# Towards Bias Correction of FedAvg over Nonuniform and Time-Varying Communications

Ming Xiang\*, Stratis Ioannidis\*, Edmund Yeh\*, Carlee Joe-Wong<sup>†</sup>, and Lili Su\*

Abstract—Federated learning (FL) is a decentralized learning framework wherein a parameter server (PS) and a collection of clients collaboratively trains a model via minimizing a global objective. Communication bandwidth is a scarce resource; in each round, the PS aggregates the updates from a subset of clients only. In this paper, we focus on non-convex minimization that is vulnerable to non-uniform and time-varying communication failures between the PS and the clients. Specifically, in each round t, the link between the PS and client i is active with probability  $p_i^t$ , which is unknown to both the PS and the clients. This arises when the channel conditions are heterogeneous across clients and are changing over time.

We show that when the  $p_i^t$ 's are not uniform, Federated Average (FedAvg) – the most widely adopted FL algorithm – fails to minimize the global objective. Observing this, we propose Federated Postponed Broadcast (FedPBC) which is a simple variant of FedAvg. It differs from FedAvg in that the PS postpones broadcasting the global model till the end of each round. We show that FedPBC converges to a stationary point of the original objective. The introduced staleness is mild and there is no noticeable slowdown. Both theoretical analysis and numerical results are provided. On the technical front, postponing the global model broadcasts enables implicit gossiping among the clients with active links at round t. Despite  $p_i^t$ 's are time-varying, we are able to bound the perturbation of the global model dynamics via the techniques of controlling the gossip-type information mixing errors.

#### I. INTRODUCTION

Federated learning (FL) is a distributed learning paradigm wherein a parameter server (PS) and a large collection of clients collaboratively learn a machine learning model with clients' local data undisclosed [1], [2] to the PS. The global objetives are often non-convex. Communication bandwidth is a scarce resource. In each round, the PS aggregates the updates from a subset of clients only – either proactively [1], [2] or passively [3]–[5]. A FL system is often deployed in a uncontrolled environment, wherein the channel conditions between the PS and the clients could be highly heterogeneous and time-varying [1]. To capture this, in this paper, we consider non-convex minimization that is vulnerable to nonuniform and time-varying link failures between the PS and the clients. Specifically, in each round, the link between the PS and client i is active with probability  $p_i^t$ , which is unknown to both the PS and the clients. A generic FL

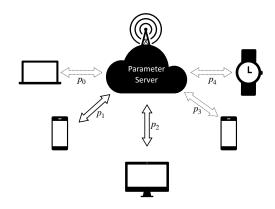


Fig. 1: A federated learning system with heterogeneous devices: Solid arrows indicate active links and dashed arrows are inactive links.

system of interest is illustrated in Fig. 1. To the best of our knowledge, the convergence of FL in the presence of non-uniform and time-varying communication is overall under-explored.

Our setup can be viewed as a special case of the general client unavailability, has received intensive attention recently [2]. Nevertheless, existing methods are not applicable to our problem. In the seminal works [1], [3], the PS chooses K clients either uniformly at random or proportionally to clients' local data volume. Neither of theses client selection methods is feasible when  $p_i^t$ 's are unknown and time-varying. In [2]–[4], [6], the PS waits for the K fastest responses. The correctness of their algorithms crucially relies on the fact that the response probability of each client is known. Ruan et al. [7] considered a generalized random client unavailability, yet required the response probability to be fixed. Timevarying response rates are also considered in [5], [8], [9]. For the methods in [5] to converge to stationary points, the response rates need to be "balanced" in the sense that either (1) the  $p_i^t$ 's are deterministic and satisfy the regularized participation, i.e.,  $\sum_{\tau=1}^P p_i^{t_0+\tau} = \mu$  for all clients at all  $t_0 \in \{0, P, 2P, \cdots\}$  where P is some carefully chosen integer; or (2)  $p_i^t$ 's are random and satisfy  $\mathbb{E}[p_i^t] = \mu$ for all clients and sufficiently many t. In contrast, we do not require such rate "balanceness". Perazzone et al. [8] analyzed the convergence of FedAvg under time-varying client participation rates. Nevertheless, they assumed (1) a uniform participation rate in each round, i.e.,  $p_i^t = p_i^t$  for any pair of clients, and (2) bounded stochastic gradient. Gu

<sup>\*</sup> Department of Electrical and Computer Engineering, Northeastern University, Boston, MA 02215, USA xiang.mi@northeastern.edu, {ioannidis,eyeh}@cce.neu.edu, l.su@northeastern.edu. † Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA cjoewong@andrew.cmu.edu. The work was supported in part by ARO Grant W911NF-23-2-0014 and National Science Foundation under Grants 2107062 and 2106891.

et al. [9] considered general client unavailability patterns for both strongly convex and non-convex global objectives. For non-convex objectives (which is our focus), they required that the consecutive unavailability rounds of a client to be deterministically upper bounded, which does not hold even for the simple uniform and time-invariant response rates. Moreover, they required the noise of the stochastic gradient to be uniformly upper bounded with probability 1.

#### Contributions. Our contributions is three-fold:

- We identify simple instances and show both analytically and numerically that when the p<sub>i</sub><sup>t</sup>'s are not uniform Federated Average (FedAvg) – the most widely adopted FL algorithm – fails to minimize the global objective.
- We propose Federated Postponed Broadcast (FedPBC). It differs from FedAvg in that the PS postpones broadcasting the global model till the end of each round. We show in Theorem 1 that, in expectation, FedPBC converges to a stationary point of the global objective. The correctness of our FedPBC neither impose any "balancedness" requirement on  $p_i^t$ 's nor require the stochastic gradients or their noises to be bounded. Moreover, compared with [5], [9], FedPBC works under a much relaxed bounded-dissimilarity assumption.
  - On the technical front, postponing the global model broadcasts enables implicit gossiping among the clients with active links. Hence, we mitigate the perturbation caused by non-uniform and time-varying  $p_i^t$  via the techniques of controlling information mixing errors.
- We validate our results empirically both on the counterexample and by using Synthetic (1,1) dataset [10]. The numerical results in the former show that FedPBC successfully corrects the bias when  $p_i^t$ 's are static but non-uniform (i.e.,  $p_i^t = p_i$ ) while FedAvg does not. In the latter, we further investigate *time-varying* link activation rates such that the responsive rates follow a uniform distribution and thus are bounded below. The results show FedPBC outperforms FedAvg.

### II. PROBLEM FORMULATION

A FL system consists of one central PS and m clients that collaboratively minimize

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} F(\boldsymbol{x}) = \frac{1}{m} \sum_{i \in [m]} F_i(\boldsymbol{x}), \qquad (1)$$

where  $F_i(x) = \mathbb{E}_{\xi_i \in \mathcal{D}_i}[\ell_i(x; \xi_i)]$  is the local objective,  $\mathcal{D}_i$  is the local distribution,  $\xi_i$  is a stochastic sample that client i has access to, and  $\ell_i$  is the local loss function. The loss function can be non-convex. We are interested in solving Eq. (1) over unreliable communication links between the PS and the clients. In each round t, the communication link between the PS and client i is active with probability  $p_i^t$ , which could be **time-varying** and is **unknown** to both the PS and the clients. We assume that  $p_i(t) \geq c$  for all t and all i, where  $c \in (0,1)$ .

# III. A CASE STUDY ON THE OBJECTIVE INCONSISTENCY OF FEDAVG

In this section, we use a simple example (a similar setup as in [11]) to illustrate FedAvg fails to minimize the global objective in Eq. (1) when  $p_i$ 's are not uniform. For completeness, we formally describe FedAvg in Algorithm 1. Notably, in Algorithm 1, all the clients (regardless of whether

# Algorithm 1: Federated Average (FedAvg) [1]

```
1 Input: T, x^0, s, \{\eta_t\}_{t=0,\cdots,T-1}
2 The PS and each client initialize parameter x^0;
  3 for t = 0, \dots, T - 1 do
                 /\star Let \mathcal{A}^t denote all the clients
                            with active communication
                            links.
                 The PS broadcasts x^t to each client;
                 for i \in [m] do
  5
                          Draw a fresh sample \xi_i^t;
  6
                         \begin{aligned} & \textbf{if} \ i \in \mathcal{A}^t \ \textbf{then} \\ & \boldsymbol{x}_i^{(t,0)} \leftarrow \boldsymbol{x}^t; \end{aligned}
  7
  8
                                  oldsymbol{x}_i^{(t,0)} \leftarrow oldsymbol{x}_i^t;
10
11
                          \begin{aligned} & \textbf{for } k = 0, \cdots, s-1 \textbf{ do} \\ & \boldsymbol{x}_i^{(t,k+1)} \leftarrow \boldsymbol{x}_i^{(t,k)} - \eta_t \nabla \ell_i(\boldsymbol{x}_i^{(t,k)}; \boldsymbol{\xi}_i^t); \end{aligned}
12
13
14
                         egin{aligned} oldsymbol{x}_i^{t+1} \leftarrow oldsymbol{x}_i^{(t,s)}; \ 	ext{Report } oldsymbol{x}_i^{t+1} 	ext{ to the PS;} \end{aligned}
15
16
17
                 /* On the PS.
                \begin{array}{c} \text{if } \mathcal{A}^t \neq \emptyset \text{ then} \\ \boldsymbol{x}^{t+1} \leftarrow \frac{1}{|\mathcal{A}^t|} \sum_{i \in \mathcal{A}^t} \boldsymbol{x}_i^{t+1}; \end{array}
18
19
20
                          oldsymbol{x}^{t+1} \leftarrow oldsymbol{x}^t:
21
22
                 end
23 end
```

the corresponding links are active or not) compute locally in Algorithm 1 in each round. This is *logically equivalent* to the usual setting where only clients in  $\mathcal{A}^t$  do the local steps because in line 20 the summation is taken over the clients in  $\mathcal{A}^t$ . Similar equivalence is observed in [5]. We present the FedAvg in the form of Algorithm 1 for ease of comparison with our FedPBC – an algorithmic fix to FedAvg for bias correction.

