence property outstandingly improves. To the best of our knowledge, this is the first work to analyze the effectiveness of large batches under partial client participation. Our theoretical results indicate that, with the proper setting for the probability, FedAMD can achieve a convergence rate of  $O(\frac{M}{AT})$  under non-convex objective, and linear convergence under Polyak-Łojasiewicz (PL) condition (Polyak, 1963; Łojasiewicz, 1963).
- **Empirically**, we conduct extensive experiments to compare FedAMD with state-of-the-art approaches. The results provide evidence of the superiority of the proposed algorithm. Achieving the same test accuracy, FedAMD utilizes less computational power measured by the cumulative gradient complexity.

## 2. Related Work

In this section, we discuss the state-of-the-art works that are most relevant to our research. A more comprehensive review is provided in Appendix A.

**Mini-batch SGD vs. Local SGD.** Distributed optimization aims to train large-scale deep learning systems. Local SGD (also known as FedAvg) (Stich, 2018; Dieuleveut & Patel, 2019; Haddadpour et al., 2019; Haddadpour & Mahdavi, 2019) performs multiple (i.e.,  $K \geq 1$ ) local updates with  $K$  small batches, while mini-batch SGD computes the gradients averaged by  $K$  small batches (Woodworth et al., 2020b;a) (or with a large batch (Shallue et al., 2019; You et al., 2018; Goyal et al., 2017)) on a given model. There has been a long discussion on which one is better (Lin et al., 2019; Woodworth et al., 2020a;b; Yun et al., 2021), but there is no existing work applying both techniques at the same time to speed up model training.

**Variance Reduction in FL.** The variance reduction techniques have critically driven the advent of FL algorithms (Karimireddy et al., 2020b; Wu et al., 2021; Liang et al., 2019; Karimireddy et al., 2020a; Murata & Suzuki, 2021; Mitra et al., 2021) by correcting each local computed gradient with respect to the estimated global orientation. Roughly, the methods fall into two categories based on types of correction, namely, static and recursive. The former methods

(Karimireddy et al., 2020b; Wu et al., 2021; Liang et al., 2019; Karimireddy et al., 2020a; Mitra et al., 2021) use a consistent gradient correction within a training round. As the local updates go further, the static correction on a client is still dramatically biased to its local optimal solution, which hinders the training efficiency, especially under partial client participation. Recursive correction can avoid the challenge because it recursively corrects the local update based on the latest local model (Murata & Suzuki, 2021). However, existing works require all workers to attain the gradient at the starting point of each round, which is infeasible for federated learning settings. To the best of our knowledge, our work is the first to achieve recursive calibration under partial-client scenarios.

## 3. FedAMD

In this section, we comprehensively describe the technical details of FedAMD, a federated learning framework with anchor sampling. In specific, it separates the active participants into the anchor group and the miner group with time-varying probabilities. The pseudo-code is illustrated in Algorithm 1.

**Problem Formulation.** In an FL system with a total of  $M$  clients, the objective function is formalized as

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := \frac{1}{M} \sum_{m \in [M]} F_m(\mathbf{x}) \quad (1)$$

where we define  $[M]$  for a set of  $M$  clients.  $F_m(\cdot)$  indicates the local expected loss function for client  $m$ , which is unbiased estimated by empirical loss  $f_m(\cdot)$  using a random realization  $\mathcal{B}_m$  from the local training data  $\mathcal{D}_m$ , i.e.,  $\mathbb{E}_{\mathcal{B}_m \sim \mathcal{D}_m} f_m(\mathbf{x}, \mathcal{B}_m) = F_m(\mathbf{x})$ . We denote  $n$  by the size of a client’s local dataset, i.e.,  $|\mathcal{D}_m| = n$  for all  $m \in [M]$ , and  $n$  can be infinite large in the streaming/online cases.  $F_*$  represents the minimum loss for Equation (1).

**Algorithm Description.** In FedAMD, a global model is initialized with arbitrary parameters  $\tilde{\mathbf{x}}_0 \in \mathbb{R}^d$ . It is required to initialize the caching gradients, and there is one way to set them with  $v_0^{(m)} = \nabla f_m(\tilde{\mathbf{x}}_0, \mathcal{B}_{m,0})$  for all  $m \in [M]$ , where  $\mathcal{B}_{m,0}$  is a  $b$ -sample batch generated from  $\mathcal{D}_m$ . Other hyperparameters’ settings such as the learning rates and the minibatch sizes will be discussed in Section 4 and Appendix F from the theoretical and the empirical perspectives, respectively.

At the beginning of each round  $t \in \{0, 1, 2, \dots\}$ , the server randomly picks an  $A$ -client subset  $\mathcal{A}$  from  $M$  clients (Line 2). Since each client is independently selected without replacement, under the setting of Equation (1), clients have an equal chance to be selected with the probability of  $\frac{A}{M}$ . Subsequently, the server distributes the global model  $\tilde{\mathbf{x}}_t$

**Algorithm 1** FedAMD

**Input:** local learning rate  $\eta_l$ , global learning rate  $\eta_s$ , minibatch size  $b$ ,  $b' < b$ , local updates  $K$ , probability  $\{p_t \in [0, 1]\}_{t \geq 0}$ , initial model  $\tilde{\mathbf{x}}_0$ , initial caching gradient  $v_0 = \{v_0^{(m)}\}_{m \in [M]}$ .

- 1: **for**  $t = 0, 1, 2, \dots$  **do**
- 2:   Sample clients  $\mathcal{A} \subseteq [M]$
- 3:   Send  $\tilde{\mathbf{x}}_t$  and  $\tilde{g}_t = \text{avg}(v_t)$  to clients  $i \in \mathcal{A}$
- 4:   Initialize subsequent caching gradient  $v_{t+1} = v_t$
- 5:   **for**  $i \in \mathcal{A}$  **in parallel do**
- 6:     **if**  $\text{Bernoulli}(p_t) == 1$  **then**
- 7:        $v_{t+1}^{(i)} = \nabla f_i(\tilde{\mathbf{x}}_t, \mathcal{B}_{i,t})$  with  $|\mathcal{B}_{i,t}| = b$
- 8:       Send  $v_{t+1}^{(i)}$  to the server
- 9:     **else**
- 10:       Initialize  $\mathbf{x}_{t,-1}^{(i)} = \mathbf{x}_{t,0}^{(i)} = \tilde{\mathbf{x}}_t, g_{t,0}^{(i)} = \tilde{g}_t$
- 11:       **for**  $k = 0, \dots, K - 1$  **do**
- 12:          Generate random realization  $|\mathcal{B}'_{i,k}| = b'$
- 13:          Compute  $g_{t,k+1}^{(i)}$  via Equation (2)
- 14:           $\mathbf{x}_{t,k+1}^{(i)} = \mathbf{x}_{t,k}^{(i)} - \eta_l \cdot g_{t,k+1}^{(i)}$
- 15:       **end for**
- 16:        $\Delta \mathbf{x}_t^{(i)} = \tilde{\mathbf{x}}_t - \mathbf{x}_{t,K}^{(i)}$
- 17:       Send  $\Delta \mathbf{x}_t^{(i)}$  to the server
- 18:     **end if**
- 19:   **end for**
- 20:   Update  $\tilde{\mathbf{x}}_{t+1}$  via Equation (3)
- 21: **end for**

to the clients in the set  $\mathcal{A}$ , accompanying the global bullseye (i.e., the averaged caching gradient)  $\tilde{g}_t = \text{avg}(v_t) = \frac{1}{M} \sum_{m \in [M]} v_t^{(m)}$  (Line 3). With the probability of  $p_t$ , client  $i \in \mathcal{A}$  is classified for the anchor group (Line 7–8) or the miner group (Line 10–17), and different groups have different objectives and focus on different tasks.

**Anchor group (Line 7–8).** Clients in this group target to discover the bullseyes based on their local data distribution. According to Line 6, client  $i \in \mathcal{A}$  has the probability of  $p_t$  to become a member of this group. Then, the client utilizes a large batch  $\mathcal{B}_{i,t}$  with  $b$  samples to obtain the gradient  $v_{t+1}^{(i)} = \nabla f_i(\tilde{\mathbf{x}}_t, \mathcal{B}_{i,t})$  (Line 6). Therefore, following the gradient  $v_{t+1}^{(i)}$  can find an optimal or near-optimal solution for client  $i$ . Next, the client pushes the gradient to the server and updates the caching gradient (Line 7). In view that some clients do not participate in the anchor group for obtaining the bullseyes at round  $t$ , the server spontaneously inherits their previous calculation from  $v_t$  (Line 4). As a result,  $\tilde{g}_t$  in Line 3 indicates an approximate orientation towards global optimal parameters, which directs the local update in the miner group and affects the final global update. Besides,  $v_{t+1}^{(i)}$  influences the training from round  $t + 1$  up to the next

time when client  $i$  is a member of anchor group.

**Miner group (Line 10–17).** Guided by the global bullseye, clients in the miner group perform multiple local updates and finally drive the update of the global model. First, client  $i$  initializes the model with  $\tilde{\mathbf{x}}_t$  and the target direction with  $\tilde{g}_t$  (Line 10). Ideally, in the subsequent  $K$  updates (Line 11), client  $i$  update the model with the gradient  $\nabla F(\mathbf{x}_{t,k}^{(i)})$  for  $k \in \{0, \dots, K - 1\}$ . This is impractical because clients cannot access all others' training sets to compute the noise-free gradients. Instead, the client at  $k$ -th iteration generates a  $b'$ -sample realization  $\mathcal{B}'_{i,k}$  (Line 12) and calculates the update  $g_{t,k+1}^{(i)}$  via a variance-reduced technique, i.e.,

$$g_{t,k+1}^{(i)} = g_{t,k}^{(i)} - \nabla f_i(\mathbf{x}_{t,k-1}^{(i)}, \mathcal{B}'_{i,k}) + \nabla f_i(\mathbf{x}_{t,k}^{(i)}, \mathcal{B}'_{i,k}). \quad (2)$$

The update  $g_{t,k+1}^{(i)}$  is approximate to  $\nabla F(\mathbf{x}_{t,k}^{(i)})$  for two reasons: (i) the first term is used to estimate the global update because  $g_{t,0}^{(i)}$  stores the global bullseye; and (ii) the rest terms remove the perturbation of data heterogeneity and reflect the true update at  $\mathbf{x}_{t,k}^{(i)}$ . Therefore, the local model update follows  $\mathbf{x}_{t,k+1}^{(i)} = \mathbf{x}_{t,k}^{(i)} - \eta_l \cdot g_{t,k+1}^{(i)}$  (Line 20). After  $K$  local updates, the model changes on client  $i$  is  $\Delta \mathbf{x}_t^{(i)} = \tilde{\mathbf{x}}_t - \mathbf{x}_{t,K}^{(i)}$ . Then, the client transmits  $\Delta \mathbf{x}_t^{(i)}$  to the server for the purpose of global model update.

**Server (Line 4 and Line 20).** After the local training on the active participants, the server merges the model changes from the miner group (Line 20) and updates the caching gradients from the anchor group (Line 4). At  $t$ -th round, let  $\Delta \mathbf{x}_t$  be the subset of  $\mathcal{A}$  acting as miners, i.e.,  $\Delta \mathbf{x}_t = \{\Delta \mathbf{x}_t^{(a_1)}, \dots, \Delta \mathbf{x}_t^{(a_\nu)}\}$ , where  $\{a_1, \dots, a_\nu\} \subset \mathcal{A}$  and  $0 \leq \nu \leq A$ . Therefore, the global update follows

$$\tilde{\mathbf{x}}_{t+1} = \begin{cases} \tilde{\mathbf{x}}_t - \eta_s \cdot \sum_{j=1}^{\nu} \Delta \mathbf{x}_t^{(a_j)} / \nu, & \nu \geq 1 \\ \tilde{\mathbf{x}}_t, & \nu = 0 \end{cases} \quad (3)$$

The reason why we solely use the changes from the miner group is that clients perform multiple local updates regulated by the global target such that the model changes walk towards the global optimal solution. While directly incorporating the new gradients from the anchor group, the global model has a degraded performance because they perform a single update that aims to find out the local bullseye deviated from the global target. Implicitly, clients in the miner group take in the update of caching gradients at iteration  $k = 0$  to update the local model, which will affect the next global parameters.

**Previous Algorithms as Special Cases.** The probabilities can vary among the rounds that disjoint the participants into the anchor group and the miner group. By setting



$A = M$ , and the probability  $\{p_t\}$  following the sequence of  $\{1, 0, 1, 0, \dots\}$ , FedAMD reduces to distributed mini-batch SGD ( $K = 1$ ) or BVR-L-SGD ( $K > 1$ ). Therefore, FedAMD subsumes the existing algorithms and takes partial client participation into consideration. To obtain the best performance, we should tune the settings of  $\{p_t\}$  and  $K$ . However, accounting for the generality of FedAMD, it faces substantial additional hurdles in its convergence analysis, which is one of our main contributions, as detailed in the following section.

**Discussion on Communication Overhead.** As the anchors are not necessary to obtain the averaged caching gradient (i.e.,  $\tilde{g}_t$ ) at  $t$ -th round, the centralized server solely distributes  $\tilde{g}_t$  to the miners. Compared to FedAvg, the proposed algorithm requires  $(1 - p_t)/2$  more communication costs, but it achieves convergence with at least  $O(\frac{1}{\epsilon})$  less communication rounds (see Table 1). Therefore, from the perspective of model training progress, FedAMD is more communication efficient than FedAvg.

**Discussion on Massive-Client Settings.** A typical example of this scenario is cross-device FL (Kairouz et al., 2019). In this setting, it is not a wise option for the server to preserve all the caching gradients for clients. Therefore, the clients retain their caching gradients while the server keeps their average. Firstly, at  $t$ -th round, client  $i \in [M]$  copies their caching gradient to  $(t + 1)$ -th round, i.e.,  $v_{t+1}^{(i)} = v_t^{(i)}$ . For the client  $i$  in the anchor group, they will follow Line 13 in Algorithm 1 to update  $v_{t+1}^{(i)}$  and push  $\wedge_t^{(i)} = v_{t+1}^{(i)} - v_t^{(i)}$  to the server. After the server receives the updates of all local caching gradients, it performs  $v_{t+1} = v_t + \frac{1}{M} \sum \wedge_t$ , where  $\wedge_t$  aggregates  $\wedge_t^{(i)}$  where client  $i$  is in anchor group.

## 4. Theoretical Analysis

In this section, we analyze the convergence rate of FedAMD under non-convex objectives with respect to  $\epsilon$ -approximation, i.e.,  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$ . Specifically, when it comes to PL condition,  $\epsilon$ -approximation refers to  $F(\tilde{\mathbf{x}}_T) - F(\mathbf{x}_*) \leq \epsilon$ . In the following discussion, we particularly highlight the setting of  $\{p_t\}$  to obtain the best performance. Before showing the convergence result, we make the following assumptions, where the first two assumptions have been widely used in machine learning studies (Karimireddy et al., 2020b; Li et al., 2020), while the last one has been adopted in some recent works (Gorbunov et al., 2021; Tyurin & Richtárik, 2022; Murata & Suzuki, 2021).

**Assumption 4.1** (L-smooth). The local objective functions are Lipschitz smooth: For all  $v, \bar{v} \in \mathbb{R}^d$ ,

$$\|\nabla F_i(v) - \nabla F_i(\bar{v})\|_2 \leq L\|v - \bar{v}\|_2, \quad \forall i \in [M].$$

**Assumption 4.2** (Bounded Noise). For all  $v \in \mathbb{R}^d$ , there exists a scalar  $\sigma \geq 0$  such that

$$\mathbb{E}_{\mathcal{B} \sim \mathcal{D}_i} \|\nabla f_i(v, \mathcal{B}) - \nabla F_i(v)\|_2^2 \leq \frac{\sigma^2}{|\mathcal{B}|}, \quad \forall i \in [M].$$

**Assumption 4.3** (Average L-smooth). For all  $v, \bar{v} \in \mathbb{R}^d$ , there exists a scalar  $L_\sigma \geq 0$  such that

$$\begin{aligned} \mathbb{E}_{\mathcal{B} \sim \mathcal{D}_i} \|\nabla f_i(v, \mathcal{B}) - \nabla f_i(\bar{v}, \mathcal{B}) - (\nabla F_i(v) - \nabla F_i(\bar{v}))\|_2^2 \\ \leq \frac{L_\sigma^2}{|\mathcal{B}|} \|v - \bar{v}\|_2^2, \quad \forall i \in [M]. \end{aligned}$$

**Remark.** Assumption 4.3 definitely provides a tighter bound for the patterns of variance reduction. In fact, solely with Assumption 4.2, the term  $\mathbb{E}_{\mathcal{B} \sim \mathcal{D}_i} \|\nabla f_i(v, \mathcal{B}) - \nabla f_i(\bar{v}, \mathcal{B}) - (\nabla F_i(v) - \nabla F_i(\bar{v}))\|_2^2$  can be bounded by a constant. Additionally, the above term will approach 0 as  $v$  and  $\bar{v}$  are very close. Therefore, we could find a coefficient for  $\|v - \bar{v}\|_2^2$ , and we define it as  $L_\sigma^2/|\mathcal{B}|$ . Furthermore, if the loss function is Lipschitz smooth, e.g., cross-entropy loss (Tewari & Chaudhuri, 2015), we can derive a similar structure as presented in Assumption 4.3.

### 4.1. Sequential Probability Settings

As mentioned in Section 3, a recursive pattern appearing in the probability sequence  $\{p_t \in \{0, 1\}\}_{t \geq 0}$  can reduce FedAMD to the existing works. We assume that the caching gradient updates every  $\tau (\geq 2)$  rounds, such that

$$p_t = \begin{cases} 1, & t \bmod \tau == 0 \\ 0, & \text{Otherwise} \end{cases} \quad (4)$$

We derive the following results under sequential probability settings. The complete proofs are provided in Appendix D. In detail, Theorem 4.4 and Theorem 4.6 are proved in Appendix D.2 and Appendix D.3, respectively.

**Theorem 4.4** (Full Client Participation). *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local updates  $K \geq 1$ , the setting for the local learning rate  $\eta_l$  and the global learning rate  $\eta_s$  satisfy the following two constraints: (1)  $\eta_s \eta_l \leq \frac{1}{6\sqrt{6}KL\tau}$ ; and (2)  $\eta_l \leq \min\left(\frac{1}{6\sqrt{2}KL}, \frac{\sqrt{b'/K}}{2\sqrt{3}L_\sigma}\right)$ . Then, the convergence rate of FedAMD with the probability setting of Equation (4) for non-convex objectives should be*

$$\begin{aligned} \min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq O\left(\frac{F(\tilde{\mathbf{x}}_0) - F_*}{\eta_s \eta_l K T}\right) \\ + O\left((\eta_s + \eta_l K L) \eta_l K L \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{M b}\right) \end{aligned} \quad (5)$$

**Corollary 4.5.** *Under the setting of Theorem 4.4, we assume that the number of rounds  $T$  is sufficiently large. To find an  $\epsilon$ -approximation of non-convex objectives, i.e.,  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$ , the number of communication rounds  $T$  performed by FedAMD with the probability setting of Equation (4) is*

- $b = O(1)$ :  $T = O\left(\frac{\sigma^2}{M\epsilon^2}\right)$  by setting  $\eta_l = \frac{1}{K}\left(\frac{M}{\sigma^2 T}\right)^{1/4}$  and  $\eta_s = \left(\frac{M}{\sigma^2 T}\right)^{1/4}$
- $b = O(K)$ :  $T = O\left(\frac{\sigma^2}{MK\epsilon^2}\right)$  by setting  $\eta_l = \frac{1}{\sqrt{TK}\sigma}$  and  $\eta_s = \sqrt{KM}$
- $b = \min\{n, \frac{\sigma^2}{M\epsilon}\}$ :  $T = O\left(\frac{\tau}{\epsilon}\right)$  by setting  $K = \frac{L_\sigma^2}{6b'L^2}$ ,  $\eta_l = \frac{1}{6\sqrt{2KL}}$ , and  $\eta_s = \frac{1}{\sqrt{3\tau}}$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

**Effectiveness of Using Large Batches.** The corollary summarizes the convergence results under various settings of the large-batch sizes. When the anchors update the caching gradients with an (averaged) SGD gradient, the convergence rate remains  $O(1/\sqrt{T})$ , which is consistent with the algorithms using small batches only, e.g., SCAF-FOLD (Karimireddy et al., 2020b). When we apply large-batch gradients to model training, i.e., compute the caching gradients with the entire training set, the convergence rate significantly lifts to  $O(1/T)$ .