Let the local objective  $F_i(x) = \frac{1}{2} \|x - u_i\|_2^2$ , where  $u_i \in \mathbb{R}^d$  is an arbitrary vector. The corresponding global objective is thus

$$F(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^{m} F_i(\mathbf{x}) = \frac{1}{2m} \sum_{i=1}^{m} \|\mathbf{x} - \mathbf{u}_i\|_2^2,$$
 (2)

with unique minimizer  $x^* = \frac{1}{m} \sum_{i=1}^m u_i$ .

**Proposition 1.** Choose  $x^0 = \mathbf{0}$  and  $\eta_t = \eta \in (0,1)$  for all t. For a global objective as per Eq. (2), if  $p_i^t = p_i$  for all t,

under FedAvg with exact local gradients

$$\lim_{T \to \infty} \boldsymbol{x}^{T} = \sum_{i=1}^{m} \frac{p_{i} \boldsymbol{u}_{i} \left[ 1 + \sum_{j=2}^{m} \left( -1 \right)^{j+1} \frac{1}{j} \sum_{S \in \mathcal{B}_{j}} \prod_{z \in S} p_{z} \right]}{1 - \prod_{i=1}^{m} \left( 1 - p_{i} \right)},$$

where 
$$\mathcal{B}_{j} \triangleq \left\{ S \middle| S \subseteq [m] \setminus \{j\}, |S| = j - 1 \right\}.$$

The proof of Proposition 1 can be found in Appendix. It can be checked that if there exist  $i,i'\in[m]$  such that  $p_i\neq p_{i'}$ , then  $\lim_{t\to\infty} \boldsymbol{x}^t\neq\frac{1}{m}\sum_{i=1}^m \boldsymbol{u}_i\triangleq \boldsymbol{x}^*;$  when  $p_i=p$  for all  $i\in[m]$ , then  $\lim_{t\to\infty} \boldsymbol{x}^t=\boldsymbol{x}^*.$  In fact, the output of FedAvg may be arbitrarily away from  $\boldsymbol{x}^*$  depending on  $p_i$ 's and  $\boldsymbol{u}_i$ 's.

#### IV. ALGORITHM: FEDPBC

In this section, we propose FedPBC (*Federated Postponed Broadcast*, formally described in Algorithm 2) - a simple variant of FedAvg.

# **Algorithm 2:** FedPBC

```
1 Input: T, x^0, s, \{\eta_t\}_{t=0,\dots,T-1}
  2 The PS and each client initialize parameter x^0;
  3 for t = 0, \dots, T - 1 do
                /\star Let \mathcal{A}^t denote all the clients
                          with active communication
                          links;
  4
               for i \in [m] do
                        Draw a fresh sample \xi_i^t;
  5
                       egin{aligned} oldsymbol{x}_i^{(t,0)} &= oldsymbol{x}_i^t; \ \mathbf{for} \ k = 0, \cdots, s-1 \ \mathbf{do} \ oldsymbol{x}_i^{(t,k+1)} &= oldsymbol{x}_i^{(t,k)} - \eta_t 
abla \ell_i(oldsymbol{x}_i^{(t,k)}; oldsymbol{\xi}_i^t); \end{aligned}
  6
  8
  q
                       egin{aligned} oldsymbol{x}_i^{t+1} &= oldsymbol{x}_i^{(t,s)}; \ 	ext{Report } oldsymbol{x}_i^{t+1} 	ext{ to the PS;} \end{aligned}
10
11
12
                /* On the PS.
              \begin{array}{c} \textbf{if } \mathcal{A}^t \neq \emptyset \textbf{ then} \\ \boldsymbol{x}^{t+1} \leftarrow \frac{1}{|\mathcal{A}^t|} \sum_{i \in \mathcal{A}^t} \boldsymbol{x}_i^{t+1}; \end{array}
13
14
15
                       oldsymbol{x}^{t+1} \leftarrow oldsymbol{x}^t:
16
17
               Multi-cast x^{t+1} to each client i \in A^t;
18
                \begin{aligned} & \textbf{for} \ m \in \mathcal{A}^t \ \textbf{do} \\ & x_i^{t+1} \leftarrow \boldsymbol{x}^{t+1}; \end{aligned} 
19
20
21
                end
22 end
```

The key difference of FedPBC from FedAvg is that we postpone the global model broadcasts to  $\mathcal{A}^t$  till the end of each round. Postponing the global model broadcast introduces some staleness as the clients might start from different  $\boldsymbol{x}_i^t$  rather than  $\boldsymbol{x}^t$ . It turns out that such staleness helps in mitigating the bias caused by non-uniform link activation probabilities. Moreover, the staleness is mild and there is

no significant slowdown. Theoretical analysis and numerical results can be found in Sections V and VI, respectively.

Implicit gossiping among clients  $\mathcal{A}^t$ . From line 14 to line 22 of Algorithm 2, via the coordination of the PS, the clients in  $\mathcal{A}^t$  implicitly average their local updates with each other, i.e., there is implicit gossiping among the clients in  $\mathcal{A}^t$  at round t. Formally, we are able to construct a mixing matrix  $W^{(t)}$  as

$$W_{ij}^{(t)} = \begin{cases} \frac{1}{|\mathcal{A}^t|}, & \text{if } i, j \in \mathcal{A}^t; \\ 1, & \text{if } i = j \text{ and } \{i \notin \mathcal{A}^t\}; \\ 0, & \text{otherwise.} \end{cases}$$

The matrix is by definition doubly-stochastic and  $W^{(t)} = I$  when  $\mathcal{A}^t = \emptyset$  or  $|\mathcal{A}^t| = 1$ . We further note that this matrix can be time-varying even in expectation since the link activation probabilities  $p_i^t$ 's can be time-varying. As can be seen later, this mixing matrix bridges the gap between local and global model heterogeneity and establishes a consensus among different clients.

Let 
$$M^{(t)} := \mathbb{E}\left[\left(W^{(t)}\right)^2\right]$$
 and  $\mathbf{J} := \frac{1}{m}\mathbf{1}\mathbf{1}^{\top}$ . Define as  $\rho(t) := \lambda_2\left(M^{(t)}\right)$  and  $\rho := \max_t \rho(t)$ . (3)

**Lemma 1** (Ergodicity). Recall that  $p_i^t \geq c$  for some constant  $c \in (0,1)$ . For each  $t \geq 1$ , it holds that  $\rho \leq 1 - \frac{c^4[1-(1-c)^m]^2}{8}$ .

We defer the proof of Lemma 3 to Appendix. The following lemma will be used in the convergence analysis.

**Lemma 2.** For any matrix  $B \in \mathbb{R}^{d \times m}$ , it holds that

$$\mathbb{E}\left[\|B\left(\prod_{r=1}^t W^{(r)} - \mathbf{J}\right)\|_{\mathrm{F}}^2\right] \leq \rho^t \|B\|_F^2.$$

The proof of Lemma 2 follows the same outline as that in [12, Lemma]; it is deferred to Appendix.

Remark 1. In Algorithm 2, each client does local computations even if its communication link is not active. Continuous local updates appear to be crucial. Numerical examples in Section VI show that bias persists when only the active clients do local computations. We leave as a future direction on how to remove the bias while maintaining local computation.

#### V. Convergence Results

#### A. Assumptions

Before diving into our convergence results, we will introduce some assumptions, which are commented towards the end of this subsection.

**Assumption 1** (Smoothness). Each local gradient function  $\nabla \ell_i(\theta)$  is  $L_i$ -Lipschitz, i.e.,

$$\begin{split} \left\|\nabla \ell_i(\boldsymbol{x}_1) - \nabla \ell_i(\boldsymbol{x}_2)\right\|_2 &\leq L_i \left\|\boldsymbol{x}_1 - \boldsymbol{x}_2\right\|_2, \\ \textit{for all } \boldsymbol{x}_1, \boldsymbol{x}_2, \textit{ and } i \in [m]. \textit{ Let } L \triangleq \max_{i \in [m]} L_i. \end{split}$$

**Assumption 2** (Bounded Variance). Stochastic gradients at each client node  $i \in [m]$  are unbiased estimates of the true gradient of the local objectives, i.e.,

$$\mathbb{E}\left[\nabla \ell_i(\boldsymbol{x}_i^t) \mid \mathcal{F}^t\right] = \nabla F_i(\boldsymbol{x}_i^t),$$

and the variance of stochastic gradients at each client node  $i \in [m]$  is uniformly bounded, i.e.,

$$\mathbb{E}\left[\left\|\nabla \ell_i(\boldsymbol{x}) - \nabla F_i(\boldsymbol{x})\right\|_2^2\right] \leq \sigma^2,$$

where  $\mathcal{F}^t$  denotes the sigma algebra generated by all the randomness up to iteration t.

**Assumption 3.** There exists  $F^* \in \mathbb{R}$  such that  $F(x) \geq F^*$  for all  $x \in \mathbb{R}^d$ .

Assumption 4 (Bounded Inter-client Heterogeneity).

$$\frac{1}{m}\sum_{i=1}^{m}\left\|\nabla F_i(\boldsymbol{x}) - \nabla F(\boldsymbol{x})\right\|_2^2 \le \beta^2 \left\|\nabla F(\boldsymbol{x})\right\|_2^2 + \zeta^2.$$

Assumptions, 1, 2 and 3 are standard in FL analysis [10], [13], [14]. Assumption 4 captures the heterogeneity across different users, and it is a more relaxed version (e.g., than [10], [15], [16].) Notably, different from [9], we do not assume fresh data per local update, and the unbiasedness in Assumption 2 is imposed for global rounds only.

#### B. Results

In this section, we formally state our key lemmas and main theorem. All proofs can be found in the full version [17].

**Lemma 3** (Lemma 1 in [18]). For  $s \ge 1$ , we have for all  $x \in \mathbb{R}^d$ :

$$\left\| \sum_{k=0}^{s-1} \left[ \nabla \ell_i(\boldsymbol{x}^{(t,k)}) - \nabla \ell_i(\boldsymbol{x}^t) \right] \right\|_2 \le \kappa \eta \binom{s}{2} L_i \left\| \nabla \ell_i(\boldsymbol{x}^t) \right\|_2,$$

where  $\kappa \triangleq \max_i \frac{(1+\eta L_i)^s - 1 - s\eta L_i}{\binom{s}{2}(\eta L_i)^2}$ .

**Claim 1.** For any  $s \in \mathbb{N}$ ,  $\kappa$  is monotonic non-decreasing with respect to  $\eta > 0$ , where

$$\kappa \triangleq \frac{\left(1 + \eta L\right)^{s} - 1 - s\eta L}{\binom{s}{2} \left(\eta L\right)^{2}}.$$

**Remark 2.** Lemma 3 yields a simple upper bound on the perturbations incurred by multiple local steps. For the special case when s=1, we simply have  $\kappa=0$ . For  $s\geq 2$ , we always have  $\kappa\geq 1$ , and furthermore  $\kappa\leq \frac{e^c-1-c}{c^2/2}$ , when  $\eta\leq \frac{c}{sL}$ , which follows from Claim 1. In other words, we can treat  $\kappa$  as a constant as long as  $\eta$  is sufficiently small. Henceforward, for the special case s=1, we know that  $\kappa=0$  and treat  $\frac{\sqrt{2}}{\kappa sL}$  as  $\infty$ . In other words, it is removed from the step-size threshold set when s=1.

Let

$$\bar{\boldsymbol{x}}^t \triangleq \frac{1}{m} \sum_{i=1}^m \boldsymbol{x}_i^t. \tag{4}$$

**Lemma 4** (Descent Lemma). Suppose Assumptions 1, 2, and 4 hold, under a choice of the learning rate  $\eta \leq \frac{1}{2s}$ , the following property holds for  $t \geq 0$ :

$$\begin{split} & \mathbb{E}\left[F(\bar{\boldsymbol{x}}^{t+1}) - F(\bar{\boldsymbol{x}}^t) \mid \mathcal{F}^t\right] \\ & \leq \sigma^2 \eta^2 s^2 \left[\kappa^2 L^2 + 2L\left(\frac{1}{m} + \frac{\kappa^2 L^2}{4}\right)\right] + 3\xi^2 \eta^2 s^2 \mathfrak{C} \\ & - \left\{\frac{s\eta}{4} - 3\eta^2 s^2 \left(\beta^2 + 1\right) \mathfrak{C}\right\} \left\|\nabla F(\bar{\boldsymbol{x}}^t)\right\|_2^2 \\ & + \left\{\eta s L^2 + 3\eta^2 s^2 L^2 \mathfrak{C}\right\} \underbrace{\frac{1}{m} \sum_{i=1}^m \left\|\boldsymbol{x}_i^t - \bar{\boldsymbol{x}}^t\right\|_2^2,}_{consensus error.} \end{split}$$

where  $\mathfrak{C} \triangleq \kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)$ .

Remark 3. Lemma 4 can be proved via following the standard outline of SGD convergence analysis with non-convex functions and plugging in Lemma 3 to bound the perturbation arises from multiple local updates and non-fresh data per update. The consensus error term comes from Assumption 1 and enables us to connect our analysis of the aforementioned W matrix, where we borrow the insights from the analysis of gossiping algorithms. Formally, in matrix form, we use the following notions

$$\boldsymbol{X}^{(t)} = \begin{bmatrix} \boldsymbol{x}_1^t, \cdots, \boldsymbol{x}_m^t \end{bmatrix};$$

$$\boldsymbol{G}_0^{(t)} = \begin{bmatrix} s \nabla \ell_1(\boldsymbol{x}_1^{(t,0)}), \cdots, s \nabla \ell_m(\boldsymbol{x}_m^{(t,0)}) \end{bmatrix};$$

$$\boldsymbol{G}^{(t)} = \begin{bmatrix} \sum_{r=0}^{s-1} \nabla \ell_1(\boldsymbol{x}_1^{(t,r)}), \cdots, \sum_{r=0}^{s-1} \nabla \ell_m(\boldsymbol{x}_m^{(t,r)}) \end{bmatrix};$$

$$\nabla \boldsymbol{F}^{(t)} = \begin{bmatrix} \nabla F_1(\boldsymbol{x}_1^t), \cdots, \nabla F_m(\boldsymbol{x}_m^t) \end{bmatrix}.$$

Equivalently, we can write down the consensus error in matrix form.

$$\begin{split} \sum_{i=1}^{m} \left\| \bar{\boldsymbol{x}}^{t} - \boldsymbol{x}_{i}^{t} \right\|_{2}^{2} &= \left\| \boldsymbol{X}^{(t)} \left( \mathbf{I} - \mathbf{J} \right) \right\|_{\mathrm{F}}^{2} \\ &= \left\| \left( \boldsymbol{X}^{(t-1)} - \eta \boldsymbol{G}^{(t-1)} \right) \boldsymbol{W}^{(t-1)} \left( \mathbf{I} - \mathbf{J} \right) \right\|_{\mathrm{F}}^{2} \\ &= \eta^{2} \left\| \sum_{q=0}^{t-1} \boldsymbol{G}^{(q)} \left( \Pi_{l=q}^{t-1} \boldsymbol{W}^{(q)} - \mathbf{J} \right) \right\|_{\mathrm{F}}^{2}, \end{split}$$

where the last follows from the fact that all clients are initiated at the same weights.

**Lemma 5** (Consensus Error). Suppose the conditions in Lemma 4 are met, under a choice of the learning rate,  $\eta \leq \min\left\{\frac{1}{2s}, \frac{\sqrt{2}}{\kappa s L}, \frac{1-\sqrt{\rho}}{6\sqrt{2\rho}Ls^2}\right\}$ . The following property holds:

$$\frac{1}{mT} \sum_{t=0}^{T-1} \mathbb{E} \left[ \| \boldsymbol{X}^{(t)} \left( \mathbf{I} - \mathbf{J} \right)^{2} \|_{F} \right] 
\leq 6\eta^{2} s^{2} \sigma^{2} \left[ \frac{2\rho}{\left( 1 - \sqrt{\rho} \right)^{2}} + \frac{\rho}{1 - \rho} \right] + \frac{72\eta^{2} s^{4} \rho}{\left( 1 - \sqrt{\rho} \right)^{2}} \xi^{2} 
+ \frac{72 \left( \beta^{2} + 1 \right) \eta^{2} s^{4} \rho}{\left( 1 - \sqrt{\rho} \right)^{2}} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \| \nabla F(\bar{\boldsymbol{x}}^{t}) \|_{2}^{2} \right].$$

Now, we are ready to present our main theorem.