**Comparison with BVR-L-SGD.** As discussed in Section 3, FedAMD reduces to BVR-L-SGD (Murata & Suzuki, 2021) when  $\tau = 2$  and all clients participate in the training. In this case, Corollary 4.5 shows a total of  $T = O(1/\epsilon)$  communication rounds are needed. This result coincides with the complexity of BVR-L-SGD in Table 1 by the setting that (1)  $nM \leq \frac{1}{\epsilon}$ , and (2)  $K \geq \sqrt{n/M}$ . In other words, we theoretically prove that BVR-L-SGD still achieves  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$  with  $T = O(1/\epsilon)$  in a looser constraint. As for computation overhead, our proposed method requires  $O\left(\frac{\sigma^2}{\epsilon^2} + \frac{MK}{\epsilon}\right)$ , which is less than BVR-L-SGD (i.e.,  $O\left(\frac{\sigma^2}{\epsilon^2} + \frac{\sigma^2 + MK}{\epsilon} + MK\right)$ ).

In addition to the convergence result on full-client participation, we provide the analysis of the partial-client scenarios.

**Theorem 4.6 (Partial Client Participation).** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local updates  $K \geq 1$ , the minibatch size  $b = \min\left(\frac{\sigma^2}{M\epsilon}, n\right)$  and  $b' < b$ . Additionally, the settings for the local learning rate  $\eta_l$  and the global learning rate  $\eta_s$  satisfy the following two constraints: (1)  $\eta_s \eta_l = \frac{1}{KL} \left(1 + \frac{2M\tau}{A}\right)^{-1}$ ; and (2)  $\eta_l \leq \min\left(\frac{1}{2\sqrt{6KL}}, \frac{\sqrt{b'/K}}{4\sqrt{3L}\sigma}\right)$ . Then, to find an  $\epsilon$ -approximation*

*of non-convex objectives, i.e.,  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$ , the number of communication rounds  $T$  performed by FedAMD with the probability setting of Equation (4) is*

$$T = O\left(\left(1 + \frac{2M\tau}{A}\right) \cdot \frac{\tau}{\tau - 1} \cdot \frac{1}{\epsilon}\right)$$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

**Discussion on the selection of  $\tau$ .** According to Theorem 4.6, we notice that  $\tau = 2$  achieves  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$  with the fewest communication rounds. The following corollary discloses the relation between computation overhead and the value of  $\tau$ .

**Corollary 4.7.** *Under the setting of Theorem 4.6, FedAMD computes  $O\left(\frac{Mb}{\epsilon} + \frac{\tau MKb'}{\epsilon}\right)$  gradients and consumes a communication overhead of  $O\left(\frac{M\tau}{A\epsilon}\right)$  during the model training.*

**Remark.** FedAMD requires an increasing computation and communication cost as  $\tau$  gets larger. Therefore,  $\tau = 2$  possesses the most outstanding performance in the sequential probability settings. In this case, FedAMD requires the communication rounds of  $O\left(\frac{M}{A\epsilon}\right)$ , the communication overhead of  $O\left(\frac{M}{A\epsilon}\right)$ , and the computation cost of  $O\left(\frac{\sigma^2}{\epsilon^2} + \frac{MK}{\epsilon}\right)$  while optimizing an online scenario where the size of local dataset is infinity large.

## 4.2. Constant Probability Settings

Apparently, when we set the constant probability as 1, all participants are in the anchor group such that the model cannot be updated. Likewise, when the constant probability is 0, all participants are in the miner group such that the global target cannot be updated, leading to degraded performance. Therefore, we manually define a constant  $p \in (0, 1)$  such that  $\{p_t = p\}_{t \geq 0}$ . In this section, we derive the following results with partial client participation. Detailed proof is provided in Appendix E. Specifically, Appendix E.2 and Appendix E.3 proves the convergence rate for Theorem 4.8 and Theorem 4.11, respectively.

**Theorem 4.8.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local updates  $K \geq \max\left(1, \frac{2L_\sigma^2}{b'L^2}\right)$ , the minibatch size  $b = \min\left(\frac{\sigma^2}{M\epsilon}, n\right)$  and  $b' < b$ , the local learning rate  $\eta_l = \frac{1}{2\sqrt{6KL}}$ , and the global learning rate  $\eta_s = \frac{2\sqrt{6}}{1 + \frac{2M}{Ap}\sqrt{1-p^A}}$ . Then, to find an  $\epsilon$ -approximation of non-convex objectives, i.e.,  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$ , the number of communication rounds  $T$  performed by FedAMD*

with constant probability  $p_t = p$  is

$$T = O\left(\frac{1}{\epsilon} \left( \frac{1}{1-p^A} + \frac{M}{Ap\sqrt{1-p^A}} \right)\right)$$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

With the constant probability  $p$  approaching 0 or 1, Theorem 4.8 shows that FedAMD requires a significant number of communication rounds. Hence, there is an optimal  $p$  such that FedAMD achieves convergence with the fewest communication rounds. In view that  $M \geq A$ , the number of communication rounds is dominated by  $O\left(\frac{M}{Ap\sqrt{1-p^A}} \cdot \frac{1}{\epsilon}\right)$ . Based on this observation, the following corollary provides the settings for the constant probability  $p$  that leads to the optimal convergence result. Based on the value of  $p$ , we further refine the settings for other parameters. The following corollary takes  $b' = 1$  into consideration, i.e., the small batch size is 1.

**Corollary 4.9.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the constant probability  $p = \frac{1}{c} \left(\frac{2}{A+2}\right)^{1/A}$ , where  $c$  is a constant greater than or equal to 1, the local updates  $K \geq \max\left(1, \frac{2L^2}{L^2}\right)$ , the minibatch size  $b = \min\left(\frac{\sigma^2}{M\epsilon}, n\right)$  and  $b' = 1$ , the local learning rate  $\eta_l = \frac{1}{2\sqrt{6}KL}$ , and the global learning rate  $\eta_s = \frac{2\sqrt{6}A}{A+3Mc}$ . Then, after the communication rounds of  $T = O\left(\frac{M}{A\epsilon}\right)$ , we have  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$ . Therefore, the number of total samples called by all clients (i.e., cumulative gradient complexity) is  $O\left(\frac{\sigma^2}{\epsilon^2} + \frac{MK}{\epsilon}\right)$  when it optimizes an online scenario.*

**Discussion on the effectiveness of  $c$ .** When  $c = 1$ , we can obtain the minimum value for the term  $\left(p\sqrt{1-p^A}\right)^{-1}$  with the constant  $p$  devised in Corollary 4.9. As the number of participants (i.e.,  $A$ ) gets larger, the optimal  $p$  increases as well and tends to be 1, indicating that most participants are in the anchor group. When we optimize an online scenario, the anchors compute a gradient with massive samples. As a result, the computation overhead of a single round is not acceptable. By deducting the anchor sampling probability to its  $1/c$ , FedAMD consumes up to  $(c-1)/c$  less computation overhead, while its convergence performance remains.

**Comparison with FedAvg.** As a classical algorithm, FedAvg (Yang et al., 2020) requires  $O\left(\frac{K}{A\epsilon^2} + \frac{1}{\epsilon}\right)$  communication rounds to achieve  $\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \epsilon$  with the total computation consumption of  $O\left(\frac{K^2}{\epsilon^2} + \frac{AK}{\epsilon}\right)$ . Apparently, FedAMD needs  $O\left(\frac{1}{\epsilon}\right)$  fewer communication rounds. As mentioned in Section 3, FedAMD consumes  $(1-p)/2$

more communication overhead than FedAvg. As  $\epsilon$  is close to 0, the total communication overhead for FedAMD is far less than the cost of FedAvg. As for computation overhead, (Yang et al., 2020) implicitly assumes that  $K \geq \sigma^2$ , FedAMD is more communication friendly than FedAvg.

In addition to the generalized non-convex objectives, we investigate the convergence rate of the PL condition, a special case under non-convex objectives. The following assumption describes this case:

**Assumption 4.10** (PL Condition (Karimi et al., 2016)). The objective function  $F$  satisfies the PL condition when there exists a scalar  $\mu > 0$  such that

$$\|\nabla F(v)\|_2^2 \geq 2\mu(F(v) - F_*), \quad \forall v \in \mathbb{R}^d.$$

Under PL condition, the rest of the section draws the convergence performance of FedAMD with partial client participation.

**Theorem 4.11.** *Suppose that Assumption 4.1, 4.2, 4.3 and 4.10 hold. Let the local updates  $K \geq \max\left(1, \frac{2L^2}{b'L^2}\right)$ , the minibatch size  $b = \min\left(\frac{\sigma^2}{M\mu\epsilon}, n\right)$  and  $b' < b$ , the local learning rate  $\eta_l = \frac{1}{2\sqrt{6}KL}$ , and the global learning rate  $\eta_s = \min\left(\frac{2\sqrt{6}LAp}{M\mu(1-p^A)}, \frac{2\sqrt{6}}{1+\frac{16ML}{\mu Ap}}\right)$ . Then, to find an  $\epsilon$ -approximation of PL condition, i.e.,  $F(\tilde{\mathbf{x}}_T) - F(\mathbf{x}_*) \leq \epsilon$ , the number of communication rounds  $T$  performed by FedAMD with constant probability  $p_t = p$  is*

$$T = O\left(\frac{1}{\mu(1-p^A)} \left(1 + \frac{M}{\mu Ap} + \frac{M\mu(1-p^A)}{Ap}\right) \log \frac{1}{\epsilon}\right)$$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

Similar to Theorem 4.8, the number of communication rounds of FedAMD is mainly occupied by  $O\left(\frac{M}{\mu^2 Ap(1-p^A)} + \frac{M}{Ap}\right)$ . According to such an approximation, we provide a mathematical expression for the setting of  $p$  in Corollary 4.12. Subsequently, we adjust the value of the hyper-parameters such that we can obtain the best result for Theorem 4.11.

**Corollary 4.12.** *Suppose that Assumption 4.1, 4.2, 4.3 and 4.10 hold. Let the constant probability  $p = \frac{1}{c} \left( \left(1 + \frac{A+1}{2\mu^2}\right) - \sqrt{\left(\frac{A+1}{2\mu^2}\right)^2 + \frac{A}{\mu^2}} \right)^{1/A}$ , where  $c$  is a constant greater than or equal to 1, the local updates  $K \geq \max\left(1, \frac{2L^2}{L^2}\right)$ , the minibatch size  $b = \min\left(\frac{\sigma^2}{M\mu\epsilon}, n\right)$  and  $b' = 1$ , the local learning rate  $\eta_l = \frac{1}{2\sqrt{6}KL}$ , and the global learning rate  $\eta_s =$*

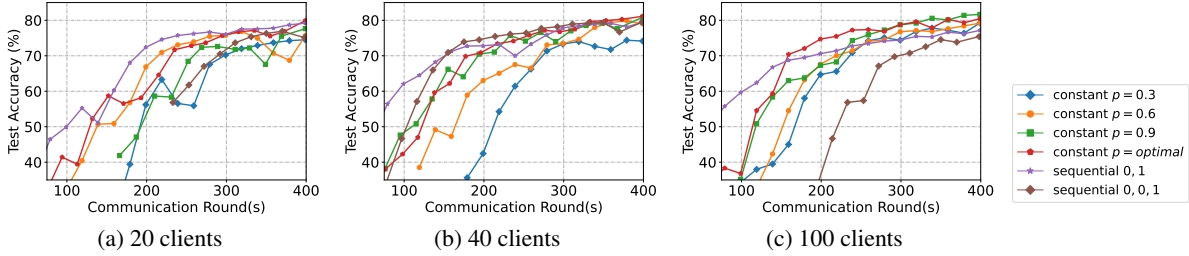


Figure 1. Comparison of different probability settings using test accuracy against the communication rounds for FedAMD.

Method	20 clients				40 clients				100 clients			
	Grad.	Comm.	Round	Acc.	Grad.	Comm.	Round	Acc.	Grad.	Comm.	Round	Acc.
BVR-L-SGD	101.9	853.9	310	77.0	177.4	1492.3	275	78.6	463.8	3908.7	<u>291</u>	<u>78.0</u>
FedAvg	<u>40.8</u>	<u>566.9</u>	318	76.4	<u>78.3</u>	<u>1087.5</u>	305	76.7	<u>186.9</u>	<u>2594.5</u>	<u>291</u>	77.0
FedPAGE	148.6	1670.4	<u>271</u>	<u>79.2</u>	164.3	1622.4	<u>203</u>	<u>80.6</u>	525.5	4540.3	339	77.1
SCAFFOLD	47.5	1318.6	370	75.9	91.4	2537.8	356	76.0	250.9	6966.0	391	75.1
FedAMD (constant)	<b>35.5</b>	<b>489.8</b>	<b>259</b>	<b>80.6</b>	<b>55.0</b>	<b>776.3</b>	<b>209</b>	<b>82.3</b>	<b>153.9</b>	<b>2147.8</b>	<b>229</b>	<b>83.4</b>
FedAMD (sequential)	<i>40.2</i>	<i>475.3</i>	<i>213</i>	<i>79.5</i>	<i>80.4</i>	<i>904.5</i>	<i>190</i>	<i>80.8</i>	<i>253.8</i>	<i>2998.8</i>	<i>269</i>	<i>78.7</i>

Table 2. Comparison among baselines in terms of cumulative gradient complexity ( $\times 10^5$  samples), communication costs ( $\times 32$  Mbits), and rounds reaching the accuracy of 75%, and the final accuracy (%) after 400 rounds. **Bold**: The best result in each column; underline: The best result of the baselines in each column; *Italy*: The results that FedAMD outperforms all baselines in each column.

$\min \left( \frac{\sqrt{6}AL}{M\mu c}, 2\sqrt{6} \left( 1 + \frac{32Mc}{\mu A} \right)^{-1} \right)$ . Then, after the communication rounds of  $T = O \left( \left( \frac{1}{\mu} + \frac{M}{\mu^2 A} + \frac{M}{A} \right) \log \frac{1}{\epsilon} \right)$ , we have  $F(\tilde{\mathbf{x}}_T) - F(\mathbf{x}_*) \leq \epsilon$ . Therefore, the number of total samples called by all clients (i.e., cumulative gradient complexity) is  $O \left( \left( \frac{A}{\mu} + \frac{M}{\mu^2} + M \right) \left( \frac{\sigma^2}{M\mu\epsilon} + K \right) \log \frac{1}{\epsilon} \right)$  when it optimizes an online scenario.

**Remark.** When  $c = 1$ , the probability  $p$  for anchor sampling approaches 100% as the number of participants is increasing. Likewise, it is necessary to use  $c \geq 1$  to reduce the computation consumption of each round. Besides, FedAMD achieves a linear convergence under PL conditions. In view that strongly-convex objectives possess a looser setting than PL conditions, FedAMD can also achieve linear convergence under strongly-convex objectives.

## 5. Experiments

This section presents the experiments of our proposed approach and other existing baselines that are most relative to this work. We also investigate the effectiveness of probability  $\{p_t\}_{t \geq 0}$ . Account for the limited space, numerical analysis on other factors like the number of local updates could be found in Appendix F. Our code is available at <https://github.com/HarliWu/FedAMD>.

**Experimental setup.** We train a convolutional neural network LeNet-5 (LeCun et al., 2015; 1989) using Fashion MNIST (Xiao et al., 2017) and 2-layer MLP on EMNIST digits (Cohen et al., 2017). We conduct the experiments with a total of 100 clients. To simulate the non-i.i.d. features, each client holds the data from 2 classes. The default setting in this section is  $K = 10$ ,  $b' = 64$  and  $b$  is the size of the local training set. For different experiments, unless specified, we leverage the best setting (e.g., learning rate) to obtain the best results. All the numerical results in this section represent the average performance of three experiments using different random seeds.

**Effectiveness of probability  $\{p_t\}_{t \geq 0}$ .** Figure 1 demonstrates the performance of various probability settings under the scenarios of different participants. In Figure 1a with 20 clients, both sequential probability setting and constant probability setting achieve the best performance. In Figure 1b and 1c, where 40 and 100 clients are selected in each round, constant probability setting outperforms sequential probability setting. In all three scenarios, with sequential probability settings, the pattern of  $\{0, 0, 1\}$  has a much worse performance than the pattern of  $\{0, 1\}$ . This empirically validates Theorem 4.6 for the best setting  $\tau = 2$  in terms of communication complexity. Similarly, with constant probability settings, the best performance is achieved when  $p$  approximates or equals optimal, which validates the statement in Corollary 4.9.



**Comparison with the state-of-the-art works.** Table 2 compares FedAMD with the existing works under partial/full client participation. At first glance, FedAMD outperforms other baselines because the texts in bold are all appeared in FedAMD. With 20-client participation, our proposed method surpasses four baselines all aroundness. As for 40-client participation, FedAMD with constant probability saves at least 30% computation and communication cost, and the final accuracy realizes up to 6% improvement. In terms of the training with full client participation, FedAMD requires 10%–20% fewer communication rounds, and its final accuracy has significant improvement, i.e., within the range of 0.7%–8.3%. Also, it is well noted that BVR-L-SGD has a similar performance as FedAMD using sequential probability in terms of communication rounds and test accuracy, but the former needs more computation and communication overhead. This is because BVR-L-SGD computes the bullseye using multiple  $b'$ -size batches rather than a large batch.

## 6. Conclusions

In this work, we investigate a federated learning framework FedAMD that separates the partial participants into anchor and miner groups. We provide the convergence analysis of the proposed algorithm for constant and sequential probability settings. Under the partial-client scenario, FedAMD achieves sublinear speedup under non-convex objectives and linear speedup under the PL condition. To the best of our knowledge, this is the first work to analyze the effectiveness of large batches under partial client participation. Experimental results demonstrate that FedAMD is superior to state-of-the-art works.

## Acknowledgements

The authors would like to thank the anonymous reviewers for their constructive comments. This work is supported in part by the US National Science Foundation under grants NSF IIS-2226108 and NSF IIS-1747614, in part by the Key-Area Research and Development Program of Guangdong Province under grant No. 2021B0101400003, in part by Hong Kong RGC Research Impact Fund under grant No. R5060-19, in part by Areas of Excellence Scheme under grant AoE/E-601/22-R, in part by General Research Fund under grants No. 152203/20E, 152244/21E, 152169/22E, in part by Shenzhen Science and Technology Innovation Commission under grant JCYJ20200109142008673, in part by the National Natural Science Foundation of China under grant 62102131, and in part by Natural Science Foundation of Jiangsu Province under grant BK20210361. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

## References

- Allen-Zhu, Z. and Hazan, E. Variance reduction for faster non-convex optimization. In *International conference on machine learning*, pp. 699–707. PMLR, 2016.
- Avdiukhin, D. and Kasiviswanathan, S. Federated learning under arbitrary communication patterns. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 425–435. PMLR, 2021.
- Bietti, A. and Mairal, J. Stochastic optimization with variance reduction for infinite datasets with finite sum structure. *Advances in Neural Information Processing Systems*, 30, 2017.
- Blum, A., Haghtalab, N., Phillips, R. L., and Shao, H. One for one, or all for all: Equilibria and optimality of collaboration in federated learning. *arXiv preprint arXiv:2103.03228*, 2021.
- Bottou, L. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pp. 177–186. Springer, 2010.
- Bottou, L. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pp. 421–436. Springer, 2012.
- Bottou, L. and Cun, Y. Large scale online learning. *Advances in neural information processing systems*, 16, 2003.
- Bottou, L., Curtis, F. E., and Nocedal, J. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- Cohen, G., Afshar, S., Tapson, J., and Van Schaik, A. Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pp. 2921–2926. IEEE, 2017.
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Mao, M., Ranzato, M., Senior, A., Tucker, P., Yang, K., et al. Large scale distributed deep networks. *Advances in neural information processing systems*, 25, 2012.
- Defazio, A., Bach, F., and Lacoste-Julien, S. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in neural information processing systems*, 27, 2014.
- Diao, E., Ding, J., and Tarokh, V. Heterofl: Computation and communication efficient federated learning for heterogeneous clients. *arXiv preprint arXiv:2010.01264*, 2020.
- Dieuleveut, A. and Patel, K. K. Communication trade-offs for local-sgd with large step size. *Advances in Neural Information Processing Systems*, 32, 2019.