**Theorem 1.** Suppose all the assumptions hold, and choose a learning rate  $\eta = c_0 \sqrt{\frac{m}{sT}}$ , where  $c_0$  is a constant, for sufficiently large T such that

$$\begin{split} \eta \leq \min \left\{ \frac{1}{24 \left(\beta^2 + 1\right) \mathfrak{C}\left[1 + \frac{18 s^2 L^2 \rho}{\left(1 - \sqrt{\rho}\right)^2}\right] + \frac{144 \left(\beta^2 + 1\right) s^2 L^2 \rho}{\left(1 - \sqrt{\rho}\right)^2}}, \\ \frac{1}{2s}, \frac{\sqrt{2}}{\kappa s L}, \frac{1}{\rho s^3}, \frac{1 - \sqrt{\rho}}{6\sqrt{2\rho} L s^2} \right\}, \end{split}$$

the following property holds for Algorithm 2:

$$\begin{split} &\frac{1}{T}\sum_{k=0}^{T-1}\mathbb{E}\left[\left\|\nabla F(\bar{\boldsymbol{x}}^t)\right\|_2^2\right] \leq O\left(\frac{8F(\bar{\boldsymbol{x}}^0) - 8F^\star}{\sqrt{msT}}\right) \\ &+ \underbrace{O\left(16L\sqrt{\frac{s}{mT}}\sigma^2 + 8\sqrt{\frac{ms}{T}}\kappa^2L^2\left(1 + \frac{L}{2}\right)\sigma^2\right)}_{Stochastic \ gradient \ noise} \\ &+ \underbrace{O\left(24\sqrt{\frac{ms}{T}}\left[\mathfrak{C} + \frac{24L^2}{\left(1 - \sqrt{\rho}\right)^2}\right]\xi^2 + \frac{1728\mathfrak{C}L^2\xi^2}{\left(1 - \sqrt{\rho}\right)^2}\frac{ms}{T}\right)}_{Client \ drift \ error} \\ &+ O\left(\frac{144\rho}{\left(1 - \sqrt{\rho}\right)^2}\left(L^2 + 3L^2\mathfrak{C}\right)\sigma^2\frac{ms}{T}\right), \end{split}$$

Intermittent participation error

where 
$$\mathfrak{C} \triangleq \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right)$$
.

Remark 4. Here, we remark on Theorem 1:

- 1) On the structures. Except for the first term, the remained terms can be grouped into three parts: the noise introduced by stochastic gradient, and the errors due to client drift (heterogeneity) and intermittent participation, each scaling with a different rate. To control the errors, we need a sufficiently small learning rate η that meets all the conditions mentioned above.
- 2) On stationary points of F. Theorem 1 says that  $\bar{x}^t$  in FedPBC converges to a stationary point of F asymptotically. In other words, the bias will be corrected towards the end. In contrast, we show in Proposition 1 that  $\bar{x}^t$  in FedAvg converges to a point that could be arbitrarily far away from the true optimum depending on  $p_i^t$  and data heterogeneity.
- 3) On the role of the activation lower bound c. It has been shown in Lemma 1 that  $\rho \leq 1 \frac{c^4[1-(1-c)^m]^2}{8}$ . A greater c leads to a smaller  $\rho$  and thus a tighter bound on  $\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}[\|\nabla F\left(\bar{\boldsymbol{x}}^t\right)\|_2]$ . Note that FedPBC reduces to FedAvg with full-client participation, i.e., when c=1. In that case, our convergence rate becomes

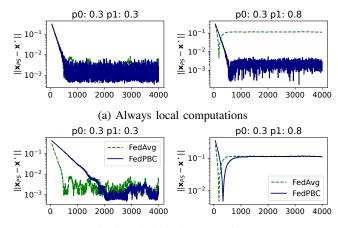
$$O\left(\frac{1}{\sqrt{msT}} + \sqrt{\frac{ms}{T}} + \frac{ms}{T}\right),\tag{5}$$

- which matches the FedAvg literature (see e.g., in [11]). We further note that because  $\frac{\kappa}{s}$  can be treated as a constant, the order of convergence rate does not change.
- 4) On linear speedup. It is trivial to see that the first two terms in Eq. (5) dominate when T is sufficiently large (e.g.,  $T \ge c_1 m^3 s^3$ , where  $c_1$  is some positive constant.) We shall see linear speedup w.r.t. the first term; however, the second term ultimately dominates all. Thus, it is unlikely that our algorithm achieves linear speedup, which is consistent with FedAvg literature, see e.g., in [3].

#### VI. NUMERICAL EXPERIMENTS

In this section, we present the numerical evaluations of the proposed algorithm and FedAvg. In each round, the PS will send an update request to each client. Client i will respond with probability  $p_i$ , which is unknown to both the PS and clients. This simulates unstable communications.

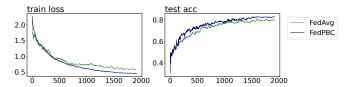
**Counterexample.** Here, we have m=100 clients, each doing 30-steps local computations, communicating for 4000 rounds, and holding a local loss function  $F_i(\boldsymbol{x}_i) = \frac{1}{2} \|\boldsymbol{x}_i - \boldsymbol{u}_i\|_2^2$ , where  $\boldsymbol{x}_i, \boldsymbol{u}_i \in \mathbb{R}^{100}$ ,  $\boldsymbol{u}_i \sim \mathcal{N}\left((i/1000)\mathbf{1}, 0.01\mathbf{I}\right)$ , and  $\boldsymbol{x}_i^0 = \mathbf{0}$  for all  $i \in [m]$ . The learning rate  $\eta = 0.0003$ . In addition, we let the first 50 clients respond with probability  $p_0$ , whereas the second half with  $p_1$  (to be specified later.)



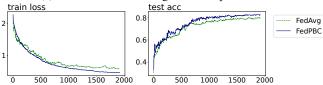
(b) Sampled local computations

Fig. 2: Distance to the optimum  $||x_{PS} - x^*||_2$  in the counterexample in logarithmic scale.

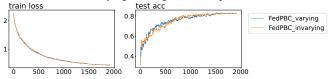
For ease of presentation, we plot the distance to the optimum  $\|x_{PS} - x^*\|_2$  after the first 50 communication rounds in Fig. 2, where  $x_{PS} \triangleq x$  in Algorithm 2. As illustrated in Fig. 2a, FedPBC is unbiased and converges to the global optimum  $x^* \triangleq \frac{1}{m} \sum_{i=1}^m u_i$  in all the combinations of  $p_0$  and  $p_1$ , matching our analysis, while FedAvg will instead converge to a different point observed from  $\|x_{PS} - x^*\|_2$  when  $p_0 \neq p_1$ . When  $p_0 = p_1$ , the two algorithms will converge to the same point, which matches our analysis. In a sharp contrast, if we let only the sampled clients do local computations, the bias persists, which we leave as a future direction.



(a) Time-invariant heterogeneous responsive rates.



(b) Time-varying heterogeneous responsive rates.



(c) FedPBC evaluations under time-invariant and time-varying responsive rates.

Fig. 3: Synthetic (1,1) evaluations.

**Synthetic** (1,1) **data**. In this simulation, we first follow [10] and construct Synthetic (1,1) dataset as follows: we generate samples  $(X_i,Y_i)$  for each client i according to the model  $y=\arg\max\left(\operatorname{softmax}\left(Wx+b\right)\right)$ , where  $x\in\mathbb{R}^{60}$ ,  $W\in\mathbb{R}^{10\times60},\ b\in\mathbb{R}^{10}$ . To characterize the non-i.i.d. data, we let  $W_i\sim N\left(u_i,1\right),\ b_i\sim N\left(u_i,1\right),\ u_i\sim N\left(0,\alpha=1\right),$  and  $x_i\sim \mathcal{N}\left(v_i,\Sigma\right)$ , where the covariance matrix is diagonal with  $\sum_{j,j}=j^{-1.2}$ . Each element in the mean vector  $v_i$  is drawn from  $N\left(B_i,1\right)$ , where  $B_i\sim (0,\beta=1)$ .

For the non-uniform link activation probabilities  $p_i$ s, we consider two scenarios:

- 1) Time-invariant heterogeneous rates. Let  $p_i^t=p_i^0=0.05$  for  $1\leq i\leq m/2$  and  $p_j^t=p_j^0=0.9$  for  $(m/2)+1\leq j\leq m$ . In other words, we have two groups of clients, one responding with probability  $p_i^0=0.05$ , while the other one with probability  $p_i^0=0.9$ ;
- 2) Time-varying heterogeneous rates. A uniformly distributed random variable, which is independent across clients and communication rounds, is imposed on each responsive rate per communication round. Formally, let  $p_i^t = p_i^0 + X_i^t \text{ and } p_j^t = p_j^0 + X_j^t, \text{ where } X_i^t, X_j^t \sim \mathcal{U}\left(-0.02, 0.02\right) \text{ for } 1 \leq i \leq (m/2) \text{ and } (m/2) + 1 \leq j \leq m.$  This ensures  $c \triangleq \min_{t \in [T], i \in [m]} p_i^t = 0.03$ .

The other auxiliary hyper-parameters are set as: client size m=30, a constant learning rate  $\eta_0$  tuned from  $\{0.1,0.5,0.01,\ldots,0.001,0.005\}$ , batch size: 100, local computation rounds: 25 for each  $i\in[m]$ , communication rounds: 1900.

Fig. 3a and Fig. 3b show that FedPBC consistently outperforms FedAvg. Moreover, Fig. 3c says that FedPBC converges to the same optimum in either settings, showing its ability to rectify the bias.

#### REFERENCES

- B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [2] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings et al., "Advances and open problems in federated learning," Foundations and Trends® in Machine Learning, vol. 14, no. 1–2, pp. 1–210, 2021.
- [3] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *International Conference on Learning Representations*, 2020. [Online]. Available: https://openreview.net/forum?id=HJxNAnVtDS
- [4] C. Philippenko and A. Dieuleveut, "Bidirectional compression in heterogeneous settings for distributed or federated learning with partial participation: tight convergence guarantees," *arXiv preprint arXiv:2006.14591*, 2020.
- [5] S. Wang and M. Ji, "A unified analysis of federated learning with arbitrary client participation," in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available: https://openreview.net/ forum?id=qSs7C7c4G8D
- [6] D. Jhunjhunwala, P. Sharma, A. Nagarkatti, and G. Joshi, "Fedvarp: Tackling the variance due to partial client participation in federated learning," in *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, ser. Proceedings of Machine Learning Research, J. Cussens and K. Zhang, Eds., vol. 180. PMLR, 01–05 Aug 2022, pp. 906–916. [Online]. Available: https://proceedings.mlr.press/v180/jhunjhunwala22a.html
- [7] Y. Ruan, X. Zhang, S.-C. Liang, and C. Joe-Wong, "Towards flexible device participation in federated learning," in *International Conference* on Artificial Intelligence and Statistics. PMLR, 2021, pp. 3403–3411.
- [8] J. Perazzone, S. Wang, M. Ji, and K. S. Chan, "Communication-efficient device scheduling for federated learning using stochastic optimization," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*. IEEE, 2022, pp. 1449–1458.
- [9] X. Gu, K. Huang, J. Zhang, and L. Huang, "Fast federated learning in the presence of arbitrary device unavailability," *Advances in Neural Information Processing Systems*, vol. 34, pp. 12052–12064, 2021.
- [10] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Proceedings of Machine Learning and Systems*, vol. 2, pp. 429–450, 2020.
- [11] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," *Advances in neural information processing systems*, vol. 33, pp. 7611–7623, 2020.
- [12] J. Wang, A. K. Sahu, G. Joshi, and S. Kar, "Matcha: A matching-based link scheduling strategy to speed up distributed optimization," *IEEE Transactions on Signal Processing*, vol. 70, pp. 5208–5221, 2022.
- [13] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [14] X. Yuan and P. Li, "On convergence of fedprox: Local dissimilarity invariant bounds, non-smoothness and beyond," in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available: https://openreview.net/forum?id=\_33ynl9VgCX
- [15] H. Yu, S. Yang, and S. Zhu, "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*, vol. 33, no. 01, 2019, pp. 5693–5700.
- [16] J. Wang and G. Joshi, "Cooperative sgd: A unified framework for the design and analysis of local-update sgd algorithms," *The Journal of Machine Learning Research*, vol. 22, no. 1, pp. 9709–9758, 2021.
- [17] M. Xiang, S. Ioannidis, E. Yeh, C. Joe-Wong, and L. Su, "Towards bias correction of fedavg over nonuniform and time-varying communications," arXiv preprint arXiv:2306.00280, 2023.
- [18] L. Su, J. Xu, and P. Yang, "Federated learning in the presence of adversarial client unavailability," arXiv preprint arXiv:2305.19971, 2023.
- [19] M. Jerrum and A. Sinclair, "Conductance and the rapid mixing property for markov chains: the approximation of permanent resolved," in *Proceedings of the twentieth annual ACM symposium on Theory of* computing, 1988, pp. 235–244.

**Proof of Proposition 1.** At each client  $i \in A^t$ , we have

$$\mathbf{x}_{i}^{(t,k+1)} = (1-\eta)^{k+1} \mathbf{x}^{t} + \eta \mathbf{u}_{i} \left[ \sum_{r=0}^{k} (1-\eta)^{r} \right].$$

Using the convention that  $\frac{0}{0} = 0$ , we get

$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^{t} \mathbf{1}_{\{\mathcal{A}_{t} = \emptyset\}} + (1 - \eta)^{s} \boldsymbol{x}^{t} \mathbf{1}_{\{\mathcal{A}^{t} \neq \emptyset\}} + \frac{\eta \sum_{i \in \mathcal{A}^{t}} \boldsymbol{u}_{i} \left[\sum_{r=0}^{s-1} (1 - \eta)^{r}\right] \mathbf{1}_{\{\mathcal{A}^{t} \neq \emptyset\}}}{|\mathcal{A}^{t}|}$$

$$= \left[\mathbf{1}_{\{\mathcal{A}^{t} = \emptyset\}} + (1 - \eta)^{s} \mathbf{1}_{\{\mathcal{A}^{t} \neq \emptyset\}}\right] \boldsymbol{x}^{t} + \eta \left[\sum_{k=0}^{s-1} (1 - \eta)^{k}\right] \frac{\mathbf{1}_{\{\mathcal{A}^{t} \neq \emptyset\}}}{|\mathcal{A}^{t}|} \sum_{i \in \mathcal{A}^{t}} \boldsymbol{u}_{i}.$$

Let  $\eta_t = \eta$  for all t. Since  $p_i^t = p_i$  for all  $i \in [m]$ ,

$$\mathbb{E}\left[rac{1}{|\mathcal{A}^t|}\sum_{i\in\mathcal{A}^t}oldsymbol{u}_i\Big|\mathcal{A}^t
eq\emptyset
ight]=\mathbb{E}\left[rac{1}{|\mathcal{A}^1|}\sum_{i\in\mathcal{A}^1}oldsymbol{u}_i\Big|\mathcal{A}^1
eq\emptyset
ight]$$

holds for all t. Taking expectation w.r.t.  $A^t$ , we get

$$\boldsymbol{x}^{t+1} = \left[ \mathbb{P} \left\{ \mathcal{A}^t = \emptyset \right\} + (1 - \eta)^s \, \mathbb{P} \left\{ \mathcal{A}^t \neq \emptyset \right\} \right] \boldsymbol{x}^t + \eta \left[ \sum_{k=0}^{s-1} (1 - \eta)^k \right] \mathbb{E} \left[ \frac{\sum_{i \in \mathcal{A}^t} \boldsymbol{u}_i}{|\mathcal{A}^t|} \middle| \mathcal{A}^t \neq \emptyset \right] \mathbb{P} \left\{ \mathcal{A}^t \neq \emptyset \right\}$$
$$= \left( 1 - \mathbf{a}^{t+1} \right) \mathbb{E} \left[ \frac{1}{|\mathcal{A}^1|} \sum_{i \in \mathcal{A}^1} \boldsymbol{u}_i \middle| \mathcal{A}^1 \neq \emptyset \right],$$

where we use the fact that  $x^0 = 0$ , and

$$a = \prod_{i=1}^{m} (1 - p_i) + [1 - \prod_{i=1}^{m} (1 - p_i)] (1 - \eta)^s.$$

Since a < 1, we get  $\lim_{t \to \infty} 1 - a^{t+1} = 1$ . Let  $X_i = \mathbf{1}_{\{i \in \mathcal{A}^1\}}$  for each  $i \in [m]$ . In sequel, we alternatively state the event  $\sum_{i=1}^m X_i \neq 0$  as  $\mathcal{A}^1 \neq \emptyset$  since they are equivalent.

$$\mathbb{E}\left[\frac{\sum_{i \in \mathcal{A}^1} \mathbf{u}_i}{|\mathcal{A}^1|} \middle| \mathcal{A}^1 \neq \emptyset\right] = \mathbb{E}\left[\frac{\sum_{i=1}^m X_i \mathbf{u}_i}{\sum_{i=1}^m X_i} \middle| \mathcal{A}^1 \neq \emptyset\right]$$
$$= \mathbb{E}\left[\sum_{i=1}^m \frac{X_i}{\sum_{i=1}^m X_i} \mathbf{u}_i \middle| \mathcal{A}^1 \neq \emptyset\right] = \sum_{i=1}^m \mathbf{u}_i \mathbb{E}\left[\frac{X_i}{\sum_{j=1}^m X_j} \middle| \mathcal{A}^1 \neq \emptyset\right].$$