- Eichner, H., Koren, T., McMahan, B., Srebro, N., and Talwar, K. Semi-cyclic stochastic gradient descent. In *International Conference on Machine Learning*, pp. 1764–1773. PMLR, 2019.
- Fang, C., Li, C. J., Lin, Z., and Zhang, T. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in Neural Information Processing Systems*, 31, 2018.
- Ghadimi, S., Lan, G., and Zhang, H. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1): 267–305, 2016.
- Gorbunov, E., Burlachenko, K. P., Li, Z., and Richtárik, P. Marina: Faster non-convex distributed learning with compression. In *International Conference on Machine Learning*, pp. 3788–3798. PMLR, 2021.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Gu, X., Huang, K., Zhang, J., and Huang, L. Fast federated learning in the presence of arbitrary device unavailability. *Advances in Neural Information Processing Systems*, 34, 2021.
- Haddadpour, F. and Mahdavi, M. On the convergence of local descent methods in federated learning. *arXiv preprint arXiv:1910.14425*, 2019.
- Haddadpour, F., Kamani, M. M., Mahdavi, M., and Cadambe, V. Local sgd with periodic averaging: Tighter analysis and adaptive synchronization. *Advances in Neural Information Processing Systems*, 32, 2019.
- He, S., Yan, Q., Wu, F., Wang, L., Lécuyer, M., and Beschastnikh, I. Gluefl: Reconciling client sampling and model masking for bandwidth efficient federated learning. *arXiv preprint arXiv:2212.01523*, 2022.
- Horváth, S. and Richtárik, P. Nonconvex variance reduced optimization with arbitrary sampling. In *International Conference on Machine Learning*, pp. 2781–2789. PMLR, 2019.
- Horváth, S., Lei, L., Richtárik, P., and Jordan, M. I. Adaptivity of stochastic gradient methods for nonconvex optimization. *arXiv preprint arXiv:2002.05359*, 2020.
- Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26, 2013.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., and et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 795–811. Springer, 2016.
- Karimireddy, S. P., Jaggi, M., Kale, S., Mohri, M., Reddi, S. J., Stich, S. U., and Suresh, A. T. Mime: Mimicking centralized stochastic algorithms in federated learning. *arXiv preprint arXiv:2008.03606*, 2020a.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A. T. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020b.
- Khaled, A. and Richtárik, P. Better theory for sgd in the non-convex world. *arXiv preprint arXiv:2002.03329*, 2020.
- Konečný, J., McMahan, B., and Ramage, D. Federated optimization: Distributed optimization beyond the datacenter. *arXiv preprint arXiv:1511.03575*, 2015.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Lai, F., Zhu, X., Madhyastha, H. V., and Chowdhury, M. Oort: Efficient federated learning via guided participant selection. In *15th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 21)*, pp. 19–35, 2021.
- Lan, G. and Zhou, Y. An optimal randomized incremental gradient method. *Mathematical programming*, 171(1): 167–215, 2018a.
- Lan, G. and Zhou, Y. Random gradient extrapolation for distributed and stochastic optimization. *SIAM Journal on Optimization*, 28(4):2753–2782, 2018b.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- LeCun, Y. et al. Lenet-5, convolutional neural networks. URL: <http://yann.lecun.com/exdb/lenet>, 20(5):14, 2015.
- Lei, L., Ju, C., Chen, J., and Jordan, M. I. Non-convex finite-sum optimization via scsg methods. *Advances in Neural Information Processing Systems*, 30, 2017.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous

- networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020.
- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations*, 2019.
- Li, Z. Ssrgd: Simple stochastic recursive gradient descent for escaping saddle points. *Advances in Neural Information Processing Systems*, 32, 2019.
- Li, Z., Bao, H., Zhang, X., and Richtárik, P. Page: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In *International Conference on Machine Learning*, pp. 6286–6295. PMLR, 2021.
- Lian, X., Wang, M., and Liu, J. Finite-sum composition optimization via variance reduced gradient descent. In *Artificial Intelligence and Statistics*, pp. 1159–1167. PMLR, 2017.
- Liang, X., Shen, S., Liu, J., Pan, Z., Chen, E., and Cheng, Y. Variance reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019.
- Lin, T., Stich, S. U., Patel, K. K., and Jaggi, M. Don’t use large mini-batches, use local sgd. In *International Conference on Learning Representations*, 2019.
- Lojasiewicz, S. A topological property of real analytic subsets. *Coll. du CNRS, Les équations aux dérivées partielles*, 117(87-89):2, 1963.
- Malinovskiy, G., Kovalev, D., Gasanov, E., Condat, L., and Richtarik, P. From local sgd to local fixed-point methods for federated learning. In *International Conference on Machine Learning*, pp. 6692–6701. PMLR, 2020.
- Mason, L., Baxter, J., Bartlett, P., and Frean, M. Boosting algorithms as gradient descent. *Advances in neural information processing systems*, 12, 1999.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Mitra, A., Jaafar, R., Pappas, G. J., and Hassani, H. Linear convergence in federated learning: Tackling client heterogeneity and sparse gradients. *Advances in Neural Information Processing Systems*, 34:14606–14619, 2021.
- Murata, T. and Suzuki, T. Bias-variance reduced local sgd for less heterogeneous federated learning. *arXiv preprint arXiv:2102.03198*, 2021.
- Nesterov, Y. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- Nguyen, L. M., Liu, J., Scheinberg, K., and Takáč, M. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pp. 2613–2621. PMLR, 2017.
- Pathak, R. and Wainwright, M. J. Fedsplit: an algorithmic framework for fast federated optimization. *Advances in Neural Information Processing Systems*, 33:7057–7066, 2020.
- Philippenko, C. and Dieuleveut, A. Bidirectional compression in heterogeneous settings for distributed or federated learning with partial participation: tight convergence guarantees. *arXiv preprint arXiv:2006.14591*, 2020.
- Polyak, B. T. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- Reddi, S. J., Hefny, A., Sra, S., Póczos, B., and Smola, A. Stochastic variance reduction for nonconvex optimization. In *International conference on machine learning*, pp. 314–323. PMLR, 2016.
- Robbins, H. and Monro, S. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Rothchild, D., Panda, A., Ullah, E., Ivkin, N., Stoica, I., Braverman, V., Gonzalez, J., and Arora, R. Fetchsgd: Communication-efficient federated learning with sketching. In *International Conference on Machine Learning*, pp. 8253–8265. PMLR, 2020.
- Roux, N., Schmidt, M., and Bach, F. A stochastic gradient method with an exponential convergence rate for finite training sets. *Advances in neural information processing systems*, 25, 2012.
- Ruan, Y., Zhang, X., Liang, S.-C., and Joe-Wong, C. Towards flexible device participation in federated learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 3403–3411. PMLR, 2021.
- Ruder, S. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- Shallue, C. J., Lee, J., Antognini, J., Sohl-Dickstein, J., Frostig, R., and Dahl, G. E. Measuring the effects of data parallelism on neural network training. *Journal of Machine Learning Research*, 20:1–49, 2019.
- Shamsian, A., Navon, A., Fetaya, E., and Chechik, G. Personalized federated learning using hypernetworks. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 9489–9502. PMLR, 2021.
- Stich, S. U. Local sgd converges fast and communicates little. In *International Conference on Learning Representations*, 2018.

- Tewari, A. and Chaudhuri, S. Generalization error bounds for learning to rank: Does the length of document lists matter? In *International Conference on Machine Learning*, pp. 315–323. PMLR, 2015.
- Tyurin, A. and Richtárik, P. Dasha: Distributed nonconvex optimization with communication compression, optimal oracle complexity, and no client synchronization. *arXiv preprint arXiv:2202.01268*, 2022.
- Wang, J., Liu, Q., Liang, H., Joshi, G., and Poor, H. V. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in Neural Information Processing Systems*, 33, 2020.
- Wang, J., Liu, Q., Liang, H., Joshi, G., and Poor, H. V. A novel framework for the analysis and design of heterogeneous federated learning. *IEEE Transactions on Signal Processing*, 69:5234–5249, 2021.
- Wang, W. and Srebro, N. Stochastic nonconvex optimization with large minibatches. In *Algorithmic Learning Theory*, pp. 857–882. PMLR, 2019.
- Wang, Z., Ji, K., Zhou, Y., Liang, Y., and Tarokh, V. Spiderboost: A class of faster variance-reduced algorithms for nonconvex optimization. *arXiv*, 2018, 2018.
- Woodworth, B., Patel, K. K., Stich, S., Dai, Z., Bullins, B., McMahan, B., Shamir, O., and Srebro, N. Is local sgd better than minibatch sgd? In *International Conference on Machine Learning*, pp. 10334–10343. PMLR, 2020a.
- Woodworth, B. E., Patel, K. K., and Srebro, N. Minibatch vs local sgd for heterogeneous distributed learning. *Advances in Neural Information Processing Systems*, 33: 6281–6292, 2020b.
- Wu, F., Guo, S., Wang, H., Qu, Z., Zhang, H., Zhang, J., and Liu, Z. From deterioration to acceleration: A calibration approach to rehabilitating step asynchronism in federated optimization. *arXiv preprint arXiv:2112.09355*, 2021.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Yan, Y., Niu, C., Ding, Y., Zheng, Z., Wu, F., Chen, G., Tang, S., and Wu, Z. Distributed non-convex optimization with sublinear speedup under intermittent client availability. *arXiv preprint arXiv:2002.07399*, 2020.
- Yang, H., Fang, M., and Liu, J. Achieving linear speedup with partial worker participation in non-iid federated learning. In *International Conference on Learning Representations*, 2020.
- You, Y., Zhang, Z., Hsieh, C.-J., Demmel, J., and Keutzer, K. Imagenet training in minutes. In *Proceedings of the 47th International Conference on Parallel Processing*, pp. 1–10, 2018.
- Yu, H., Jin, R., and Yang, S. On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization. In *International Conference on Machine Learning*, pp. 7184–7193. PMLR, 2019a.
- Yu, H., Yang, S., and Zhu, S. Parallel restarted sgd with faster convergence and less communication: Demystifying why model averaging works for deep learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 5693–5700, 2019b.
- Yuan, H. and Ma, T. Federated accelerated stochastic gradient descent. *arXiv preprint arXiv:2006.08950*, 2020.
- Yuan, Z., Guo, Z., Xu, Y., Ying, Y., and Yang, T. Federated deep auc maximization for heterogeneous data with a constant communication complexity. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 12219–12229. PMLR, 2021.
- Yue, K., Jin, R., Pilgrim, R., Wong, C.-W., Baron, D., and Dai, H. Neural tangent kernel empowered federated learning. In *International Conference on Machine Learning*, pp. 25783–25803. PMLR, 2022.
- Yun, C., Rajput, S., and Sra, S. Minibatch vs local sgd with shuffling: Tight convergence bounds and beyond. *arXiv preprint arXiv:2110.10342*, 2021.
- Zhang, H., J Reddi, S., and Sra, S. Riemannian svrg: Fast stochastic optimization on riemannian manifolds. *Advances in Neural Information Processing Systems*, 29, 2016.
- Zhang, J., Guo, S., Ma, X., Wang, H., Xu, W., and Wu, F. Parameterized knowledge transfer for personalized federated learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Zhao, H., Li, Z., and Richtárik, P. Fedpage: A fast local stochastic gradient method for communication-efficient federated learning. *arXiv preprint arXiv:2108.04755*, 2021.
- Zhao, Y., Li, M., Lai, L., Suda, N., Civin, D., and Chandra, V. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- Zhou, D., Xu, P., and Gu, Q. Stochastic nested variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Zinkevich, M., Weimer, M., Li, L., and Smola, A. Parallelized stochastic gradient descent. *Advances in neural information processing systems*, 23, 2010.



## A. Related Work

**Mini-batch SGD vs. Local SGD.** Distributed optimization is required to train large-scale deep learning systems. Local SGD (also known as FedAvg) (Stich, 2018; Dieuleveut & Patel, 2019; Haddadpour et al., 2019; Haddadpour & Mahdavi, 2019) performs multiple (i.e.,  $K \geq 1$ ) local updates with  $K$  small batches, while mini-batch SGD computes the gradients averaged by  $K$  small batches (Woodworth et al., 2020b;a) (or a large batches (Shallue et al., 2019; You et al., 2018; Goyal et al., 2017)) on a given model. There has been a long discussion on which one is better (Lin et al., 2019; Woodworth et al., 2020a;b; Yun et al., 2021), but no existing work considers how to disjoint the nodes such that both can be trained at the same time.

**Federated Learning.** FL was proposed to ensure data privacy and security (Kairouz et al., 2019), and now it has become a hot field in the distributed system (Yuan & Ma, 2020; Shamsian et al., 2021; Zhang et al., 2021; Avdiukhin & Kasiviswanathan, 2021; Yuan et al., 2021; Diao et al., 2020; Blum et al., 2021; He et al., 2022). The FL training methods in the past few years usually require all trainers to participate in each training session (Kairouz et al., 2019), but this is obviously impractical when facing the increase in FL clients. To enhance the systems' feasibility, this work assumes that a fixed number of clients are sampled at each round, which is widely adopted in (Li et al., 2019; Philippenko & Dieuleveut, 2020; Gorbunov et al., 2021; Karimireddy et al., 2020b; Yang et al., 2020; Li et al., 2020; Eichner et al., 2019; Ruan et al., 2021). Therefore, the server collects the data from this participation every synchronization to update the model parameters (Li et al., 2019; Philippenko & Dieuleveut, 2020; Gorbunov et al., 2021; Karimireddy et al., 2020b; Yang et al., 2020; Li et al., 2020; Eichner et al., 2019; Yan et al., 2020; Ruan et al., 2021; Lai et al., 2021; Gu et al., 2021).

**Variance Reduction in Finite-sum Problems.** Variance reduction techniques (Johnson & Zhang, 2013; Defazio et al., 2014; Nguyen et al., 2017; Li et al., 2021; Lan & Zhou, 2018a;b; Allen-Zhu & Hazan, 2016; Reddi et al., 2016; Lei et al., 2017; Zhou et al., 2018; Horváth & Richtárik, 2019; Horváth et al., 2020; Fang et al., 2018; Wang et al., 2018; Li, 2019; Roux et al., 2012; Lian et al., 2017; Zhang et al., 2016) was once proposed for traditional centralized machine learning to optimize finite-sum problems (Bietti & Mairal, 2017; Bottou & Cun, 2003; Robbins & Monro, 1951) by mitigating the estimation gap between small-batch (Bottou, 2012; Ghadimi et al., 2016; Khaled & Richtárik, 2020) and large-batch (Nesterov, 2003; Ruder, 2016; Mason et al., 1999). SGD randomly samples a small-batch and computes the gradient in order to approach the optimal solution. Since the data are generally noisy, an insufficiently large batch results in convergence rate degradation. By utilizing all data in every update, GD can remove the noise affecting the training process. However, it is time-consuming because the period for a single GD step can implement multiple SGD updates. Based on the trade-off, variance-reduced methods periodically perform GD steps while correcting SGD updates with reference to the most recent GD steps.

**Variance Reduction in FL.** The variance reduction techniques have critically driven the advent of FL algorithms (Karimireddy et al., 2020b; Wu et al., 2021; Liang et al., 2019; Karimireddy et al., 2020a; Murata & Suzuki, 2021; Mitra et al., 2021) by correcting each local computed gradient with respect to the estimated global orientation. Roughly, the methods fall into two categories based on types of correction, namely, static and recursive. The former methods (Karimireddy et al., 2020b; Wu et al., 2021; Liang et al., 2019; Karimireddy et al., 2020a; Mitra et al., 2021) use a consistent gradient correction within a training round. As the local updates go further, the static correction on a client is still dramatically biased to its local optimal solution, which hinders the training efficiency, especially under partial client participation. Recursive correction can avoid the challenge because it recursively corrects the local update based on the latest local model (Murata & Suzuki, 2021). However, existing works require all workers to attain the gradient at the starting point of each round, which is infeasible for federated learning settings. To the best of our knowledge, our work is the first to achieve recursive calibration under partial-client scenarios.

## B. Useful Lemmas

Prior to giving detailed proofs of the theorems, we cover some technical lemmas in this section, and all of them are valid in general cases.

**Lemma B.1.** *Let  $\varepsilon = \{\varepsilon_1, \dots, \varepsilon_a\}$  be the set of random variables in  $\mathbb{R}^{a \times d}$ . Every element in  $\varepsilon$  is independent with others. For  $i \in \{1, \dots, a\}$ , the value for  $\varepsilon_i$  follows the setting below:*

$$\varepsilon_i = \begin{cases} e_i, & \text{probability} = q \\ \mathbf{0}, & \text{otherwise} \end{cases} \quad (6)$$

where  $q$  is a constant real number between 0 and 1, i.e.,  $q \in [0, 1]$ . Let  $|\cdot|$  indicate the length of a set,  $\varepsilon \setminus \{\mathbf{0}\}$  represent a set in which an element is in  $\varepsilon$  but not  $\mathbf{0}$ . Then, there is a probability of  $(1 - q)^a$  for  $|\varepsilon \setminus \{\mathbf{0}\}| = 0$ , let  $\text{avg}(\varepsilon)$  be the averaged result with the exception of zero vectors, i.e.,

$$\text{avg}(\varepsilon) = \begin{cases} \frac{1}{|\varepsilon \setminus \{\mathbf{0}\}|} \sum_{i=1}^a \varepsilon_i, & |\varepsilon \setminus \{\mathbf{0}\}| \neq 0 \\ 0, & |\varepsilon \setminus \{\mathbf{0}\}| = 0 \end{cases} \quad (7)$$

Then, the following formulas hold for  $\mathbb{E}(\text{avg}(\varepsilon))$  and its second norm  $\mathbb{E} \|\text{avg}(\varepsilon)\|_2^2$ :

$$\mathbb{E}(\text{avg}(\varepsilon)) = (1 - (1 - q)^a) \cdot \frac{1}{a} \sum_{i=1}^a e_i; \quad \mathbb{E} \|\text{avg}(\varepsilon)\|_2^2 \leq (1 - (1 - q)^a) \cdot \frac{1}{a} \sum_{i=1}^a \|e_i\|_2^2 \quad (8)$$

*Proof.* When  $q = 0$ , the formulas in Equation 8 obviously hold because  $\mathbb{E}(\text{avg}(\varepsilon)) = 0$  and  $\mathbb{E} \|\text{avg}(\varepsilon)\|_2^2 = 0$ . As for  $q = 1$ , since  $\text{avg}(\varepsilon) = \frac{1}{a} \sum_{i=1}^a e_i$ , we leverage Cauchy–Schwarz inequality and get  $\mathbb{E} \|\text{avg}(\varepsilon)\|_2^2 = \|\frac{1}{a} \sum_{i=1}^a e_i\|_2^2 \leq \frac{1}{a} \sum_{i=1}^a \|e_i\|_2^2$ , which is consistent with the formulas in Equation 8. In addition to the preceding cases, we consider some general cases for the probability  $q$  within 0 and 1, i.e.,  $q \in (0, 1)$ .

Firstly, we show the proof details for  $\mathbb{E}(\text{avg}(\varepsilon))$ . For all  $i$  in  $\{1, \dots, a\}$ , given that  $\varepsilon_i$  is not a zero vector, the coefficient of  $e_i$  is based on the binomial distribution on how many non-zero elements in the set  $\{\varepsilon_1, \dots, \varepsilon_{i-1}\} \cup \{\varepsilon_{i+1}, \dots, \varepsilon_a\}$ . Therefore, with the probability  $q$  that  $\varepsilon_i$  is equal to  $e_i$ , the coefficient of  $e_i$  in the expected form is

$$q \left( \underbrace{\frac{1}{a} \cdot \binom{a-1}{a-1} q^{a-1}}_{(a-1) \text{ non-zero elements}} + \dots + \frac{1}{1} \cdot \underbrace{\binom{a-1}{0} (1-q)^{a-1}}_{0 \text{ non-zero element}} \right)$$

Then, the coefficient of  $\frac{1}{a} e_i$  can be expressed and simplified for

$$q \left( \frac{a}{a} \cdot \binom{a-1}{a-1} q^{a-1} + \dots + \frac{a}{1} \cdot \binom{a-1}{0} (1-q)^{a-1} \right) \quad (9)$$

$$= q \left( \binom{a}{a} q^{a-1} + \dots + \binom{a}{1} (1-q)^{a-1} \right) \quad (10)$$

$$= \binom{a}{a} q^a + \dots + \binom{a}{1} q (1-q)^{a-1} \quad (11)$$

$$= 1 - (1-q)^a \quad (12)$$

where Equation (11) follows

$$\binom{\alpha}{\beta} = \frac{\alpha}{\beta} \cdot \frac{(\alpha-1) \times \dots \times (\alpha-\beta+1)}{1 \times \dots \times (\beta-1)} = \frac{\alpha}{\beta} \binom{\alpha-1}{\beta-1}, \quad \forall \alpha \geq \beta > 0$$

and Equation (12) follows

$$(q + (1-q))^a = \binom{a}{a} q^a + \dots + \binom{a}{0} (1-q)^a.$$

Thus, the equation  $\mathbb{E}(\text{avg}(\varepsilon)) = (1 - (1 - q)^a) \cdot \frac{1}{a} \sum_{i=1}^a e_i$  holds.