Using the convention that  $\frac{0}{0} = 0$ , we know that

$$\mathbb{E}\left[\frac{X_i}{\sum_{j=1}^{M} X_j} \middle| \sum_{j=1}^{M} X_j \neq 0\right]$$

$$= \frac{\mathbb{E}\left[\frac{X_i}{\sum_{i=1}^{M} X_i} \middle| \mathcal{A}^1 \neq \emptyset\right] \mathbb{P}\left\{\mathcal{A}^1 \neq \emptyset\right\} + 0 \times \mathbb{P}\left\{\mathcal{A}^1 = \emptyset\right\}}{\mathbb{P}\left\{\mathcal{A}^1 \neq \emptyset\right\}}$$

$$= \frac{1}{1 - \prod_{i=1}^{m} (1 - p_i)} \mathbb{E}\left[\frac{X_i}{\sum_{j=1}^{m} X_j}\right].$$

Additionally,

$$\mathbb{E}\left[\frac{X_i}{\sum_{i=1}^m X_i}\right] = \mathbb{P}\left\{X_i = 1\right\} \mathbb{E}\left[\frac{X_i}{\sum_{j=1}^m X_j} \middle| X_i = 1\right]$$

$$+ \mathbb{P}\left\{X_i = 0\right\} \mathbb{E}\left[\frac{X_i}{\sum_{j=1}^m X_j} \middle| X_i = 0\right]$$

$$= p_i \mathbb{E}\left[\frac{1}{1 + \sum_{j \in [m] \setminus \{i\}} X_j} \middle| X_i = 1\right]$$

$$= p_i + \sum_{j=2}^m (-1)^{j+1} \frac{p_i}{j} \sum_{S \in \mathcal{B}_i} \prod_{z \in S} p_z,$$

where  $\mathcal{B}_j \triangleq \left\{S \middle| S \subseteq [m] \setminus \{i\}, |S| = j-1\right\}$ , and the last follows from the definition of a binomial distribution and can be seen through inspection of the terms.

**Proof of Lemma 1.** For ease of exposition, in this proof we drop the time index.

We first get the explicit expression for  $\mathbb{E}\left[W_{ij'}^2 \mid \mathcal{A} \neq \emptyset\right]$ . For  $j' \neq j$ , we have

$$W_{jj'}^{2} = \sum_{k=1}^{m} W_{jk} W_{j'k}$$

$$= W_{jj} W_{j'j} + W_{jj'} W_{j'j'} + \sum_{k \in [m] \setminus \{j,j'\}} W_{jk} W_{j'k}.$$

When  $k \neq j$  and  $k \neq j'$ , we have

$$W_{jk}W_{j'k} = \frac{1}{|\mathcal{A}|^2} \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}} \mathbf{1}_{\{k \in \mathcal{A}\}}.$$

In addition, we have

$$W_{jj}W_{j'j} = \frac{1}{|\mathcal{A}|} \left( 1 - \mathbf{1}_{\{j \in \mathcal{A}\}} \right) \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}}$$
$$+ \frac{1}{|\mathcal{A}|^2} \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}},$$

and

$$W_{j'j'}W_{jj'} = \frac{1}{|\mathcal{A}|} \left( 1 - \mathbf{1}_{\{j' \in \mathcal{A}\}} \right) \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}}$$
$$+ \frac{1}{|\mathcal{A}|^2} \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}}.$$

Thus,

$$W_{jj'}^{2} = \sum_{k=1}^{m} W_{jk} W_{j'k}$$

$$= \frac{1}{|\mathcal{A}|} \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}} + \frac{1}{|\mathcal{A}|} \left( 1 - \mathbf{1}_{\{j \in \mathcal{A}\}} \right) \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}}$$

$$+ \frac{1}{|\mathcal{A}|} \left( 1 - \mathbf{1}_{\{j \in \mathcal{A}\}} \right) \mathbf{1}_{\{j \in \mathcal{A}\}} \mathbf{1}_{\{j' \in \mathcal{A}\}}.$$

For j=j', we have  $W_{jj}^2=\frac{1}{|\mathcal{A}|}\mathbf{1}_{\{j\in\mathcal{A}\}}+\left(1-\mathbf{1}_{\{j\in\mathcal{A}\}}\right)$ . Taking expectation, we get

$$\mathbb{E}\left[W_{jj}^{2} \mid \mathcal{A}^{t} \neq \emptyset\right] = \mathbb{E}\left[\frac{1}{|\mathcal{A}|} \mathbf{1}_{\{j \in \mathcal{A}\}} + \left(1 - \mathbf{1}_{\{j \in \mathcal{A}\}}\right)\right]$$
$$= \mathbb{E}\left[\frac{1}{1 + |\mathcal{A}\setminus\{j\}|}\right] p_{j} + 1 \cdot (1 - p_{j}).$$

Note that  $A \setminus \{j\}$  is random and could be empty.

Let  $X_i = \mathbf{1}_{\{i \in \mathcal{A}\}}$ . We have

$$\mathbb{E}\left[\frac{1}{1+|\mathcal{A}\setminus\{j\}|}\right] = \mathbb{E}\left[\frac{1}{1+\sum_{i\in[m]\setminus\{j\}}X_i}\right] = \int_0^1 \prod_{k\neq j} \left[(1-p_k) + p_k s\right] ds$$
$$\geq \int_0^1 \prod_{k\neq j} \left[(1-p_k)s + p_k s\right] ds = \frac{1}{m}.$$

Thus,  $\mathbb{E}\left[W_{jj}^2 \mid \mathcal{A} \neq \emptyset\right] \geq \frac{1}{m}p_j + (1-p_j) \geq \frac{1}{m}$ . Similarly,

$$\mathbb{E}\left[W_{jj'}^2 \mid \mathcal{A} \neq \emptyset\right] \geq p_j p_{j'} \mathbb{E}\left[\frac{1}{2 + \sum_{k \in [m] \setminus \{j,j'\}} X_k}\right] \geq \frac{p_j p_{j'}}{m} \geq \frac{c^2}{m}.$$

Then,

$$M_{jj'} = \mathbb{E}\left[W_{jj'}^2\right] = \mathbb{E}\left[W_{jj'}^2 \mid \mathcal{A} \neq \emptyset\right] \mathbb{P}\left\{\mathcal{A} \neq \emptyset\right\}$$
$$+ \mathbb{E}\left[W_{jj'}^2 \mid \mathcal{A} = \emptyset\right] \mathbb{P}\left\{\mathcal{A} = \emptyset\right\}$$
$$\geq \frac{c^2}{m} \left[1 - (1 - c)^m\right].$$

$$M_{jj} = \mathbb{E}\left[W_{jj}^2\right] = \mathbb{E}\left[W_{jj}^2 \mid \mathcal{A} \neq \emptyset\right] \mathbb{P}\left\{\mathcal{A} \neq \emptyset\right\}$$

$$+ \mathbb{E}\left[W_{jj}^2 \mid \mathcal{A} = \emptyset\right] \mathbb{P}\left\{\mathcal{A} = \emptyset\right\}$$

$$\geq \frac{1}{m}\left[1 - (1 - c)^m\right] + (1 - c)^m$$

$$\geq \frac{1}{m} \geq \frac{c^2}{m}\left[1 - (1 - c)^m\right].$$

We first show that  $\rho(t) = \lambda_2(M)$ . We denote by  $\lambda_i$  and  $v_i$  the non-increasing eigenvalues and the associated eigenvectors of matrix M for  $i \in [m]$  with  $\lambda_1 = 1$  and  $v_1 = \frac{1}{\sqrt{m}} \mathbf{1}$ . By spectral decomposition

$$M - \frac{1}{m} \mathbf{1} \mathbf{1}^{\top} = \sum_{i=1}^{m} \lambda_i v_i v_i^{\top} - \frac{1}{m} \mathbf{1} \mathbf{1}^{\top} = \sum_{i=2}^{m} \lambda_i v_i v_i^{\top},$$

showing  $\rho(t) = \lambda_2$ .

Next, we show that a Markov chain with M as the transition matrix is ergodic. This is indeed true as the chain is (1) irreducible:  $M_{jj'} \geq \frac{c^2}{m} \left[1 - \left(1 - c\right)^m\right] > 0$  for  $j, j' \in [m]$  and (2) aperiodic (it has self-loops.) Moreover, it has a stationary distribution  $\pi = \frac{1}{m} \mathbf{1}^{\top}$ . Furthermore, this irreducible Markov chain is reversible since the following property is satisfied for all the states  $\pi_i M_{ij} = \pi_j M_{ji}$ .