Secondly, we provide the analysis for  $\mathbb{E} \|\text{avg}(\varepsilon)\|_2^2$ . Based on the definition for  $\text{avg}(\varepsilon)$  in Equation (7), we discuss the case  $|\varepsilon \setminus \{\mathbf{0}\}| \neq 0$ . By means of Cauchy-Schwarz inequality, we can obtain the following inequality:

$$\left\| \frac{1}{|\varepsilon \setminus \{\mathbf{0}\}|} \sum_{i=1}^a \varepsilon_i \right\|_2^2 = \left\| \frac{1}{|\varepsilon \setminus \{\mathbf{0}\}|} \sum_{i, \varepsilon_i \neq \mathbf{0}} \varepsilon_i \right\|_2^2 \leq \frac{1}{|\varepsilon \setminus \{\mathbf{0}\}|} \sum_{i, \varepsilon_i \neq \mathbf{0}} \|\varepsilon_i\|_2^2 = \frac{1}{|\varepsilon \setminus \{\mathbf{0}\}|} \sum_{i=1}^a \|\varepsilon_i\|_2^2 \quad (13)$$

Therefore,

$$\|\text{avg}(\varepsilon)\|_2^2 \leq \begin{cases} \frac{1}{|\varepsilon \setminus \{\mathbf{0}\}|} \sum_{i=1}^a \|\varepsilon_i\|_2^2, & |\varepsilon \setminus \{\mathbf{0}\}| \neq 0 \\ 0, & |\varepsilon \setminus \{\mathbf{0}\}| = 0 \end{cases} \quad (14)$$

Apparently, Equation (14) is very similar to Equation (7) in terms of the expression. As a result, we can adopt the same proof framework in the analysis of  $\mathbb{E}(\text{avg}(\varepsilon))$ . Then, we can directly draw a conclusion  $\mathbb{E} \|\text{avg}(\varepsilon)\|_2^2 \leq (1 - (1 - q)^a) \cdot \frac{1}{a} \sum_{i=1}^a \|e_i\|_2^2$ .  $\square$

**Lemma B.2.** Let  $\varepsilon = \{\varepsilon_1, \dots, \varepsilon_a\}$  be the set of random variables in  $\mathbb{R}^d$  with the number of  $a$ . These random variables are not necessarily independent. We can suppose that  $\mathbb{E}[\varepsilon_i] = e_i$ , and the variance is bounded as  $\mathbb{E}[\|\varepsilon_i - e_i\|_2^2] \leq \sigma^2$ . After that we can get

$$\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i \right\|_2^2 \right] \leq \left\| \sum_{i=1}^a e_i \right\|_2^2 + a^2 \sigma^2 \quad (15)$$

If we make another suppose that the conditional mean of these random variables is  $\mathbb{E}[\varepsilon_i | \varepsilon_{i-1}, \dots, \varepsilon_1] = e_i$ , and the variables  $\{\varepsilon_i - e_i\}$  form a martingale difference sequence, and the bound of the variance is  $\mathbb{E}[\|\varepsilon_i - e_i\|_2^2] \leq \sigma^2$ . So we can make a much tighter bound

$$\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i \right\|_2^2 \right] \leq 2 \left\| \sum_{i=1}^a e_i \right\|_2^2 + 2a\sigma^2 \quad (16)$$

*Proof.* For any random variable  $X$ ,  $\mathbb{E}[X^2] = (\mathbb{E}[X - \mathbb{E}[X]])^2 + (\mathbb{E}[X])^2$  implying

$$\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i \right\|_2^2 \right] = \left\| \sum_{i=1}^a e_i \right\|_2^2 + \mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i - e_i \right\|_2^2 \right] \quad (17)$$

Expanding above expression using relaxed triangle inequality:

$$\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i - e_i \right\|_2^2 \right] \leq a \sum_{i=1}^a \mathbb{E}[\|\varepsilon_i - e_i\|_2^2] \leq a^2 \sigma^2 \quad (18)$$

For the second statement,  $e_i$  depends on  $[\varepsilon_{i-1}, \dots, \varepsilon_1]$ . Thus we choose to use a relaxed triangle inequality

$$\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i \right\|_2^2 \right] \leq 2 \left\| \sum_{i=1}^a e_i \right\|_2^2 + 2\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i - e_i \right\|_2^2 \right] \quad (19)$$

then we use a much tighter expansion and we can get:

$$\mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i - e_i \right\|_2^2 \right] = \sum_{i,j} \mathbb{E}[(\varepsilon_i - e_i)^\top (\varepsilon_j - e_j)] = \sum_i \mathbb{E} \left[ \left\| \sum_{i=1}^a \varepsilon_i - e_i \right\|_2^2 \right] \leq a\sigma^2 \quad (20)$$

When  $\{\varepsilon_i - e_i\}$  form a martingale difference sequence, the cross terms will have zero means.  $\square$

**Lemma B.3.** Suppose there is a sequence  $\{y_t \in \mathbb{R}^d\}_{t \geq 0}$  satisfying a recursive function  $y_{t+1} = y_t - \eta \Delta y_t$ , where  $\eta > 0$  is a constant and  $\Delta y_t \in \mathbb{R}^d$  is a vector. Given a  $L$ -smooth function  $G$ , the following inequality holds for any  $\eta$  and  $\Delta y_t$ :

$$G(y_{t+1}) \leq G(y_t) - \frac{\eta \eta'}{2} \|\nabla G(y_t)\|_2^2 - \left( \frac{1}{2\eta \eta'} - \frac{L}{2} \right) \|y_{t+1} - y_t\|_2^2 + \frac{\eta}{2\eta'} \|\Delta y_t - \eta' \nabla G(y_t)\|_2^2 \quad (21)$$

where  $\eta' > 0$  can be any constant.

*Proof.* Since  $G$  is a  $L$ -smooth function, for any  $v, \bar{v} \in \mathbb{R}^d$ , the following inequality holds:

$$G(\bar{v}) = G(v) + \int_0^1 \frac{\partial G(v + t(\bar{v} - v))}{\partial t} dt \quad (22)$$

$$= G(v) + \int_0^1 \nabla G(v + t(\bar{v} - v)) \cdot (\bar{v} - v) dt \quad (23)$$

$$= G(v) + \nabla G(v)(\bar{v} - v) + \int_0^1 (\nabla G(v + t(\bar{v} - v)) - \nabla G(v)) \cdot (\bar{v} - v) dt \quad (24)$$

$$\leq G(v) + \nabla G(v)(\bar{v} - v) + \int_0^1 L \|t(\bar{v} - v)\|_2 \|\bar{v} - v\|_2 dt \quad (25)$$

$$\leq G(v) + \nabla G(v)(\bar{v} - v) + \frac{L}{2} \|\bar{v} - v\|_2^2. \quad (26)$$

Based on the conclusion on  $L$ -smooth drawn from Equation (26), we derive Equation (21) step by step:

$$G(y_{t+1}) \leq G(y_t) + \langle \nabla G(y_t), y_{t+1} - y_t \rangle + \frac{L}{2} \|y_{t+1} - y_t\|_2^2 \quad (27)$$

$$= G(y_t) + \langle \nabla G(y_t), -\eta \Delta y_t \rangle + \frac{L}{2} \|y_{t+1} - y_t\|_2^2 \quad (28)$$

$$= G(y_t) - \frac{\eta}{\eta'} \langle \eta' \nabla G(y_t), \Delta y_t \rangle + \frac{L}{2} \|y_{t+1} - y_t\|_2^2 \quad (29)$$

$$= G(y_t) - \frac{\eta}{2\eta'} \left( \eta'^2 \|\nabla G(y_t)\|_2^2 + \|\Delta y_t\|_2^2 - \|\Delta y_t - \eta' \nabla G(y_t)\|_2^2 \right) + \frac{L}{2} \|y_{t+1} - y_t\|_2^2 \quad (30)$$

$$= G(y_t) - \frac{\eta \eta'}{2} \|\nabla G(y_t)\|_2^2 - \left( \frac{1}{2\eta \eta'} - \frac{L}{2} \right) \|y_{t+1} - y_t\|_2^2 + \frac{\eta}{2\eta'} \|\Delta y_t - \eta' \nabla G(y_t)\|_2^2 \quad (31)$$

where Equation (30) is in accordance with  $\langle \alpha, \beta \rangle = \frac{1}{2} (\alpha^2 + \beta^2 - (\alpha - \beta)^2)$ , and Equation (31) follows  $\|\Delta y_t\|_2^2 = \frac{1}{\eta^2} \|y_{t+1} - y_t\|_2^2$ .  $\square$



## C. Preliminary for FedAMD

Algorithm 1 describes FedAMD in details. The objective in this part is to find the recursive function for the sequence of models, i.e.,  $\{\tilde{\mathbf{x}}_t\}_{t \geq 0}$ . As mentioned in Line 26 in Algorithm 1, let  $\Delta \mathbf{x}_t$  aggregate  $\Delta \mathbf{x}_t^{(i)}$  where client  $i$  updates model, then the difference between  $\tilde{\mathbf{x}}_{t+1}$  and  $\tilde{\mathbf{x}}_t$  follows the recursive function written as

$$\tilde{\mathbf{x}}_{t+1} = \tilde{\mathbf{x}}_t - \eta_s \cdot \text{avg}(\Delta \mathbf{x}_t) \quad (32)$$

where  $\text{avg}()$  is same as defined in Lemma B.1. As we know, the length of  $\Delta \mathbf{x}_t$  changes over rounds but does not exceed the number of participants, i.e.,  $|\Delta \mathbf{x}_t| \leq A$ . Then, suppose that  $\Delta \mathbf{x}_t^{(m)}$  is in  $\Delta \mathbf{x}_t$ ,  $\Delta \mathbf{x}_t^{(m)}$  can be expressed as

$$\Delta \mathbf{x}_t^{(m)} = - \left( \mathbf{x}_{t,K}^{(m)} - \tilde{\mathbf{x}}_t \right) = - \sum_{k=0}^{K-1} \left( \mathbf{x}_{t,k+1}^{(m)} - \mathbf{x}_{t,k+1}^{(m)} \right) = \sum_{k=0}^{K-1} \eta_l g_{t,k+1}^{(m)} \quad (33)$$

where the last equal sign is according to Line 20 in Algorithm 1. Next, with the recursive formula in Line 19, we have

$$g_{t,k+1}^{(m)} = g_{t,k}^{(m)} - \nabla f_m \left( \mathbf{x}_{t,k-1}^{(m)}, \mathcal{B}'_{m,k} \right) + \nabla f_m \left( \mathbf{x}_{t,k}^{(m)}, \mathcal{B}'_{m,k} \right) \quad (34)$$

$$= \tilde{g}_t - \sum_{\kappa=0}^k \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) + \sum_{\kappa=0}^k \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right). \quad (35)$$

Then, Equation (33) can be rewritten as

$$\Delta \mathbf{x}_t^{(m)} = \eta_l K \tilde{g}_t - \eta_l \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) + \eta_l \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) \quad (36)$$

## D. Proofs under Sequential Probabilistic Settings

### D.1. Preliminary

**Lemma D.1.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local learning rate satisfy  $\eta_l \leq \min \left( \frac{1}{2\sqrt{3}KL}, \frac{1}{2\sqrt{3}L\sigma} \sqrt{\frac{b'}{K}} \right)$ . With FedAMD,  $\sum_{k=0}^{K-1} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2$  represents the sum of the second norm of every iteration's difference. Therefore, the bound for such a summation in the expected form should be*

$$\sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \leq 6\eta_l^2 K \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 6\eta_l^2 K \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (37)$$

*Proof.* According to Equation (35), the update at  $(k-1)$ -th iteration is

$$\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} = -\eta_l g_{t,k}^{(m)} = -\eta_l \left( \tilde{g}_t - \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) + \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) \right). \quad (38)$$

To find the bound for the expected value of its second norm, the analysis is presented as follows:

$$\mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (39)$$

$$= \eta_l^2 \mathbb{E} \left\| \tilde{g}_t - \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) + \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) \right\|_2^2 \quad (40)$$

$$\leq 3\eta_l^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \\ + 3\eta_l^2 \mathbb{E} \left\| \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) - \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) \right\|_2^2 \quad (41)$$

$$\begin{aligned}
 &= 3\eta_l^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 \mathbb{E} \left\| \sum_{\kappa=0}^{k-1} \left( \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) - \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) \right) \right\|_2^2 \\
 &\quad + 3\eta_l^2 \mathbb{E} \left\| \sum_{\kappa=0}^{k-1} \left( \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) + \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) \right) \right\|_2^2 \quad (42)
 \end{aligned}$$

$$\begin{aligned}
 &\leq 3\eta_l^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 KL^2 \sum_{\kappa=0}^{k-1} \mathbb{E} \|\mathbf{x}_{t,\kappa-1}^{(m)} - \mathbf{x}_{t,\kappa}^{(m)}\|_2^2 \\
 &\quad + 3\eta_l^2 \mathbb{E} \left\| \sum_{\kappa=0}^{k-1} \left( \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) + \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) \right) \right\|_2^2 \quad (43)
 \end{aligned}$$

$$\begin{aligned}
 &\leq 3\eta_l^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 KL^2 \sum_{\kappa=0}^{k-1} \mathbb{E} \|\mathbf{x}_{t,\kappa-1}^{(m)} - \mathbf{x}_{t,\kappa}^{(m)}\|_2^2 \\
 &\quad + 3\eta_l^2 \frac{L_\sigma^2}{b'} \sum_{\kappa=0}^{k-1} \mathbb{E} \|\mathbf{x}_{t,\kappa-1}^{(m)} - \mathbf{x}_{t,\kappa}^{(m)}\|_2^2 \quad (44)
 \end{aligned}$$

where Equation (41) is based on Cauchy-Schwarz inequality; Equation (42) is based on the variance expansion on the third term of Equation (41); Equation (43) is based on Cauchy-Schwarz inequality and Assumption 4.1 on the third term of Equation (42); Equation (44) is based on Lemma B.2 and Assumption 4.3 on the fourth term of Equation (43).

Therefore, by summing Equation (44) for  $k = 1, \dots, K$ , we have

$$\sum_{k=0}^{K-1} \|\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)}\|_2^2 \leq \sum_{k=0}^{K-1} \|\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)}\|_2^2 \quad (45)$$

$$\leq 3\eta_l^2 K \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 K \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_l^2 K \left( KL^2 + \frac{L_\sigma^2}{b'} \right) \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{x}_{t,\kappa-1}^{(m)} - \mathbf{x}_{t,\kappa}^{(m)}\|_2^2 \quad (46)$$

Obviously, according to the setting of the local learning rate in the description above, the inequality  $3\eta_l^2 K \left( KL^2 + \frac{L_\sigma^2}{b'} \right) \leq \frac{1}{2}$  holds. Therefore, we can easily obtain the bound for the sum of the second norm of every iteration's difference, which is consistent with Equation (37).  $\square$

## D.2. Full Client Participation

**Theorem D.2.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local updates  $K \geq 1$ , the setting for the local learning rate  $\eta_l$  and the global learning rate  $\eta_s$  satisfy the following two constraints: (1)  $\eta_s \eta_l \leq \frac{1}{6\sqrt{6}KL\tau}$ ; and (2)*

*$\eta_l \leq \min \left( \frac{1}{6\sqrt{2}KL}, \frac{\sqrt{b'/K}}{2\sqrt{3}L_\sigma} \right)$ . Then, the convergence rate of FedAMD for non-convex objectives should be*

$$\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \frac{4(F(\tilde{\mathbf{x}}_0) - F_*)}{\eta_s \eta_l K T} + 16\eta_l L K (2\eta_s + 3\eta_l K L) \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{M b} \quad (47)$$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

*Proof.* When  $p_t = 1$ , according to Algorithm 1, there is no model update between two consecutive rounds, i.e.,  $\tilde{\mathbf{x}}_{t+1} = \tilde{\mathbf{x}}_t$ .

Next, we consider the case when  $p_t = 0$ . According to Assumption 4.1, we have

$$\mathbb{E} F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \leq \mathbb{E} \langle \nabla F(\tilde{\mathbf{x}}_t), \tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t \rangle + \frac{L}{2} \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \quad (48)$$

Knowing that when  $p_t = 0$  and all clients involve in the training,

$$\text{avg}(\Delta \mathbf{x}_t) = \eta_l K \tilde{g}_t - \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) + \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}), \quad (49)$$

we have the bound for the first term of Equation (48):

$$\mathbb{E} \langle \nabla F(\tilde{\mathbf{x}}_t), \tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t \rangle \quad (50)$$

$$= -\eta_s \eta_l K \mathbb{E} \left\langle \nabla F(\tilde{\mathbf{x}}_t), \mathbb{E} \tilde{g}_t - \frac{1}{MK} \sum_{m \in [M]} \sum_{k=0}^{K-1} \left( \nabla F_m(\mathbf{x}_{t,k}^{(m)}) - \nabla F_m(\tilde{\mathbf{x}}_t) \right) \right\rangle \quad (51)$$

$$\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{\eta_s \eta_l K}{2} \left\| \left( \nabla F(\tilde{\mathbf{x}}_t) - \mathbb{E} \tilde{g}_t \right) - \frac{1}{MK} \sum_{m \in [M]} \sum_{k=0}^{K-1} \left( \nabla F_m(\mathbf{x}_{t,k}^{(m)}) - \nabla F_m(\tilde{\mathbf{x}}_t) \right) \right\|_2^2 \quad (52)$$

$$\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \eta_s \eta_l K \|\nabla F(\tilde{\mathbf{x}}_t) - \mathbb{E} \tilde{g}_t\|_2^2 + \eta_s \eta_l K \cdot \frac{1}{MK} \sum_{m \in [M]} \sum_{k=0}^{K-1} \left\| \nabla F_m(\mathbf{x}_{t,k}^{(m)}) - \nabla F_m(\tilde{\mathbf{x}}_t) \right\|_2^2 \quad (53)$$

$$\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \eta_s \eta_l K \|\nabla F(\tilde{\mathbf{x}}_t) - \mathbb{E} \tilde{g}_t\|_2^2 + \frac{\eta_s \eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} L^2 \left\| \mathbf{x}_{t,k}^{(m)} - \tilde{\mathbf{x}}_t \right\|_2^2 \quad (54)$$

$$\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \eta_s \eta_l K \|\nabla F(\tilde{\mathbf{x}}_t) - \mathbb{E} \tilde{g}_t\|_2^2 + \frac{\eta_s \eta_l L^2}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^{k-1} \left\| \mathbf{x}_{t,\kappa+1}^{(m)} - \mathbf{x}_{t,\kappa}^{(m)} \right\|_2^2 \quad (55)$$

$$\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \eta_s \eta_l K \|\nabla F(\tilde{\mathbf{x}}_t) - \mathbb{E} \tilde{g}_t\|_2^2 + \eta_s \eta_l L^2 K \left( 6\eta_l^2 K \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 6\eta_l^2 K \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \quad (56)$$

where Equation (52) follows  $-\langle \alpha, \beta \rangle = -\frac{1}{2}(\alpha^2 + \beta^2 - (\alpha - \beta)^2) \leq -\frac{1}{2}\alpha^2 + \frac{1}{2}(\alpha - \beta)^2$ , Equation (53) applies Cauchy-Schwarz Inequality, and (56) uses Lemma D.1.

To bound the second term of Equation (48), we start with

$$\mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 = \mathbb{E} \|\eta_s \cdot \text{avg}(\Delta \mathbf{x}_t)\|_2^2 \quad (57)$$

$$= \mathbb{E} \|\eta_s \cdot \text{avg}(\Delta \mathbf{x}_t) - \eta_s \eta_l K \nabla F(\tilde{\mathbf{x}}_t) + \eta_s \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (58)$$

$$\leq 2\eta_s^2 \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\eta_s^2 \eta_l^2 K^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (59)$$

Then, we have the bound for  $\mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2$  according to the following derivation:

$$\mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (60)$$

$$= \mathbb{E} \left\| \eta_l K (\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)) + \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) \right) \right\|_2^2 \quad (61)$$

$$\leq 2\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\mathbb{E} \left\| \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) \right) \right\|_2^2 \quad (62)$$

$$= 2\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\mathbb{E} \left\| \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) - \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) \right) \right\|_2^2 + 2\mathbb{E} \left\| \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) - \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) + \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) \right) \right\|_2^2 \quad (63)$$

$$= 2\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K L^2}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k k \cdot \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2$$

$$+ 2\mathbb{E} \left\| \frac{\eta_l}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) - \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) - \nabla F_m \left( \mathbf{x}_{t,\kappa}^{(m)} \right) + \nabla F_m \left( \mathbf{x}_{t,\kappa-1}^{(m)} \right) \right) \right\|_2^2 \quad (64)$$

$$\leq 2\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K L^2}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k k \cdot \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \\ + \frac{2\eta_l^2}{M^2} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \frac{L_\sigma^2}{b'} \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \quad (65)$$

$$\leq 2\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K^3 L^2}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \\ + \frac{2\eta_l^2 K L_\sigma^2}{M^2 b'} \sum_{m \in [M]} \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (66)$$

$$\leq 2\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 12\eta_l^4 K^2 \left( \frac{L_\sigma^2}{M b'} + K^2 L^2 \right) \left( \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \quad (67)$$

where Equation (62) follows  $(\alpha + \beta)^2 \leq 2\alpha^2 + 2\beta^2$ ; Equation (63) is based on variance expansion; Equation (64) is based on Cauchy-Schwarz inequality and Assumption 4.1; Equation (65) is based on Lemma B.2 and Assumption 4.3; Equation (67) is based on Lemma D.1. According to the constraints on the local learning rate, we can further simplify Equation (67) as

$$\mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq 4\eta_l^2 K^2 \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{\eta_l^2 K^2}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2. \quad (68)$$

Plugging Equation (56), Equation (59), and Equation (68) into Equation (48), we reorganize the formula and obtain

$$\mathbb{E} F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \leq -\frac{\eta_s \eta_l K}{4} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \eta_s \eta_l K \|\nabla F(\tilde{\mathbf{x}}_t) - \mathbb{E} \tilde{g}_t\|_2^2 + \eta_s \eta_l^2 K^2 L (6\eta_l L + 4\eta_s) \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (69)$$

Let  $\Lambda(t)$  indicate the most recent round where  $p_{\Lambda(t)} = 1$  and  $\Lambda(t) \neq t$ . It is noted that recursively using  $\Lambda(\cdot)$  can achieve the value of 0, i.e.,  $\underbrace{\Lambda(\Lambda(\dots\Lambda(t)))}_{\text{multiple } \Lambda} = 0$ . As we know,

$$\mathbb{E} \|\tilde{g}_\theta - \nabla F(\tilde{\mathbf{x}}_\theta)\|_2^2 = \mathbb{E} \|\tilde{g}_{\Lambda(\theta)} - \nabla F(\tilde{\mathbf{x}}_\theta)\|_2^2 \quad (70)$$

$$= \mathbb{E} \|\tilde{g}_{\Lambda(\theta)} - \nabla F(\tilde{\mathbf{x}}_{\Lambda(\theta)})\|_2^2 + \mathbb{E} \|\nabla F(\tilde{\mathbf{x}}_{\Lambda(\theta)}) - \nabla F(\tilde{\mathbf{x}}_\theta)\|_2^2 \quad (71)$$

$$\leq \mathbb{E} \|\tilde{g}_{\Lambda(\theta)} - \nabla F(\tilde{\mathbf{x}}_{\Lambda(\theta)})\|_2^2 + L^2 \mathbb{E} \|\tilde{\mathbf{x}}_\theta - \tilde{\mathbf{x}}_{\Lambda(\theta)}\|_2^2 \quad (72)$$

$$\leq \mathbb{E} \|\tilde{g}_{\Lambda(\theta)} - \nabla F(\tilde{\mathbf{x}}_{\Lambda(\theta)})\|_2^2 + L^2 \tau \sum_{\Xi=\Lambda(\theta)}^{\theta-1} \mathbb{E} \|\tilde{\mathbf{x}}_{\Xi+1} - \tilde{\mathbf{x}}_\Xi\|_2^2 \quad (73)$$

$$\leq \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{M b} + L^2 \tau \sum_{\Xi=\Lambda(\theta)}^{\theta-1} \mathbb{E} \|\tilde{\mathbf{x}}_{\Xi+1} - \tilde{\mathbf{x}}_\Xi\|_2^2 \quad (74)$$

where Equation (71) is based on the variance expansion; Equation (72) is based on Assumption 4.1; Equation (73) is according to Cauchy-Schwarz inequality and  $\theta - \Lambda(\theta) \leq \tau$ ; Equation (74) follows Assumption 4.2. Therefore, Equation (69) can be further simplified as

$$\mathbb{E} F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \leq -\frac{\eta_s \eta_l K}{3} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s \eta_l K L^2 \tau \sum_{\Xi=\Lambda(\theta)}^{\theta-1} \mathbb{E} \|\tilde{\mathbf{x}}_{\Xi+1} - \tilde{\mathbf{x}}_\Xi\|_2^2 + \eta_s \eta_l^2 K^2 L (6\eta_l L + 4\eta_s) \cdot \frac{\sigma^2}{M b} \quad (75)$$



By summing  $\mathbb{E} \|\tilde{\mathbf{x}}_{\theta+1} - \tilde{\mathbf{x}}_{\theta}\|_2^2$  from  $\theta = \Lambda(t)$  to  $t-1$ , we can easily obtain the following bound for all  $t \in \{\Lambda(t), \dots, \Lambda(t) + \tau - 1\}$

$$\sum_{\theta=\Lambda(t)}^{t-1} \mathbb{E} \|\tilde{\mathbf{x}}_{\theta+1} - \tilde{\mathbf{x}}_{\theta}\|_2^2 \leq 16\eta_s^2 \eta_l^2 K^2 \tau \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + 6\eta_s^2 \eta_l^2 K^2 \sum_{\theta=\Lambda(t)}^{t-1} \|\nabla F(\tilde{\mathbf{x}}_{\theta})\|_2^2 \quad (76)$$

By summing Equation (75) from  $\Lambda(t)$  to  $t-1$  and applying the bound in Equation (76), we have

$$\mathbb{E}F(\tilde{\mathbf{x}}_t) - F(\tilde{\mathbf{x}}_{\Lambda(t)}) = \sum_{\theta=\Lambda(t)}^{t-1} (\mathbb{E}F(\tilde{\mathbf{x}}_{\theta+1}) - F(\tilde{\mathbf{x}}_{\theta})) \quad (77)$$

$$\leq -\frac{\eta_s \eta_l K}{4} \sum_{\theta=\Lambda(t)}^{t-1} \|\nabla F(\tilde{\mathbf{x}}_{\theta})\|_2^2 + \eta_s \eta_l^2 L K^2 ((4\eta_s + 6\eta_l K L)(t - \Lambda(t)) + 48\eta_s^2 \eta_l^2 K \tau^2 L) \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb}. \quad (78)$$

Therefore, based on the equation above, by summing up all  $t \in \{T, \Lambda(T), \dots, 0\}$ , we can obtain the following inequality:

$$F_* - F(\tilde{\mathbf{x}}_0) \leq \mathbb{E}F(\tilde{\mathbf{x}}_T) - F(\tilde{\mathbf{x}}_0) \leq -\frac{\eta_s \eta_l K}{4} \sum_{t=0}^{T-1} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 4\eta_s \eta_l^2 L K^2 T (2\eta_s + 3\eta_l K L) \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (79)$$

Thus, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \frac{4(F(\tilde{\mathbf{x}}_0) - F_*)}{\eta_s \eta_l K T} + 16\eta_l L K (2\eta_s + 3\eta_l K L) \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (80)$$

By using the settings of the local learning rate and the global learning rate in the description, we can obtain the desired result.  $\square$

### D.3. Partial Client Participation

**Theorem D.3.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local updates  $K \geq 1$ , and the local learning rate  $\eta_l$  and the global learning rate  $\eta_s$  be  $\eta_s \eta_l = \frac{1}{KL} \left(1 + \frac{2M\tau}{A}\right)^{-1}$ , where  $\eta_l \leq \min\left(\frac{1}{2\sqrt{6}KL}, \frac{\sqrt{b'/K}}{4\sqrt{3}L\sigma}\right)$ . Therefore, the convergence rate of FedAMD for non-convex objectives should be*

$$\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq O\left(\frac{1}{T - \lfloor T/\tau \rfloor} \left(1 + \frac{2M\tau}{A}\right)\right) + O\left(\mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb}\right) \quad (81)$$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

*Proof.* When  $p_t = 1$ , according to Algorithm 1, there is no model update between two consecutive rounds, i.e.,  $\tilde{\mathbf{x}}_{t+1} = \tilde{\mathbf{x}}_t$ .

Next, we consider the case when  $p_t = 0$ . Based on Lemma B.3, we have

$$\begin{aligned} \mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) &\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left(\frac{1}{2\eta_s \eta_l K} - \frac{L}{2}\right) \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \\ &\quad + \frac{\eta_s}{2\eta_l K} \mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \end{aligned} \quad (82)$$

Knowing that when  $p_t = 0$  and a set of clients  $\mathcal{A}$  involve in the training,

$$\text{avg}(\Delta \mathbf{x}_t) = \eta_l K \tilde{\mathbf{g}}_t - \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \nabla f_m(\mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa}) + \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \nabla f_m(\mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa}), \quad (83)$$

we have the bound for  $\mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2$  according to the following derivation:

$$\mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (84)$$

$$= \mathbb{E} \left\| \eta_l K (\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)) + \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_i(\mathbf{x}_{t,\kappa}^{(i)}, \mathcal{B}'_{i,\kappa}) - \nabla f_i(\mathbf{x}_{t,\kappa-1}^{(i)}, \mathcal{B}'_{i,\kappa}) \right) \right\|_2^2 \quad (85)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\mathbb{E} \left\| \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_i(\mathbf{x}_{t,\kappa}^{(i)}, \mathcal{B}'_{i,\kappa}) - \nabla f_i(\mathbf{x}_{t,\kappa-1}^{(i)}, \mathcal{B}'_{i,\kappa}) \right) \right\|_2^2 \quad (86)$$

$$= 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\mathbb{E} \left\| \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla F_i(\mathbf{x}_{t,\kappa}^{(i)}) - \nabla F_i(\mathbf{x}_{t,\kappa-1}^{(i)}) \right) \right\|_2^2 \\ + 2\mathbb{E} \left\| \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_i(\mathbf{x}_{t,\kappa}^{(i)}, \mathcal{B}'_{i,\kappa}) - \nabla f_i(\mathbf{x}_{t,\kappa-1}^{(i)}, \mathcal{B}'_{i,\kappa}) - \nabla F_i(\mathbf{x}_{t,\kappa}^{(i)}) + \nabla F_i(\mathbf{x}_{t,\kappa-1}^{(i)}) \right) \right\|_2^2 \quad (87)$$

$$= 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K L^2}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k k \cdot \mathbb{E} \|\mathbf{x}_{t,\kappa}^{(i)} - \mathbf{x}_{t,\kappa-1}^{(i)}\|_2^2 \\ + 2\mathbb{E} \left\| \frac{\eta_l}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla f_i(\mathbf{x}_{t,\kappa}^{(i)}, \mathcal{B}'_{i,\kappa}) - \nabla f_i(\mathbf{x}_{t,\kappa-1}^{(i)}, \mathcal{B}'_{i,\kappa}) - \nabla F_i(\mathbf{x}_{t,\kappa}^{(i)}) + \nabla F_i(\mathbf{x}_{t,\kappa-1}^{(i)}) \right) \right\|_2^2 \quad (88)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K L^2}{A} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k k \cdot \mathbb{E} \|\mathbf{x}_{t,\kappa}^{(i)} - \mathbf{x}_{t,\kappa-1}^{(i)}\|_2^2 \\ + \frac{2\eta_l^2}{A^2} \sum_{i \in \mathcal{A}} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \frac{L_\sigma^2}{b'} \mathbb{E} \|\mathbf{x}_{t,\kappa}^{(i)} - \mathbf{x}_{t,\kappa-1}^{(i)}\|_2^2 \quad (89)$$

$$= 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K L^2}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k k \cdot \mathbb{E} \|\mathbf{x}_{t,\kappa}^{(i)} - \mathbf{x}_{t,\kappa-1}^{(i)}\|_2^2 \\ + \frac{2\eta_l^2}{AM} \sum_{m \in [M]} \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \frac{L_\sigma^2}{b'} \mathbb{E} \|\mathbf{x}_{t,\kappa}^{(i)} - \mathbf{x}_{t,\kappa-1}^{(i)}\|_2^2 \quad (90)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{2\eta_l^2 K^3 L^2}{M} \sum_{m \in [M]} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)}\|_2^2 \\ + \frac{2\eta_l^2 K L_\sigma^2}{AM b'} \sum_{m \in [M]} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)}\|_2^2 \quad (91)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 12\eta_l^4 K^2 \left( \frac{L_\sigma^2}{Ab'} + K^2 L^2 \right) \left( \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \quad (92)$$

where Equation (86) follows  $(\alpha + \beta)^2 \leq 2\alpha^2 + 2\beta^2$ ; Equation (87) is based on variance expansion; Equation (88) is based on Cauchy-Schwarz inequality and Assumption 4.1; Equation (89) is based on Lemma B.2 and Assumption 4.3; Equation (90) is based on the setting of client selection, where each client is selected with a probability of  $A/M$ ; Equation (92) is based on Lemma D.1. According to the constraints on the local learning rate, we can further simplify Equation (92) as

$$\mathbb{E} \|\text{avg}(\Delta \mathbf{x}_t) - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq 4\eta_l^2 K^2 \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{\eta_l^2 K^2}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2. \quad (93)$$

Plugging Equation (93) into Equation (82), we have

$$\mathbb{E} F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \quad (94)$$

$$\leq -\frac{\eta_s \eta_l K}{4} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} \right) \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 + 2\eta_s \eta_l K \mathbb{E} \|\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (95)$$

$$= -\frac{\eta_s \eta_l K}{4} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} \right) \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2$$

$$+ 2\eta_s\eta_l K \left( \mathbb{E} \|\tilde{g}_t - \mathbb{E}\tilde{g}_t\|_2^2 + \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \quad (96)$$

$$\leq -\frac{\eta_s\eta_l K}{4} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left( \frac{1}{2\eta_s\eta_l K} - \frac{L}{2} \right) \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 + 2\eta_s\eta_l K \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \\ + 2\eta_s\eta_l K \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (97)$$

where Equation (96) is based on variance expansion, and Equation (97) is based on Assumption 4.2.

By summing Equation (97) for all  $t \in \{0, \dots, T\}$ , we have

$$F_* - F(\tilde{\mathbf{x}}_0) \leq \mathbb{E}F(\tilde{\mathbf{x}}_{T+1}) - F(\tilde{\mathbf{x}}_0) = \sum_{t=0}^T (\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t)) \quad (98)$$

$$\leq -\frac{\eta_s\eta_l K}{4} \sum_{t=0; t \bmod \tau=0}^T \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left( \frac{1}{2\eta_s\eta_l K} - \frac{L}{2} \right) \sum_{t=0; t \bmod \tau=0}^T \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \\ + 2\eta_s\eta_l K \sum_{t=0; t \bmod \tau=0}^T \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\eta_s\eta_l K (T - \lfloor T/\tau \rfloor) \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (99)$$

Let  $\Lambda(t)$  indicate the most recent round where  $p_{\Lambda(t)} = 1$  and  $\Lambda(t) \neq t$ . It is noted that recursively using  $\Lambda(\cdot)$  can achieve the value of 0, i.e.,  $\underbrace{\Lambda(\Lambda(\dots\Lambda(t)))}_{\text{multiple } \Lambda} = 0$ .

To find the bound for  $\sum_{t=0; t \bmod \tau=0}^T \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2$ , the first step is to provide the bound for  $\mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2$ . When  $p_t = 1$ , a client updates the caching gradient with a probability of  $A/M$ , and therefore,  $\mathbb{E}\tilde{g}_t = (1 - \frac{A}{M}) \mathbb{E}\tilde{g}_{\Lambda(t)} + \frac{A}{M} \nabla F(\tilde{\mathbf{x}}_t)$ . Based on this fact, the bound for  $\mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2$  can be derived as follows:

$$\mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 = \mathbb{E} \|\mathbb{E}\tilde{g}_{\Lambda(t)} - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (100)$$

$$= \mathbb{E} \left\| \left( 1 - \frac{A}{M} \right) \mathbb{E} (\tilde{g}_{\Lambda(t)} - \nabla F(\tilde{\mathbf{x}}_{\Lambda(t)})) + (\nabla F(\tilde{\mathbf{x}}_{\Lambda(t)}) - \nabla F(\tilde{\mathbf{x}}_t)) \right\|_2^2 \quad (101)$$

$$\leq \left( 1 - \frac{A}{M} \right) \mathbb{E} \|\mathbb{E}\tilde{g}_{\Lambda(t)} - \nabla F(\tilde{\mathbf{x}}_{\Lambda(t)})\|_2^2 + \frac{M}{A} \mathbb{E} \|\nabla F(\tilde{\mathbf{x}}_{\Lambda(t)}) - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (102)$$

$$\leq \sum_{\theta=0}^{\lfloor t/\tau \rfloor - 1} \left( 1 - \frac{A}{M} \right)^{\lfloor t/\tau \rfloor - \theta} \cdot \frac{M}{A} L^2 \mathbb{E} \|\tilde{\mathbf{x}}_{\theta\tau} - \tilde{\mathbf{x}}_{(\theta+1)\tau}\|_2^2 + \frac{M}{A} L^2 \mathbb{E} \|\tilde{\mathbf{x}}_{\Lambda(t)} - \tilde{\mathbf{x}}_t\|_2^2 \quad (103)$$

where Equation (102) follows  $(\alpha + \beta)^2 \leq \left(1 + \frac{1}{\gamma}\right) \alpha^2 + (1 + \gamma) \beta^2 - \left(\frac{1}{\sqrt{\gamma}}\alpha + \sqrt{\gamma}\beta\right)^2 \leq \left(1 + \frac{1}{\gamma}\right) \alpha^2 + (1 + \gamma) \beta^2$  and  $\gamma = \frac{M-A}{A}$ . With Equation (103), we sum up all  $t \in \{1, \dots, T\}$  and obtain the following result:

$$\sum_{t=0}^T \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (104)$$

$$\leq \sum_{t=0}^T \left( \sum_{\theta=0}^{\lfloor t/\tau \rfloor - 1} \left( 1 - \frac{A}{M} \right)^{\lfloor t/\tau \rfloor - \theta} \cdot \frac{M}{A} L^2 \mathbb{E} \|\tilde{\mathbf{x}}_{\theta\tau} - \tilde{\mathbf{x}}_{(\theta+1)\tau}\|_2^2 + \frac{M}{A} L^2 \mathbb{E} \|\tilde{\mathbf{x}}_{\Lambda(t)} - \tilde{\mathbf{x}}_t\|_2^2 \right) \quad (105)$$

$$\leq \sum_{\theta=0}^{\lfloor T/\tau \rfloor - 1} \frac{M(M-A)}{A^2} L^2 \tau \mathbb{E} \|\tilde{\mathbf{x}}_{\theta\tau} - \tilde{\mathbf{x}}_{(\theta+1)\tau}\|_2^2 + \frac{M}{A} L^2 \sum_{t=0}^T \mathbb{E} \|\tilde{\mathbf{x}}_t - \tilde{\mathbf{x}}_{\Lambda(t)}\|_2^2 \quad (106)$$

$$\leq \sum_{\theta=0}^{\lfloor T/\tau \rfloor - 1} \frac{M(M-A)}{A^2} L^2 \tau^2 \sum_{\Xi=\theta\tau+1}^{(\theta+1)\tau-1} \mathbb{E} \|\tilde{\mathbf{x}}_{\Xi+1} - \tilde{\mathbf{x}}_{\Xi}\|_2^2 \\ + \frac{M}{A} L^2 \sum_{t=0}^T (t - \Lambda(t) - 1) \sum_{\Xi=\Lambda(t)+1}^{t-1} \mathbb{E} \|\tilde{\mathbf{x}}_{\Xi+1} - \tilde{\mathbf{x}}_{\Xi}\|_2^2 \quad (107)$$

$$\leq \frac{M(M-A)}{A^2} L^2 \tau^2 \sum_{t=0; t \bmod \tau=0}^{T-1} \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 + \frac{M}{A} L^2 \tau^2 \sum_{t=0; t \bmod \tau=0}^{T-1} \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \quad (108)$$

where Equation (106) follows that, for all  $\theta \in \{0, \dots, \lfloor T/\tau \rfloor - 1\}$ , the coefficient for  $\frac{M}{A} L^2$  includes  $(1 - \frac{A}{M}), \dots, (1 - \frac{A}{M})^{\lfloor T/\tau \rfloor - \theta}$ , and each of them has a maximum of  $\tau$  ts, meaning that the upper bound of the coefficient should be

$$\tau \left( \left(1 - \frac{A}{M}\right) + \dots + \left(1 - \frac{A}{M}\right)^{\lfloor T/\tau \rfloor - \theta} \right) \leq \tau \cdot \frac{M}{2A} \left(1 - \frac{A}{M}\right); \quad (109)$$

Equation (107) follows Cauchy-Schwarz inequality.