Following [19], the conductance of reversible Markov chain with underlying graph  $\mathcal G$  is defined by  $\Phi(\mathcal G)=\min_{\sum_{i\in\mathcal S},j\notin\mathcal S} \frac{w_{ij}}{\sum_{i\in\mathcal S}\pi_i}$ , where the vertices of the graph are the states of the M Markov chain, and for each pair  $i,j\in\mathcal V$ , the edge weight  $w_{ij}=M_{ij}\pi_i=M_{ji}\pi_j$ . From Cheeger's inequality, we know that  $\frac{1-\lambda_2}{2}\leq\Phi(\mathcal G)\leq\sqrt{2(1-\lambda_2)}$ , where  $\lambda_2$  is the second largest eigenvalue of M. It remains to bound  $\Phi(\mathcal G)$ ,

$$\Phi(\mathcal{G}) = \min_{\sum_{i \in \mathcal{S}} \pi_i \leq \frac{1}{2}} \frac{\pi_i \sum_{i \in \mathcal{S}, j \notin \mathcal{S}} M_{ij}}{\sum_{i \in \mathcal{S}} \pi_i}$$

$$\geq \frac{\left(\frac{c}{m}\right)^2 \left[1 - \left(1 - c\right)^m\right] |\mathcal{S}| |\bar{\mathcal{S}}|}{\frac{|\mathcal{S}|}{m}} = \frac{c^2 \left[1 - \left(1 - c\right)^m\right] |\bar{\mathcal{S}}|}{m} |\bar{\mathcal{S}}|,$$

where the inequality follows from (1)  $\mathcal{G}$  is fully-connected (2)  $M_{jj'} \geq \frac{c^2}{m} [1 - (1-c)^m]$  for  $j, j' \in [m]$ . Meanwhile,  $|\bar{\mathcal{S}}| = m - |\mathcal{S}| \geq \frac{m}{2}$ . Plug it back in, we get

$$\Phi(\mathcal{G}) \ge \frac{c^2 \left[1 - (1 - c)^m\right]}{m} \left| \bar{\mathcal{S}} \right| \ge \frac{c^2 \left[1 - (1 - c)^m\right]}{2}.$$

Thus, 
$$\rho(t) = \lambda_2 \le 1 - \frac{\Phi^2(\mathcal{G})}{2} \le 1 - \frac{c^4[1 - (1 - c)^m]^2}{8}$$
.

**Proof of Lemma 2.** Similar to the proof in [12], let us define  $A_{r,t} \triangleq \prod_{l=r}^t W^{(\ell)} - \mathbf{J}$  and use  $\boldsymbol{b}_i^{\top}$  to denote the *i*-th row vector of B. Since for  $\ell \in \mathbb{N}$ , we have  $(W^{(\ell)})^{\top} = W^{(\ell)}$  and  $W^{(\ell)}\mathbf{J} = \mathbf{J}W^{(\ell)} = \mathbf{J}$ . Thus, one can obtain

$$A_{1,t} = \prod_{\ell=1}^{t} \left( W^{(\ell)} - \mathbf{J} \right) = A_{1,t-1} \left( W^{(t)} - \mathbf{J} \right).$$

Then, by taking expectation w.r.t.  $W^{(t)}$ , we have

$$\begin{split} & \mathbb{E}_{W^{(t)}} \left[ \| B A_{1,t} \|_{\mathrm{F}}^2 \right] \\ & = \sum_{i=1}^d \mathbb{E}_{W^{(t)}} \left[ \left\| b_i^\top A_{1,t} \right\|_2^2 \right] \\ & = \sum_{i=1}^d \mathbb{E}_{W^{(t)}} \left[ b_i^\top A_{1,t-1} \left( \left( W^{(t)} \right)^\top W^{(t)} - \mathbf{J} \right) A_{1,t-1}^\top b_i \right] \\ & = \sum_{i=1}^d b_i^\top A_{1,t-1} \mathbb{E}_{W^{(t)}} \left[ \left( (W^{(t)})^\top W^{(t)} - \mathbf{J} \right) \right] A_{1,t-1}^\top b_i. \end{split}$$

Let  $C_t = \mathbb{E}_{W^{(t)}} \left[ (W^{(t)})^\top W^{(t)} - \mathbf{J} \right]$  and  $v_i = A_{1,t-1}^\top b_i$ , then

$$\mathbb{E}_{W(t)} \left[ \|BA_{1,t}\|_{\mathrm{F}}^{2} \right] = \sum_{i=1}^{d} v_{i}^{\top} C_{t} v_{i} \leq \sigma_{\max} \left( C_{t} \right) \sum_{i=1}^{d} v_{i}^{\top} v_{i}$$
$$\leq \rho \|BA_{1,t-1}\|_{\mathrm{F}}^{2}.$$

Repeat the above procedures, since  $W^{(\ell)}$ 's are independent matrices, we have

$$\mathbb{E}\left[\|BA_{1,t}\|_{\mathrm{F}}^{2}\right] = \mathbb{E}_{W^{(1)}}\left[\mathbb{E}_{W^{(2)}}\left[\cdots\mathbb{E}_{W^{(t-1)}}\left[\mathbb{E}_{W^{(t)}}\left[\|BA_{1,t}\|_{\mathrm{F}}^{2}\right]\right]\right]\right] \leq \rho^{t}\|B\|_{\mathrm{F}}^{2}.$$

**Proof of Lemma 3.** By the definition of  $\kappa$ ,

$$\kappa \eta \binom{s}{2} L_i \ge \frac{(1 + \eta L_i)^s - 1 - s \eta_t L_i}{\eta L_i}.$$

Hence it suffices to show

$$\left\| \sum_{k=0}^{s-1} \left[ \nabla \ell_i(\boldsymbol{x}^{(t,k)}) - \nabla \ell_i(\boldsymbol{x}^t) \right] \right\|_2 \le \frac{(1 + \eta L_i)^s - 1 - s\eta L_i}{\eta L_i} \left\| \nabla \ell_i(\boldsymbol{x}_t) \right\|_2.$$
 (6)

We prove (6) holds for all  $s \ge 1$  by induction. The base case s = 1 follows from the definition. Suppose (6) holds true for  $s = 1, \ldots, n-1$ , where  $n \ge 2$ . Next we prove (6) for s = n. We have

$$\left\|\nabla \ell_{i}(\boldsymbol{x}^{(t,n-1)}) - \nabla \ell_{i}(\boldsymbol{x}^{t})\right\|_{2} \leq L_{i} \left\|\boldsymbol{x}^{(t,n-1)} - \boldsymbol{x}^{t}\right\|_{2}$$

$$\leq L_{i} \eta \left\|\sum_{k=0}^{n-2} \left[\nabla \ell_{i}\left(\boldsymbol{x}^{(t,k)}\right) - \nabla \ell_{i}\left(\boldsymbol{x}^{t}\right)\right]\right\|_{2} + L_{i} \eta \left(n-1\right) \left\|\nabla \ell_{i}\left(\boldsymbol{x}^{t}\right)\right\|_{2}$$

$$\stackrel{\text{(a)}}{\leq} \left[\left(1 + \eta L_{i}\right)^{n-1} - 1\right] \left\|\nabla \ell_{i}\left(\boldsymbol{x}^{t}\right)\right\|_{2},$$

$$(7)$$

where (a) follows from the induction hypothesis.

Plug Eq. (7) back in, use the induction hypothesis and triangle inequality, we get

$$\begin{split} \left\| \sum_{k=0}^{n-1} \left[ \nabla \ell_{i}(\boldsymbol{x}^{(t,k)}) - \nabla \ell_{i}(\boldsymbol{x}^{t}) \right] \right\|_{2} &\leq \left\| \sum_{k=0}^{n-2} \left[ \nabla \ell_{i}(\boldsymbol{x}^{(t,k)}) - \nabla \ell_{i}(\boldsymbol{x}^{t}) \right] \right\|_{2} + \left\| \nabla \ell_{i}(\boldsymbol{x}^{(t,n-1)}) - \nabla \ell_{i}(\boldsymbol{x}^{t}) \right\|_{2} \\ &\leq \left[ \frac{\left(1 + \eta L_{i}\right)^{n-1} - 1 - \left(n - 1\right) \eta L_{i}}{\eta L_{i}} + \left(1 + \eta L_{i}\right)^{n-1} - 1 \right] \left\| \nabla \ell_{i}\left(\boldsymbol{x}^{t}\right) \right\|_{2} \\ &= \frac{\left(1 + \eta L_{i}\right)^{n} - 1 - n \eta L_{i}}{\eta L_{i}} \left\| \nabla \ell_{i}\left(\boldsymbol{x}^{t}\right) \right\|_{2}. \end{split}$$

The proof is completed.