Plugging Equation (108) back to Equation (99), we have:

$$\begin{aligned} F_* - F(\tilde{\mathbf{x}}_0) &\leq -\frac{\eta_s \eta_l K}{4} \sum_{t=0; t \bmod \tau=0}^T \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \\ &\quad - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} - \eta_s \eta_l K L^2 \tau^2 \frac{M^2}{A^2} \right) \sum_{t=0; t \bmod \tau=0}^T \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \\ &\quad + 2\eta_s \eta_l K (T - \lfloor T/\tau \rfloor) \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (110)$$

Since  $\eta_s \eta_l = \frac{1}{KL} \left(1 + \frac{2M\tau}{A}\right)^{-1}$ ,  $\frac{1}{2\eta_s \eta_l K} - \frac{L}{2} - \eta_s \eta_l K L^2 \tau^2 \frac{M^2}{A^2} \geq 0$  such that the term  $\sum_{t=0; t \bmod \tau=0}^T \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2$  can be omitted in Equation (110). Hence, we can easily obtain the following inequality:

$$\frac{1}{T - \lfloor T/\tau \rfloor} \sum_{t=0; t \bmod \tau=0}^T \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \frac{4(F(\tilde{\mathbf{x}}_0) - F_*)}{\eta_s \eta_l K (T - \lfloor T/\tau \rfloor)} + 8 \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (111)$$

By using the settings of the local learning rate and the global learning rate in the description, we can obtain the desired result.  $\square$



## E. Proofs under Constant Probabilistic Settings

### E.1. Preliminary

**Lemma E.1.** *Suppose that Assumption 4.1 holds, and  $p_t \in (0, 1)$ . Let  $\tilde{g}_t$  be the definition of Line 9 of Algorithm 1, i.e., the average of the caching gradients. Therefore, the recursive expression for  $\{\tilde{g}_t\}_{t \geq 0}$  in the expected form is*

$$\mathbb{E}\tilde{g}_t = \begin{cases} (1 - \frac{A}{M}p_{t-1}) \cdot \mathbb{E}\tilde{g}_{t-1} + \frac{A}{M}p_{t-1} \cdot \nabla F(\tilde{\mathbf{x}}_{t-1}), & t > 0 \\ \nabla F(\tilde{\mathbf{x}}_0), & t = 0 \end{cases} \quad (112)$$

Furthermore, when  $t > 0$  we can obtain the following inequality:

$$\mathbb{E}\|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \left(1 - \frac{A}{M}p_{t-1}\right) \cdot \mathbb{E}\|\mathbb{E}\tilde{g}_{t-1} - \nabla F(\tilde{\mathbf{x}}_{t-1})\|_2^2 + \frac{M}{Ap_{t-1}} \cdot L^2 \cdot \mathbb{E}\|\tilde{\mathbf{x}}_t - \tilde{\mathbf{x}}_{t-1}\|_2^2. \quad (113)$$

As for  $t = 0$ , we have  $\mathbb{E}\|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 = 0$ .

*Proof.* According to the definition of Line 9 of Algorithm 1,  $\tilde{g}_{t+1} = \text{avg}(v_{t+1}) = \frac{1}{M} \sum_{m \in [M]} v_{t+1}^{(m)}$ . Hence, for each element in  $v_{t+1}$ , i.e.,  $v_{t+1}^{(m)}$ , where  $m \in [M]$ , they have a probability of  $(1 - \frac{A}{M}p_t)$  to retain the previous value, or otherwise update as anchor clients using large batches. Thus, the expected value for  $\mathbb{E}v_{t+1}^{(m)}$  is:

$$\mathbb{E}v_{t+1}^{(m)} = \frac{A}{M}p_t \cdot \mathbb{E}\nabla f_m(\tilde{\mathbf{x}}_t, \mathcal{B}_{m,t}) + \left(1 - \frac{A}{M}p_t\right) \cdot \mathbb{E}v_t^{(m)} \quad (114)$$

$$= \frac{A}{M}p_t \cdot \nabla F_m(\tilde{\mathbf{x}}_t) + \left(1 - \frac{A}{M}p_t\right) \cdot \mathbb{E}v_t^{(m)} \quad (115)$$

Therefore, we have

$$\mathbb{E}\tilde{g}_{t+1} = \frac{1}{M} \sum_{m=1}^M \mathbb{E}v_{t+1}^{(m)} = \frac{A}{M}p_t \cdot \nabla F(\tilde{\mathbf{x}}_t) + \left(1 - \frac{A}{M}p_t\right) \cdot \mathbb{E}\tilde{g}_t \quad (116)$$

It is worth noting that  $\mathbb{E}\tilde{g}_0 = \nabla F(\tilde{\mathbf{x}}_0)$  as it is initialized at the beginning of the training, i.e., Line 2 – 4 in Algorithm 1. Therefore,  $\mathbb{E}\|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 = 0$ .

Next, we find the recursive bound for  $\mathbb{E}\|\mathbb{E}\tilde{g}_{t+1} - \nabla F(\tilde{\mathbf{x}}_{t+1})\|_2^2$ :

$$\mathbb{E}\|\mathbb{E}\tilde{g}_{t+1} - \nabla F(\tilde{\mathbf{x}}_{t+1})\|_2^2 \quad (117)$$

$$= \mathbb{E}\left\|\left(1 - \frac{A}{M}p_t\right) \cdot (\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)) + \nabla F(\tilde{\mathbf{x}}_t) - \nabla F(\tilde{\mathbf{x}}_{t+1})\right\|_2^2 \quad (118)$$

$$\leq \left(1 + \frac{Ap_t}{M - Ap_t}\right) \left(1 - \frac{A}{M}p_t\right)^2 \mathbb{E}\|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \left(1 + \frac{M - Ap_t}{Ap_t}\right) \mathbb{E}\|\nabla F(\tilde{\mathbf{x}}_t) - \nabla F(\tilde{\mathbf{x}}_{t+1})\|_2^2 \quad (119)$$

$$\leq \left(1 - \frac{A}{M}p_t\right) \mathbb{E}\|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \frac{M}{Ap_t} L^2 \mathbb{E}\|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \quad (120)$$

where Equation (119) follows  $(\alpha + \beta)^2 \leq \left(1 + \frac{1}{\gamma}\right) \alpha^2 + (1 + \gamma) \beta^2 - \left(\frac{1}{\gamma} \alpha + \sqrt{\gamma} \beta\right)^2 \leq \left(1 + \frac{1}{\gamma}\right) \alpha^2 + (1 + \gamma) \beta^2$ , and Equation (120) follows Assumption 4.1.  $\square$

**Lemma E.2.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local learning rate satisfy  $\eta_l \leq \min\left(\frac{1}{2\sqrt{3}KL}, \frac{1}{2\sqrt{3}L_g^2} \sqrt{\frac{b'}{K}}\right)$ . With FedAMD,  $\sum_{k=0}^{K-1} \|\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)}\|_2^2$  represents the sum of the second norm of every iteration's difference. Therefore, the bound for such a summation in the expected form should be*

$$\sum_{k=0}^{K-1} \mathbb{E}\left\|\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)}\right\|_2^2 \leq 2\eta_l^2(K+1) \frac{\sigma^2}{Mb} + 6\eta_l^2(K+1) \mathbb{E}\|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 6\eta_l^2(K+1) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (121)$$

*Proof.* According to Equation (35), the update at  $(k-1)$ -th iteration is

$$\mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} = -\eta_l g_{t,k}^{(m)} = -\eta_l \left( g_{t,0}^{(m)} - \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) + \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) \right). \quad (122)$$

To find the bound for the expected value of its second norm, the analysis is presented as follows:

$$\mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (123)$$

$$= \eta_l^2 \mathbb{E} \left\| \tilde{g}_t - \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) + \sum_{\kappa=0}^{k-1} \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) \right\|_2^2 \quad (124)$$

$$= \eta_l^2 \mathbb{E} \left\| \tilde{g}_t - \mathbb{E} \tilde{g}_t - \sum_{\kappa=0}^{k-1} \left( \nabla f_m \left( \mathbf{x}_{t,\kappa-1}^{(m)}, \mathcal{B}'_{m,\kappa} \right) - \nabla f_m \left( \mathbf{x}_{t,\kappa}^{(m)}, \mathcal{B}'_{m,\kappa} \right) - \nabla F_m \left( \mathbf{x}_{t,\kappa-1}^{(m)} \right) + \nabla F_m \left( \mathbf{x}_{t,\kappa}^{(m)} \right) \right) \right\|_2^2$$

$$+ \eta_l^2 \mathbb{E} \left\| \mathbb{E} \tilde{g}_t - \sum_{\kappa=0}^{k-1} \nabla F_m \left( \mathbf{x}_{t,\kappa-1}^{(m)} \right) + \sum_{\kappa=0}^{k-1} \nabla F_m \left( \mathbf{x}_{t,\kappa}^{(m)} \right) \right\|_2^2 \quad (125)$$

$$= \eta_l^2 \left( \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + \sum_{\kappa=0}^{k-1} \frac{L_\sigma^2}{b'} \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \right)$$

$$+ \eta_l^2 \mathbb{E} \left\| \mathbb{E} \tilde{g}_t - \sum_{\kappa=0}^{k-1} \left( \nabla F_m \left( \mathbf{x}_{t,\kappa-1}^{(m)} \right) - \nabla F_m \left( \mathbf{x}_{t,\kappa}^{(m)} \right) \right) \right\|_2^2 \quad (126)$$

$$= \eta_l^2 \left( \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + \sum_{\kappa=0}^{k-1} \frac{L_\sigma^2}{b'} \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \right)$$

$$+ \eta_l^2 \mathbb{E} \left\| \mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t) + \nabla F(\tilde{\mathbf{x}}_t) + \sum_{\kappa=0}^{k-1} \left( \nabla F_m \left( \mathbf{x}_{t,\kappa}^{(m)} \right) + \nabla F_m \left( \mathbf{x}_{t,\kappa-1}^{(m)} \right) \right) \right\|_2^2 \quad (127)$$

$$\leq \eta_l^2 \left( \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + \sum_{\kappa=0}^{k-1} \frac{L_\sigma^2}{b'} \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \right)$$

$$+ 3\eta_l^2 \cdot \mathbb{E} \left\| \mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 + 3\eta_l^2 \left\| \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 + 3\eta_l^2 K \sum_{\kappa=0}^{k-1} L^2 \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \quad (128)$$

$$= \eta_l^2 \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + 3\eta_l^2 \cdot \mathbb{E} \left\| \mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 + 3\eta_l^2 \left\| \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2$$

$$+ 3\eta_l^2 \left( KL^2 + \frac{L_\sigma^2}{b'} \right) \sum_{\kappa=0}^{k-1} \mathbb{E} \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \quad (129)$$

where Equation (125) is based on the variance expansion on the first term of Equation (124); Equation (128) is based on Cauchy-Schwarz inequality.

Therefore, by summing Equation (129) for  $k = 1, \dots, K$ , we have

$$\mathbb{E} \sum_{k=0}^{K-1} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \leq \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (130)$$

$$\leq \eta_l^2 (K+1) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + 3\eta_l^2 (K+1) \mathbb{E} \left\| \mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 + 3\eta_l^2 (K+1) \left\| \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2$$

$$+ 3\eta_l^2 K \left( KL^2 + \frac{L_\sigma^2}{b'} \right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (131)$$

Obviously, according to the setting of the local learning rate in the description above, the inequality  $3\eta_l^2 K \left( KL^2 + \frac{L^2}{b'} \right) \leq \frac{1}{2}$  holds. Therefore, we can easily obtain the bound for the sum of the second norm of every iteration's difference, which is consistent with Equation (121).  $\square$

## E.2. Proofs for Non-convex Objectives

The following lemma provides a recursive expression on  $\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t)$  for time-varying probability settings.

**Lemma E.3.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold, and the time-varying probability sequence  $\{p_t \in (0, 1)\}_{t \geq 0}$ . Let the local updates  $K \geq 1$ , and the local learning rate  $\eta_l \leq \min\left(\frac{1}{2\sqrt{6}KL}, \frac{\sqrt{b'/K}}{2\sqrt{3}L\sigma}\right)$ . With the model training using FedAMD, the recursive function between  $F(\tilde{\mathbf{x}}_{t+1})$  and  $F(\tilde{\mathbf{x}}_t)$  in expected form is*

$$\begin{aligned} \mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) &\leq -\frac{\eta_s \eta_l K}{4} \left(1 - (p_t)^A\right) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left(\frac{1}{2\eta_s \eta_l K} - \frac{L}{2}\right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \\ &\quad + 4\eta_s \eta_l K \left(1 - (p_t)^A\right) \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s \eta_l K \left(1 - (p_t)^A\right) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (132)$$

*Proof.* According to Lemma D.1, we have:

$$\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \quad (133)$$

$$\leq -\frac{\eta_s \eta_l K}{2} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left(\frac{1}{2\eta_s \eta_l K} - \frac{L}{2}\right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 + \frac{\eta_s}{2\eta_l K} \mathbb{E} \|\Delta \mathbf{x}_t - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (134)$$

When  $|\Delta \mathbf{x}_t| = 0$ , the probability will be  $(1 - q)^a$ , and when  $|\Delta \mathbf{x}_t| \neq 0$ , the probability will be  $1 - (1 - q)^a$ . Next, we find the bound for the third term of Equation (134), i.e.,  $\mathbb{E} \|\Delta \mathbf{x}_t - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2$ . By Lemma B.1, we have the following derivation:

$$\mathbb{E} \|\Delta \mathbf{x}_t - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (135)$$

$$\leq (1 - (p_t)^A) \frac{1}{M} \sum_{m=1}^M \mathbb{E} \left\| \Delta \mathbf{x}_t^{(m)} - \eta_l K \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 + (p_t)^A \|\eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (136)$$

$$\begin{aligned} &= (1 - (p_t)^A) \frac{1}{M} \sum_{m=1}^M \left( \mathbb{E} \left\| \Delta \mathbf{x}_t^{(m)} - \mathbb{E} \Delta \mathbf{x}_t^{(m)} \right\|_2^2 + \mathbb{E} \left\| \mathbb{E} \Delta \mathbf{x}_t^{(m)} - \eta_l K \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 \right) \\ &\quad + (p_t)^A \eta_l^2 K^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \end{aligned} \quad (137)$$

where Equation (137) follows variance equation. To find the bound for Equation (137), we first analyze its first term, i.e.,  $\mathbb{E} \left\| \Delta \mathbf{x}_t^{(m)} - \mathbb{E} \Delta \mathbf{x}_t^{(m)} \right\|_2^2$ . According to Section C, we have:

$$\mathbb{E} \left\| \Delta \mathbf{x}_t^{(m)} - \mathbb{E} \Delta \mathbf{x}_t^{(m)} \right\|_2^2 \quad (138)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\tilde{\mathbf{g}}_t - \mathbb{E}\tilde{\mathbf{g}}_t\|_2^2 + 2\eta_l^2 \sum_{k=0}^{K-1} (K-1) \frac{L_\sigma^2}{b'} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (139)$$

$$\leq 2\eta_l^2 K^2 \left(1 + 2\eta_l^2 \frac{L_\sigma^2}{b'}\right) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + 12\eta_l^4 K^2 \frac{L_\sigma^2}{b'} \left( \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \quad (140)$$

where Equation (140) follows Lemma E.2. According to the local learning rate setting in the description, we have

$$\mathbb{E} \left\| \Delta \mathbf{x}_t^{(m)} - \mathbb{E} \Delta \mathbf{x}_t^{(m)} \right\|_2^2 \leq 4\eta_l^2 K^2 \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (141)$$

$$+ 12\eta_l^4 K^2 \frac{L_\sigma^2}{b'} \left( \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \quad (142)$$

After finding the bound for the first term of Equation (137), we now give the bound for its second term, i.e.,

$$\mathbb{E} \left\| \mathbb{E} \Delta \mathbf{x}_t^{(m)} - \eta_l K \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2.$$

$$\mathbb{E} \left\| \mathbb{E} \Delta \mathbf{x}_t^{(m)} - \eta_l K \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 \quad (143)$$

$$= \mathbb{E} \left\| \eta_l K (\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)) + \eta_l \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) - \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) \right) \right\|_2^2 \quad (144)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\eta_l^2 \left\| \sum_{k=0}^{K-1} \sum_{\kappa=0}^k \left( \nabla F_m(\mathbf{x}_{t,\kappa}^{(m)}) - \nabla F_m(\mathbf{x}_{t,\kappa-1}^{(m)}) \right) \right\|_2^2 \quad (145)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\eta_l^2 K L^2 \sum_{k=0}^{K-1} \sum_{\kappa=0}^k k \left\| \mathbf{x}_{t,\kappa}^{(m)} - \mathbf{x}_{t,\kappa-1}^{(m)} \right\|_2^2 \quad (146)$$

$$\leq 2\eta_l^2 K^2 \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 2\eta_l^2 K \frac{K(K-1)}{2} L^2 \sum_{k=0}^{K-1} \left\| \mathbf{x}_{t,k}^{(m)} - \mathbf{x}_{t,k-1}^{(m)} \right\|_2^2 \quad (147)$$

$$\leq 2\eta_l^4 K^4 L^2 \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + 2\eta_l^2 K^2 (1 + 3\eta_l^2 K^2 L^2) \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 6\eta_l^4 K^4 L^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (148)$$

where Equation (145) follows  $(\alpha + \beta)^2 \leq 2\alpha^2 + 2\beta^2$ ; Equation (146) follows Cauchy–Schwarz inequality and Assumption 4.1; Equation (148) is based on Lemma E.2. Then, according to the setting for the local learning rate in the description above, we can further simplify Equation (148):

$$\mathbb{E} \left\| \mathbb{E} \Delta \mathbf{x}_t^{(m)} - \eta_l K \nabla F(\tilde{\mathbf{x}}_t) \right\|_2^2 \leq 2\eta_l^4 K^4 L^2 \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} + 4\eta_l^2 K^2 \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 6\eta_l^4 K^4 L^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (149)$$

Plugging Equation (142) and Equation (149) back to Equation (137), we can primarily obtain the inequality below:

$$\begin{aligned} \mathbb{E} \|\Delta \mathbf{x}_t - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 &\leq 2\eta_l^2 K^2 \left(1 - (p_t)^A\right) \left(2 + \eta_l^2 K^2 L^2\right) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \\ &\quad + 4\eta_l^2 K^2 \left(1 - (p_t)^A\right) \left(1 + 3\eta_l^2 \frac{L^2 \sigma}{b'}\right) \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \\ &\quad + 6\eta_l^4 K^2 \left(1 - (p_t)^A\right) \left(\frac{2L^2 \sigma}{b'} + K^2 L^2\right) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \end{aligned} \quad (150)$$

$$+ (p_t)^A \eta_l^2 K^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (151)$$

With the setting described in the Lemma, we have:

$$\begin{aligned} \mathbb{E} \|\Delta \mathbf{x}_t - \eta_l K \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 &\leq 6\eta_l^2 K^2 \left(1 - (p_t)^A\right) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \\ &\quad + 8\eta_l^2 K^2 \left(1 - (p_t)^A\right) \mathbb{E} \|\mathbb{E} \tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \\ &\quad + 6\eta_l^4 K^2 \left(1 - (p_t)^A\right) \left(\frac{2L^2 \sigma}{b'} + K^2 L^2\right) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \end{aligned} \quad (152)$$

$$+ (p_t)^A \eta_l^2 K^2 \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \quad (153)$$

Therefore, according to the upper bound analyzed in the previous inequalities, Equation (134) can be reformulated as

$$\begin{aligned} &\mathbb{E} F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \quad (154) \\ &\leq -\frac{\eta_s \eta_l K}{2} \left(1 - (p_t)^A\right) \left(1 - 6\eta_l^2 \left(\frac{2L^2 \sigma}{b'} + K^2 L^2\right)\right) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left(\frac{1}{2\eta_s \eta_l K} - \frac{L}{2}\right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \end{aligned}$$

$$+ 4\eta_s\eta_l K (1 - (p_t)^A) \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s\eta_l K (1 - (p_t)^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (155)$$

We can obtain the desired conclusion through the setting in the description above.  $\square$

**Theorem E.4.** *Suppose that Assumption 4.1, 4.2 and 4.3 hold. Let the local updates  $K \geq 1$ , and the local learning rate  $\eta_l$  and the global learning rate  $\eta_s$  be  $\eta_s\eta_l = \frac{1}{KL} \left(1 + \frac{2M}{Ap} \sqrt{1 - p^A}\right)^{-1}$ , where  $\eta_l \leq \min\left(\frac{1}{2\sqrt{6}KL}, \frac{\sqrt{b'}/K}{4\sqrt{3}L\sigma}\right)$ . Therefore, the convergence rate of FedAMD for non-convex objectives should be*

$$\min_{t \in [T]} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq O\left(\frac{1}{T} \left(\frac{1}{1 - p^A} + \frac{M}{Ap\sqrt{1 - p^A}}\right)\right) + O\left(\mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb}\right) \quad (156)$$

where we treat  $F(\tilde{\mathbf{x}}_0) - F_*$  and  $L$  as constants.

*Proof.* With Lemma E.1 and Lemma E.3, we can find the following recursive function under the constant probability settings:

$$\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) + \frac{4\eta_s\eta_l K (1 - p^A) M}{Ap} \mathbb{E} \|\mathbb{E}\tilde{g}_{t+1} - \nabla F(\tilde{\mathbf{x}}_{t+1})\|_2^2 \quad (157)$$

$$\begin{aligned} &\leq \mathbb{E}F(\tilde{\mathbf{x}}_t) + \frac{4\eta_s\eta_l K (1 - p^A) M}{Ap} \mathbb{E} \|\mathbb{E}\tilde{g}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \frac{\eta_s\eta_l K}{4} (1 - p^A) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \\ &\quad - \left(\frac{1}{2\eta_s\eta_l K} - \frac{L}{2} - \frac{4\eta_s\eta_l K (1 - p^A) M}{Ap} \frac{M}{Ap} L^2\right) \mathbb{E} \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \\ &\quad + 3\eta_s\eta_l K (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (158)$$

Since  $\eta_s\eta_l = \frac{1}{KL} \left(1 + \frac{2M}{Ap} \sqrt{1 - p^A}\right)^{-1}$ , we have:

$$F_* \leq \mathbb{E}F(\tilde{\mathbf{x}}_T) \leq \mathbb{E}F(\tilde{\mathbf{x}}_T) + \frac{4\eta_s\eta_l K (1 - p^A) M}{Ap} \mathbb{E} \|\mathbb{E}\tilde{g}_T - \nabla F(\tilde{\mathbf{x}}_T)\|_2^2 \quad (159)$$

$$\begin{aligned} &\leq \mathbb{E}F(\tilde{\mathbf{x}}_{T-1}) + \frac{4\eta_s\eta_l K (1 - p^A) M}{Ap} \mathbb{E} \|\mathbb{E}\tilde{g}_{T-1} - \nabla F(\tilde{\mathbf{x}}_{T-1})\|_2^2 \\ &\quad - \frac{\eta_s\eta_l K}{4} (1 - p^A) \|\nabla F(\tilde{\mathbf{x}}_{T-1})\|_2^2 + 3\eta_s\eta_l K (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (160)$$

$$\begin{aligned} &\leq F(\tilde{\mathbf{x}}_0) + \frac{4\eta_s\eta_l K (1 - p^A) M}{Ap} \|\mathbb{E}\tilde{g}_0 - \nabla F(\tilde{\mathbf{x}}_0)\|_2^2 \\ &\quad - \frac{\eta_s\eta_l K}{4} (1 - p^A) \sum_{t=0}^{T-1} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s\eta_l K T (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (161)$$

According to Lemma E.1,  $\|\mathbb{E}\tilde{g}_0 - \nabla F(\tilde{\mathbf{x}}_0)\|_2^2 = 0$ . Therefore, based on the derivation above, we can attain the following inequality:

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \leq \frac{4(F(\tilde{\mathbf{x}}_0) - F_*)}{\eta_s\eta_l K T (1 - p^A)} + 3\mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (162)$$

By using the settings of the local learning rate and the global learning rate in the description, we can obtain the desired result.  $\square$

### E.3. Proofs for PL Condition

**Theorem E.5.** *Suppose that Assumption 4.1, 4.2, 4.3 and 4.10 hold. Let the local updates  $K \geq 1$ , and the local learning rate  $\eta_l$  and the global learning rate  $\eta_s$  be  $\eta_s \eta_l = \min \left( \frac{Ap}{MK\mu(1-p^A)}, \frac{1}{KL(1+\frac{16M}{\mu Ap}L)} \right)$ , where  $\eta_l \leq \min \left( \frac{1}{2\sqrt{6}KL}, \frac{\sqrt{b'/K}}{4\sqrt{3}L\sigma} \right)$ . Therefore, the convergence rate of FedAMD for PL condition should be*

$$\begin{aligned} \mathbb{E}F(\tilde{\mathbf{x}}_T) - F_* &\leq \left( 1 - \frac{1}{2}\mu K (1-p^A) \min \left( \frac{Ap}{MK\mu(1-p^A)}, \frac{1}{KL(1+\frac{16M}{\mu Ap}L)} \right) \right)^T (F(\tilde{\mathbf{x}}_0) - F_*) \\ &\quad + O \left( \frac{1}{\mu} \cdot \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \right) \end{aligned} \quad (163)$$

*Proof.* With Lemma E.3, we have the recursive function on the time-varying probability settings under PL condition:

$$\begin{aligned} &\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F(\tilde{\mathbf{x}}_t) \quad (164) \\ &\leq -\frac{\eta_s \eta_l K}{4} (1 - (p_t)^A) \|\nabla F(\tilde{\mathbf{x}}_t)\|_2^2 - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} \right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \end{aligned}$$

$$+ 4\eta_s \eta_l K (1 - (p_t)^A) \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s \eta_l K (1 - (p_t)^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (165)$$

$$\begin{aligned} &\leq -\frac{\mu \eta_s \eta_l K}{2} (1 - (p_t)^A) (F(\tilde{\mathbf{x}}_t) - F_*) - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} \right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \\ &\quad + 4\eta_s \eta_l K (1 - (p_t)^A) \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s \eta_l K (1 - (p_t)^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (166)$$

According to the description, we consider the probability  $p_t = p$  and have:

$$\begin{aligned} &\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F_* \quad (167) \\ &\leq \left( 1 - \frac{\mu \eta_s \eta_l K}{2} (1 - p^A) \right) (F(\tilde{\mathbf{x}}_t) - F_*) - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} \right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 \end{aligned}$$

$$+ 4\eta_s \eta_l K (1 - p^A) \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 + 3\eta_s \eta_l K (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (168)$$

Since  $\eta_s \eta_l \leq \frac{Ap}{MK\mu(1-p^A)}$ , we have:

$$\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F_* + \frac{8}{\mu} \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_{t+1} - \nabla F(\tilde{\mathbf{x}}_{t+1})\|_2^2 \quad (169)$$

$$\begin{aligned} &\leq \left( 1 - \frac{\mu \eta_s \eta_l K}{2} (1 - p^A) \right) \left( F(\tilde{\mathbf{x}}_t) - F_* + \frac{8}{\mu} \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) \\ &\quad - \left( \frac{1}{2\eta_s \eta_l K} - \frac{L}{2} - \frac{8}{\mu} \frac{M}{Ap} L^2 \right) \|\tilde{\mathbf{x}}_{t+1} - \tilde{\mathbf{x}}_t\|_2^2 + 3\eta_s \eta_l K (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (170)$$

According to the description  $\eta_s \eta_l \leq \frac{1}{KL(1+\frac{16M}{\mu Ap}L)}$ , we have:

$$\mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F_* \leq \mathbb{E}F(\tilde{\mathbf{x}}_{t+1}) - F_* + \frac{8}{\mu} \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_{t+1} - \nabla F(\tilde{\mathbf{x}}_{t+1})\|_2^2 \quad (171)$$

$$\leq \left( 1 - \frac{\mu \eta_s \eta_l K}{2} (1 - p^A) \right) \left( F(\tilde{\mathbf{x}}_t) - F_* + \frac{8}{\mu} \mathbb{E} \|\mathbb{E}\tilde{\mathbf{g}}_t - \nabla F(\tilde{\mathbf{x}}_t)\|_2^2 \right) + 3\eta_s \eta_l K (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (172)$$

$$\begin{aligned} &\leq \left( 1 - \frac{\mu \eta_s \eta_l K}{2} (1 - p^A) \right)^{t+1} (F(\tilde{\mathbf{x}}_0) - F_*) \\ &\quad + \left( 1 + \dots + \left( 1 - \frac{\mu \eta_s \eta_l K}{2} (1 - p^A) \right)^t \right) 3\eta_s \eta_l K (1 - p^A) \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \end{aligned} \quad (173)$$



$$= \left(1 - \frac{\mu\eta_s\eta_l K}{2} (1 - p^A)\right)^{t+1} (F(\tilde{\mathbf{x}}_0) - F_*) + \frac{6}{\mu} \mathbf{1}_{\{b < n\}} \frac{\sigma^2}{Mb} \quad (174)$$

By using the settings of the local learning rate and the global learning rate in the description, we can obtain the desired result. □

## F. Additional Experiments

In the main text, we have analyzed some experimental results in Section 5. In this part, we conduct more thorough experiments by setting different numbers of local updates and different secondary mini-batch sizes.

### F.1. Detailed Experimental Setup

**Training on Fashion MNIST.** In Section 5, the experiment conducts on Fashion MNIST (Xiao et al., 2017), an image classification task to categorize a  $28 \times 28$  greyscale image into 10 labels (including T-shirt/top, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot). In the training dataset, each class owns 6K samples. Then, we follow the setting of (Konečný et al., 2016; Li et al., 2019) and partition the dataset into 100 clients ( $M = 100$ ) such that each client holds two classes with a total of 600 samples. By this means, we simulate the heterogeneous data setting. To obtain a recognizable model on the images in the test dataset, we utilize a convolutional neural network structure LeNet-5 (LeCun et al., 1989; 2015). Below comprehensively presents the structure of LeNet-5 on Fashion MNIST:

Table 3. Details for LeNet-5 on Fashion-MNIST.

Layer	Output Shape	Trainable Parameters	Activation	Hyperparameters
Input	(1, 28, 28)	0		
Conv2d	(6, 24, 24)	156	ReLU	kernel size=5
MaxPool2d	(6, 12, 12)	0		kernel size=2
Conv2d	(16, 8, 8)	2416	ReLU	kernel size=5
MaxPool2d	(16, 4, 4)	0		kernel size=2
Flatten	256	0		
Dense	120	30840	ReLU	
Dense	84	10164	ReLU	
Dense	10	850	softmax	

**Training on EMNIST digits.** In addition to Fashion MNIST, we utilize one more dataset EMNIST (Cohen et al., 2017) digits to further assess our approach efficiency. This task is to recognize 10 handwritten digits with a total of 240K training samples and 40K test samples. Similar to Fashion MNIST, we equally disjoint the dataset into 100 clients ( $M = 100$ ), and each client possesses two classes. The model is trained with a 2-layer MLP (Yue et al., 2022), i.e.,

Table 4. Details for 2-layer MLP on EMNIST digits.

Layer	Output Shape	Trainable Parameters	Activation	Hyperparameters
Input	(1, 28, 28)	0		
Flatten	784	0		
Dense	100	78500	ReLU	
Dense	10	1010	softmax	

**Validation metrics.** The training loss is calculated by the clients who perform local SGD on the average loss of all iterations. As for the test accuracy, the server utilizes the entire test dataset after the global model updates. The gradient complexity is the sum of all samples used for gradient calculation by all clients throughout the training. The communication overhead is measured by the transmission between the server and the clients.

**Miscellaneous.** Our simulation experiment runs on Ubuntu 18.04 with Intel(R) Xeon(R) Gold 6254 CPU, NVIDIA RTX A6000 GPU, and CUDA 11.2. Our code is implemented using Python and PyTorch v.1.12.1. Clients are picked randomly and uniformly, without replacement in one round but with replacement in subsequent rounds. For each baseline, the local learning ( $\eta_l$ ) rate picks the best one from the set  $\{0.1, 0.03, 0.02, 0.01, 0.008, 0.005\}$ , while the global learning rate ( $\eta_s$ ) is selected from the set  $\{1.0, 0.8, 0.1\}$ . Without the annotation, we implicitly assume Fashion MNIST follows these settings: small minibatch size  $b' = 64$ , large minibatch size  $b = full$ , and the number of local updates  $K = 10$ . As for EMNIST, we suppose the anchor nodes utilize the entire dataset for the caching gradient, i.e.,  $b = full$ , and the number of participants in each round is 20, i.e.,  $A = 20$ . Besides, to make BVR-L-SGD (Murata & Suzuki, 2021) compatible with partial participation in FL training, we only use sampled clients to compute the full gradients of local objectives instead of using all clients.

## F.2. More Numerical Results on Fashion MNIST

In addition to the empirical results in Section 5, we evaluate the performance of FedAMD by using different large minibatch  $b$  settings. Then, considering the number of local updates  $K$ , we assess the performance of the algorithm under various probability settings and compare it with other baselines.

### F.2.1. COMPARISON AMONG VARIOUS HYPER-PARAMETER SETTINGS

**The setting of large mini-batch  $b$ .** Figure 2 – 4 depict the performance of FedAMD under constant probability  $p = 0.9$ , optimal constant probability, and optimal sequential probability, respectively, when the algorithm uses different  $b$ s. Overall,  $b = full$  always outperform  $b = 256$  and  $b = 64$ . Although there is no distinct difference between  $b = 256$  and  $b = 64$  in terms of final test accuracy and training loss,  $b = 256$  is easier to attain a lower training loss during the training.

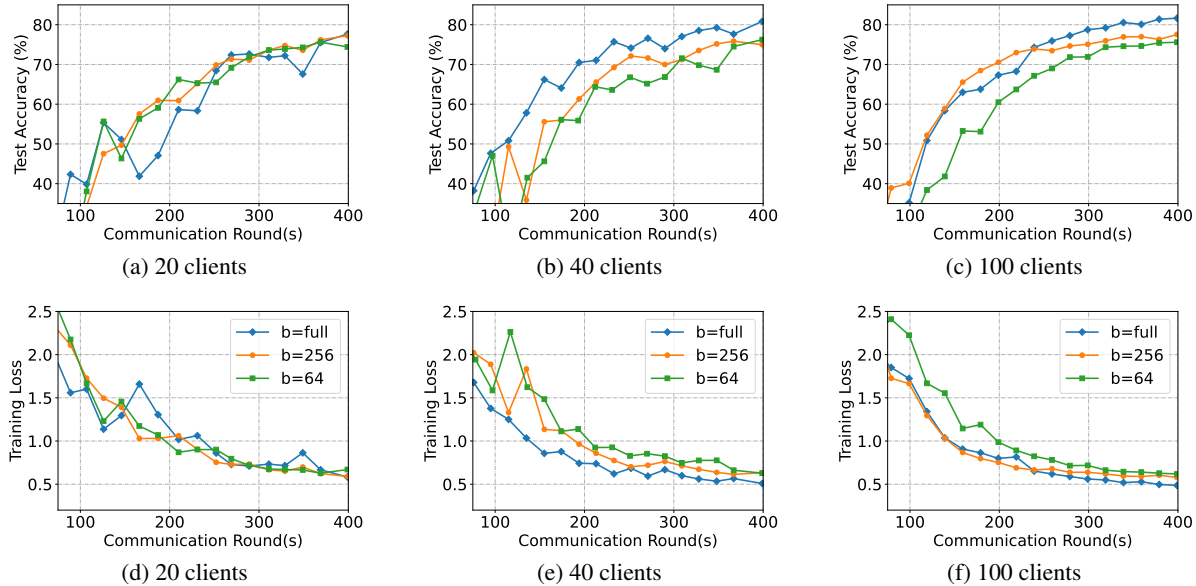


Figure 2. Comparison of test accuracy and training loss against the communication rounds for FedAMD with constant  $p = 0.9$ .

**Number of local updates  $K$ .** Figure 5 – 7 present the performance of FedAMD under the setting of  $K = 10$ ,  $K = 20$ , and  $K = 5$ , respectively. In Section 5, we present the results under  $K = 10$  (Figure 5), which manifests that: (I) The setting  $\{0, 1\}$  is the most efficient performance under the sequential probability setting; (II) The setting near the optimal probability can attain the best result under the constant probability settings. In this part, we verify whether these two statements still hold in two more examples. As for  $K = 20$  and  $K = 5$ , it can provide the best performance when the constant probability is set to be near the theoretical optimal one. However, statement (I) does not always hold in both settings. Specifically, when all clients participate in the training,  $\{0, 0, 1\}$  even outperforms  $\{0, 1\}$ . A possible reason is that  $\{0, 0, 1\}$  has more rounds to update the global model while the caching gradient does not significantly change compared to the situation running for one more round.

### F.2.2. COMPARISON AMONG VARIOUS BASELINES

In Section 5, we present the comparison in a tabular format. Then, in this part, we visualize the training progress as well as introduce more results under different  $K$ s with the help of Figures 8, 9, and 10. In specific, Figure 9a – 9f are summarized into Table 2, while the rest explore the efficiency of FedAMD under more scenarios. As described in Table 2 and the first six figures, when  $K = 10$ , the conclusions we can draw include: (I) the final test accuracy of FedAMD exceeds that of the baselines; (II) FedAMD is able to attain an accurate model with less communication and computation consumption. Next, we evaluate the performance of FedAMD when  $K = 20$  and  $K = 5$ .

- $K = 20$  (Figure 10a – 10f): In this case, the gradient computation of a miner is around twice as that of an anchor. Therefore, FedAMD may consume less computation overhead than FedAvg and SCAFFOLD. In terms of final test accuracy, these baselines achieve similar results in all cases, while FedAMD achieves the performance with less computation overhead.

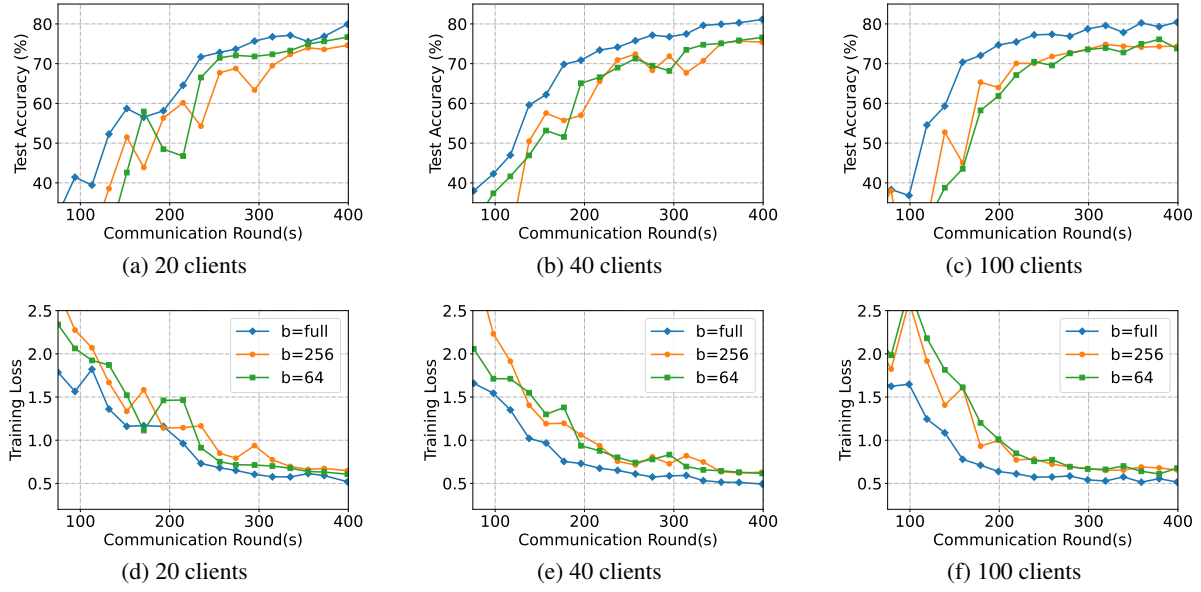


Figure 3. Comparison of test accuracy and training loss against the communication rounds for FedAMD with optimal constant probability  $p$ .

- $K = 5$  (Figure 8a – 8f): FedAMD eventually achieves the best accuracy compared to the existing works. Additionally, we can obtain a well-performed model with less computational consumption. These two phenomena are in support of the statements mentioned above.

### F.3. Numerical Results on EMNIST digits

In this section, Figure 11 analyzes our algorithm with one more dataset, i.e., EMNIST. In the first three figures, we evaluate different probability settings. As for the rest of the figures, we compare FedAMD with other baselines.

With regards to different probability settings (Figure 11a – 11c), we are still able to draw two conclusions as stated in Appendix F.2.1. As for the comparison among different algorithms, our proposed algorithm is able to outperform the state-of-the-art works when we take the test accuracy and the computation overhead into joint consideration.

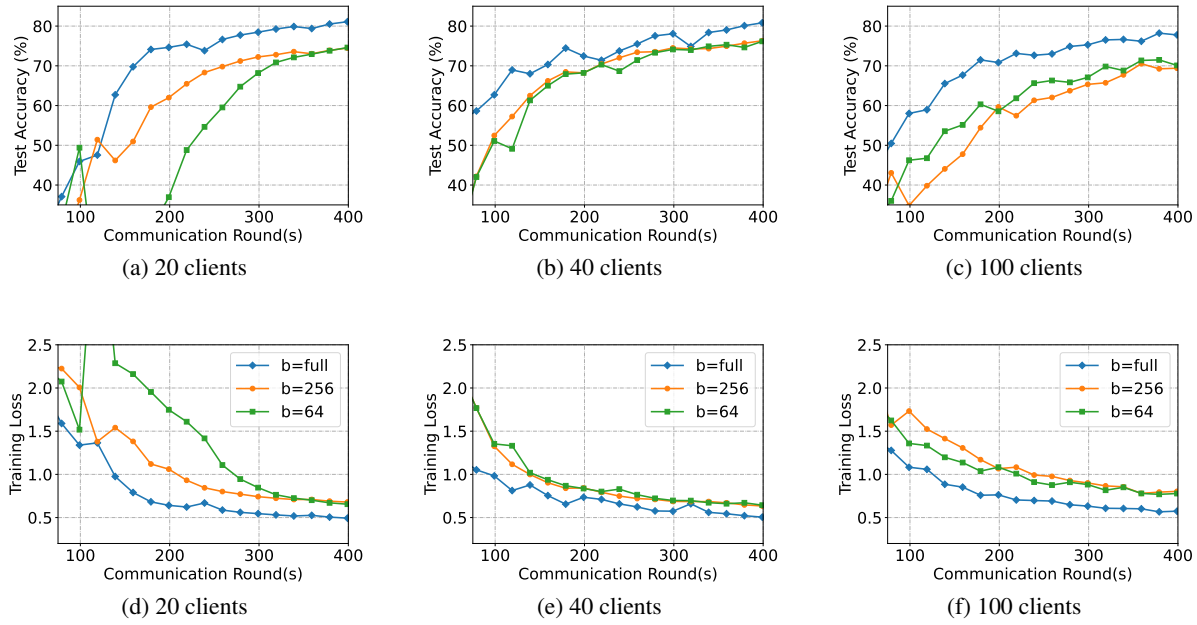


Figure 4. Comparison of test accuracy and training loss against the communication rounds for FedAMD with sequential  $\{0, 1\}$ .

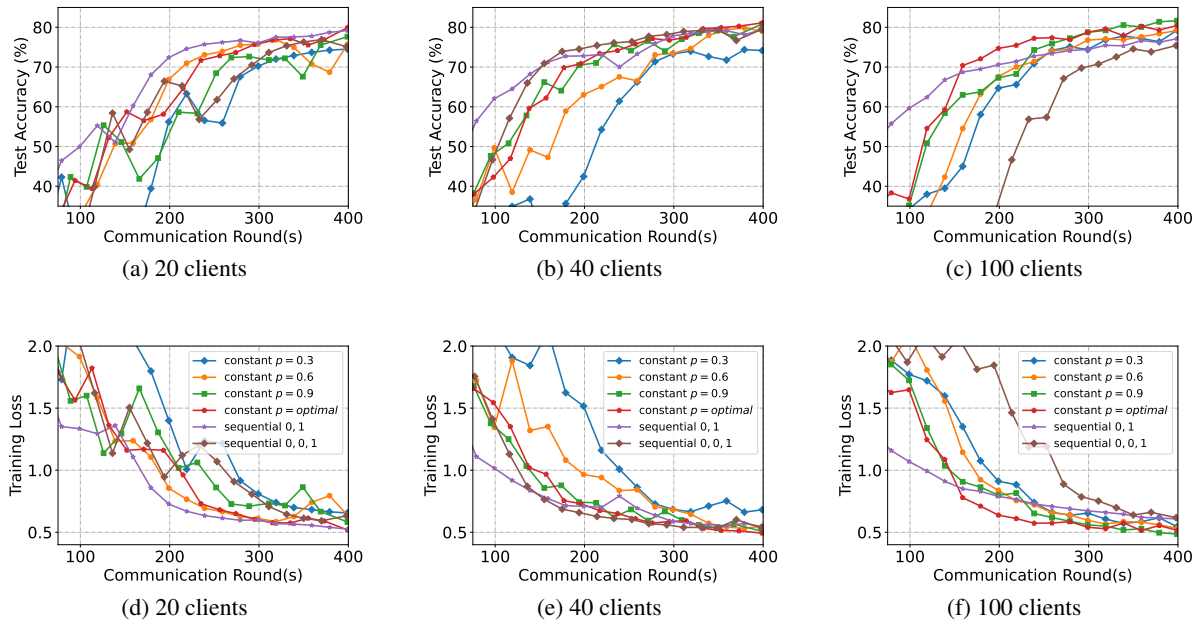


Figure 5. Comparison of different probability settings using training loss and test accuracy against the communication rounds for FedAMD by setting  $K = 10$ .

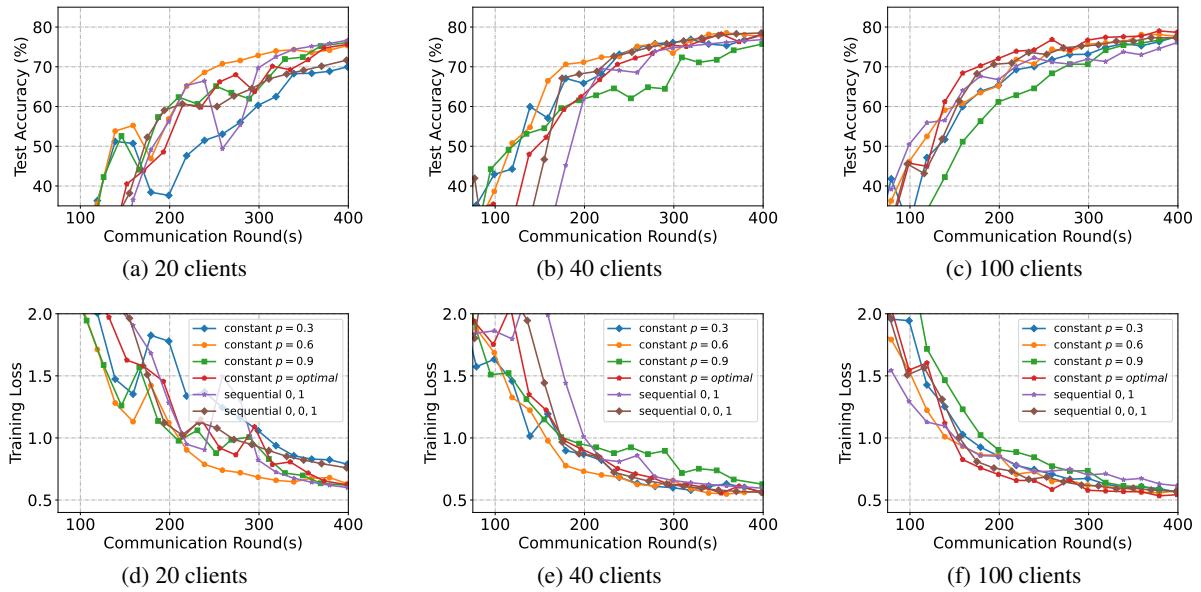


Figure 6. Comparison of different probability settings using training loss and test accuracy against the communication rounds for FedAMD by setting  $K = 20$ .

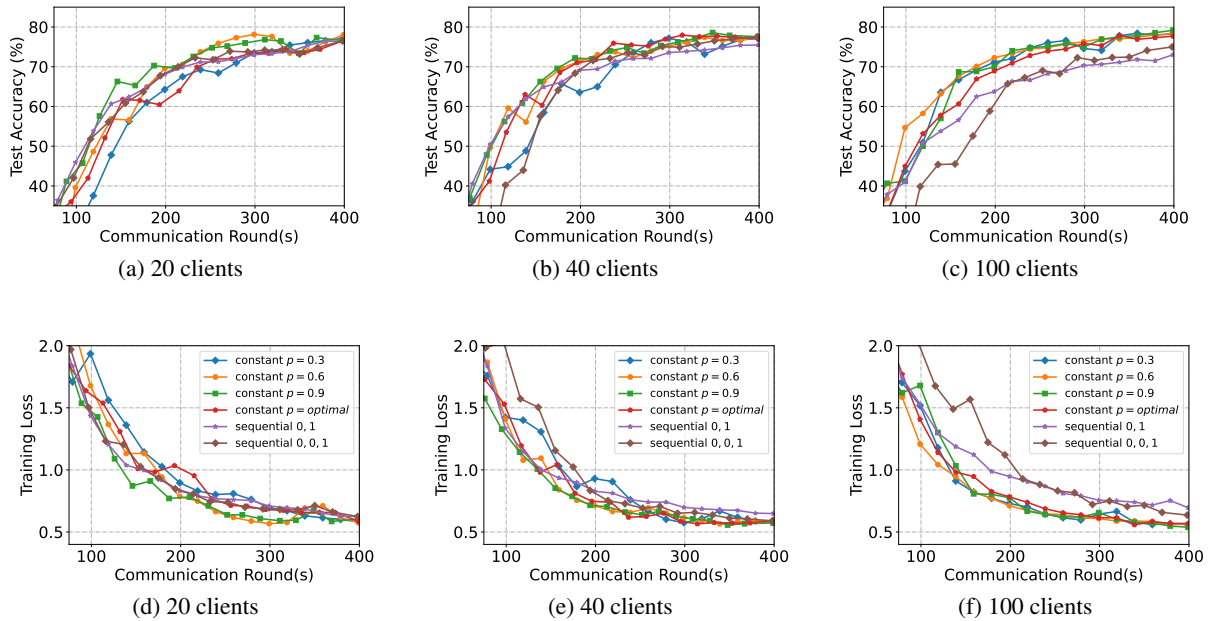


Figure 7. Comparison of different probability settings using training loss and test accuracy against the communication rounds for FedAMD by setting  $K = 5$ .



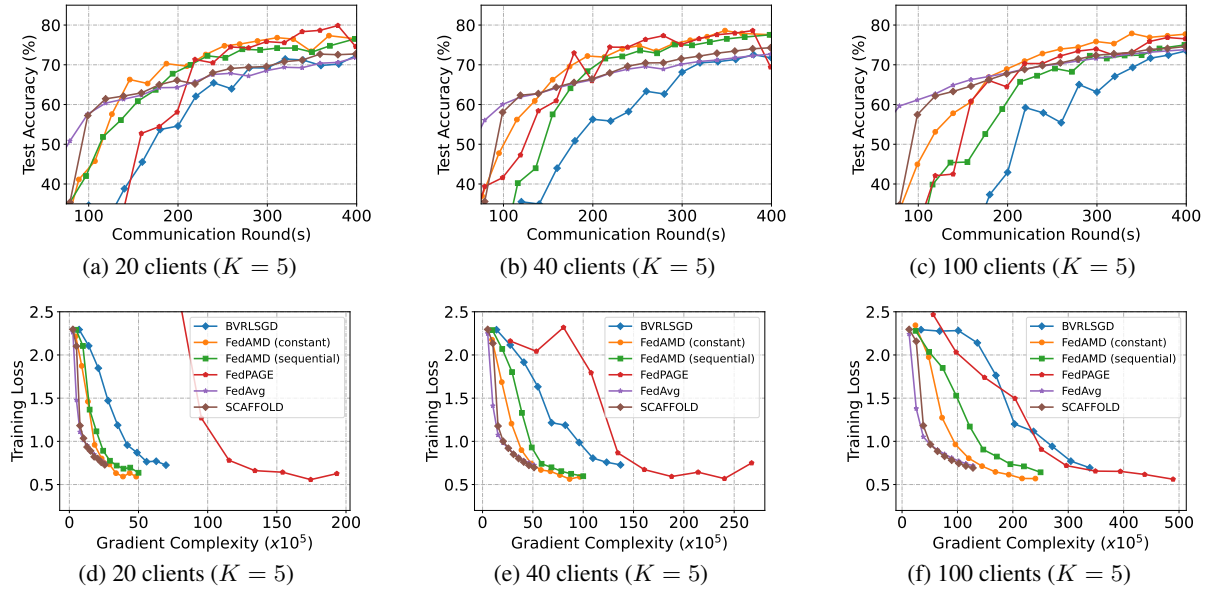


Figure 8. Comparison of different algorithms using test accuracy against the communication rounds and training loss against gradient complexity when the number of local iterations is 5 ( $K = 5$ ).

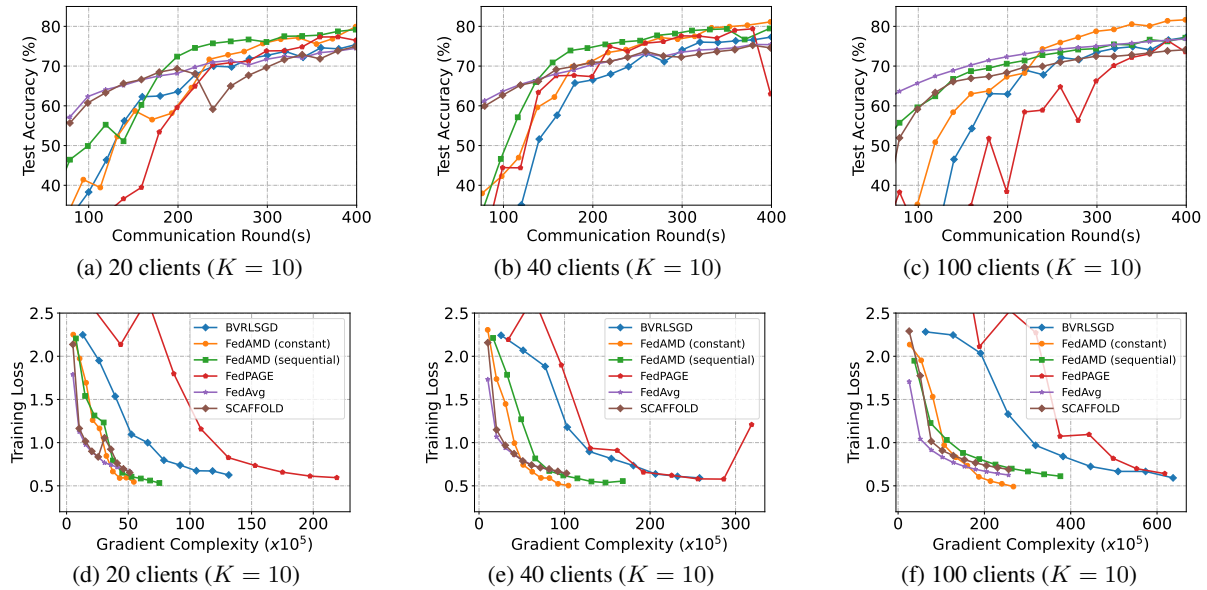


Figure 9. Comparison of different algorithms using test accuracy against the communication rounds and training loss against gradient complexity when the number of local iterations is 10 ( $K = 10$ ).

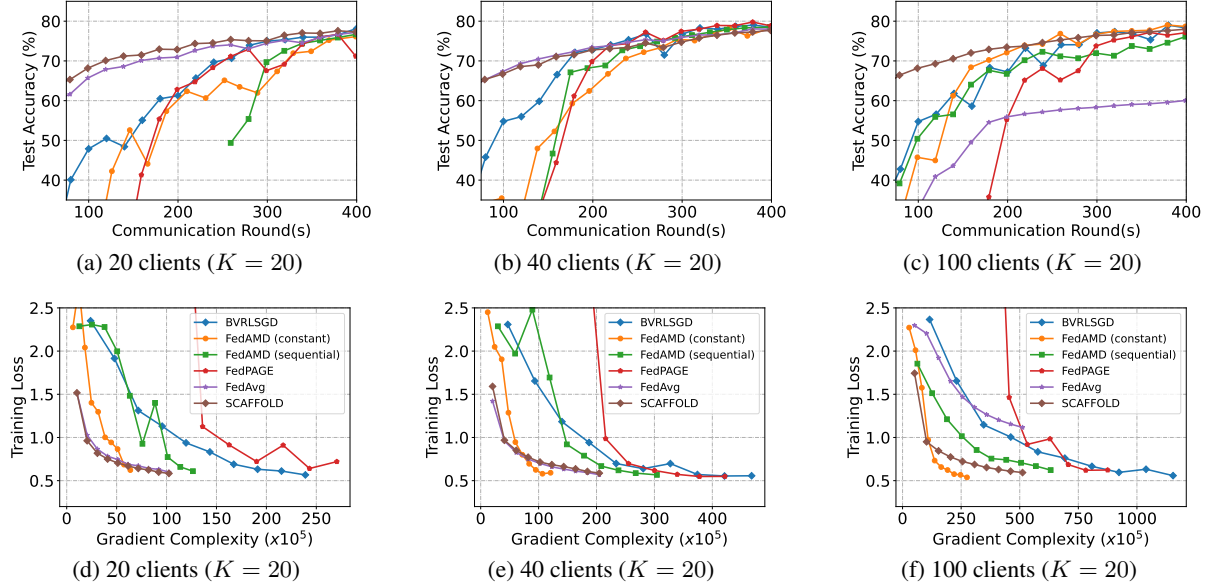


Figure 10. Comparison of different algorithms using test accuracy against the communication rounds and training loss against gradient complexity when the number of local iterations is 20 ( $K = 20$ ).

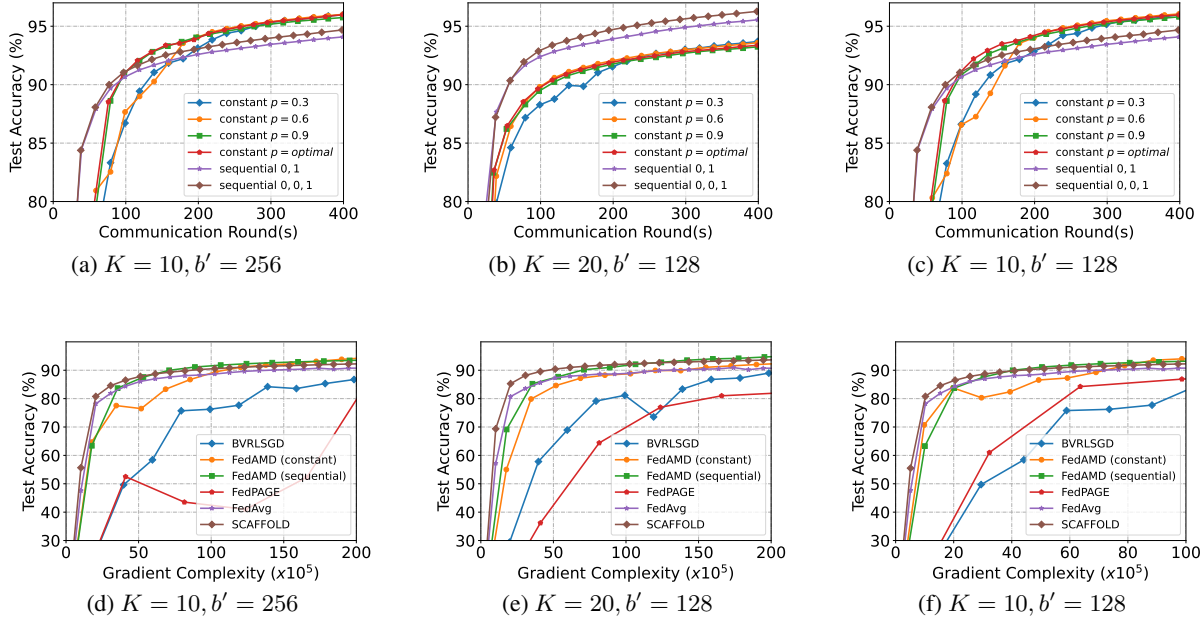


Figure 11. Comparison of different baselines and probability settings using test accuracy against communication rounds and gradient complexity, respectively.