#### Proof of Claim 1. Recall that

$$\kappa = \frac{(1 + \eta L)^s - 1 - s\eta L}{\binom{s}{2} (\eta L)^2}.$$

From binomial theorem, we know that

$$(1 + \eta L)^s = \sum_{i=0}^s \binom{s}{i} (\eta L)^i,$$

it follows that

$$\frac{\left(1+\eta L\right)^{s}-1-s\eta L}{\binom{s}{2}\left(\eta L\right)^{2}}=\frac{\sum_{i=2}^{s}\binom{s}{i}\left(\eta L\right)^{i}}{\binom{s}{2}\left(\eta L\right)^{2}}=\sum_{i=2}^{s}\frac{\binom{s}{i}}{\binom{s}{2}}\left(\eta L\right)^{i-2}.$$

Since  $i-2 \geq 0$  for  $i \geq 2$ , we can see that  $\kappa$  is a polynomial of  $(\eta L)$ . Thus, it is monotonic non-decreasing w.r.t.  $\eta > 0$ . The proof is completed.

**Proposition 2.** For any  $t \in [T-1]$ , it holds that

$$\frac{1}{m} \sum_{i=1}^{m} \left\| \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 \le \frac{3L^2}{m} \sum_{i=1}^{m} \left\| \boldsymbol{x}_i^t - \bar{\boldsymbol{x}}^t \right\|_2^2 + 3\left(\beta^2 + 1\right) \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 + 3\xi^2.$$

# Proof of Proposition 2.

$$\begin{split} \frac{1}{m} \sum_{i=1}^{m} \left\| \nabla F_{i}(\boldsymbol{x}_{i}^{t}) \right\|_{2}^{2} &= \frac{1}{m} \sum_{i=1}^{m} \left\| \nabla F_{i}(\boldsymbol{x}_{i}^{t}) - \nabla F_{i}(\bar{\boldsymbol{x}}^{t}) + \nabla F_{i}(\bar{\boldsymbol{x}}^{t}) - \nabla F(\bar{\boldsymbol{x}}^{t}) + \nabla F(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} \\ &\leq \frac{3}{m} \sum_{i=1}^{m} \left\| \nabla F_{i}(\boldsymbol{x}_{i}^{t}) - \nabla F_{i}(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} + \frac{3}{m} \sum_{i=1}^{m} \left\| \nabla F_{i}(\bar{\boldsymbol{x}}^{t}) - \nabla F(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} + 3 \left\| \nabla F(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} \\ &\leq \frac{3L^{2}}{m} \sum_{i=1}^{m} \left\| \boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t} \right\|_{2}^{2} + 3\beta^{2} \left\| \nabla F(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} + 3\xi^{2} + 3 \left\| \nabla F(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} \\ &= \frac{3L^{2}}{m} \sum_{i=1}^{m} \left\| \boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t} \right\|_{2}^{2} + 3\left(\beta^{2} + 1\right) \left\| \nabla F(\bar{\boldsymbol{x}}^{t}) \right\|_{2}^{2} + 3\xi^{2}, \end{split}$$

where inequality (a) follows from Assumptions 1 and 4.

**Proof of Lemma 4.** By L-smoothness, we have

$$F(\bar{\boldsymbol{x}}^{t+1}) - F(\bar{\boldsymbol{x}}^t) \le \left\langle \nabla F(\bar{\boldsymbol{x}}^t), \bar{\boldsymbol{x}}^{t+1} - \bar{\boldsymbol{x}}^t \right\rangle + \frac{L}{2} \left\| \bar{\boldsymbol{x}}^{t+1} - \bar{\boldsymbol{x}}^t \right\|_2^2$$
$$= \left\langle \nabla F(\bar{\boldsymbol{x}}^t), -\frac{\eta}{m} \boldsymbol{G}^{(t)} \mathbf{1} \right\rangle + \frac{L\eta^2}{2} \left\| \frac{1}{m} \boldsymbol{G}^{(t)} \mathbf{1} \right\|_2^2.$$

Taking expectations with respect to the randomness in the mini-batches at k-th rounds, we have

$$\mathbb{E}\left[F(\bar{\boldsymbol{x}}^{t+1}) - F(\bar{\boldsymbol{x}}^t) \mid \mathcal{F}^t\right] \leq \mathbb{E}\left[\left\langle \nabla F(\bar{\boldsymbol{x}}^t), -\frac{\eta}{m}\boldsymbol{G}^{(t)}\boldsymbol{1}\right\rangle + \frac{L}{2}\left\|-\frac{\eta}{m}\boldsymbol{G}^{(t)}\boldsymbol{1}\right\|_2^2\right].$$

For ease of notations, we abbreviate  $abla \ell_i \left( oldsymbol{x}_i^{(t,k)} 
ight)$  as  $abla \ell_i^{(t,k)}$ .

a) Bounding  $\langle \nabla f(\bar{\boldsymbol{x}}^t), -\frac{\eta}{m} \nabla \boldsymbol{F}^{(t)} \boldsymbol{1} \rangle$ .:

$$\begin{split} & \mathbb{E}\left[\left\langle \nabla F(\bar{\boldsymbol{x}}^t), -\frac{\eta}{m}\boldsymbol{G}^{(t)}\boldsymbol{1}\right\rangle \mid \mathcal{F}^t\right] = -\frac{\eta}{m}\mathbb{E}\left[\left\langle \nabla F(\bar{\boldsymbol{x}}^t), \sum_{i=1}^m \sum_{k=0}^{s-1} \nabla \ell_i^{(t,k)} \right\rangle \mid \mathcal{F}^t\right] \\ & = -\frac{\eta}{m}\mathbb{E}\left[\left\langle \nabla F(\bar{\boldsymbol{x}}^t), \sum_{i=1}^m s \nabla \ell_i^{(t,0)} - s \nabla \ell_i^{(t,0)} + \sum_{k=0}^{s-1} \nabla \ell_i^{(t,k)} \right\rangle \mid \mathcal{F}^t\right] \\ & = -\frac{s\eta}{m}\left\langle \nabla F(\bar{\boldsymbol{x}}^t), \sum_{i=1}^m \nabla F_i(\boldsymbol{x}_i^t) \right\rangle + \mathbb{E}\left[\sum_{i=1}^m \frac{\eta}{m} \left\langle \nabla F(\bar{\boldsymbol{x}}^t), s \nabla \ell_i^{(t,0)} - \sum_{k=0}^{s-1} \nabla \ell_i^{(t,k)} \right\rangle \mid \mathcal{F}^t\right] \\ & = \underbrace{-\frac{s\eta}{m}\left\langle \nabla F(\bar{\boldsymbol{x}}^t), \nabla \boldsymbol{F}^{(t)}\boldsymbol{1} \right\rangle}_{(\mathbf{A})} + \underbrace{\mathbb{E}\left[\frac{\eta}{m} \left\langle \nabla F(\bar{\boldsymbol{x}}^t), \sum_{i=1}^m s \nabla \ell_i^{(t,0)} - \sum_{k=0}^{s-1} \nabla \ell_i^{(t,k)} \right\rangle \mid \mathcal{F}^t\right]}_{(\mathbf{B})}. \end{split}$$

Term (A) can be bounded as

$$\begin{split} \left\langle \nabla F(\bar{\boldsymbol{x}}^t), -\frac{s\eta}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\rangle &= -s\eta \left\langle \nabla F(\bar{\boldsymbol{x}}^t), \frac{1}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\rangle \\ &= -\frac{s\eta}{2} \left( \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 + \left\| \frac{1}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\|_2^2 - \left\| \nabla F(\bar{\boldsymbol{x}}^t) - \frac{1}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\|_2^2 \right) \\ &= -\frac{s\eta}{2} \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 - \frac{s\eta}{2} \left\| \frac{1}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\|_2^2 + \frac{s\eta}{2} \left\| \nabla F(\bar{\boldsymbol{x}}^t) - \frac{1}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\|_2^2 \\ &\leq -\frac{s\eta}{2} \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 - \frac{s\eta}{2} \left\| \frac{1}{m} \nabla \boldsymbol{F}^{(t)} \mathbf{1} \right\|_2^2 + \frac{s\eta L^2}{2m} \sum_{i=1}^m \left\| \bar{\boldsymbol{x}}^t - \boldsymbol{x}_i^t \right\|_2^2. \end{split}$$

For term (B), we have

$$\begin{split} & \mathbb{E}\left[\frac{\eta}{m}\left\langle\nabla F(\bar{\boldsymbol{x}}^t), \sum_{i=1}^m s\nabla\ell_i^{(t,0)} - \sum_{k=0}^{s-1}\nabla\ell_i^{(t,k)}\right\rangle \mid \mathcal{F}^t\right] \\ & = \frac{\eta}{m}\sum_{i=1}^m \left\langle\nabla F(\bar{\boldsymbol{x}}^t), \mathbb{E}\left[s\nabla\ell_i^{(t,0)} - \sum_{k=0}^{s-1}\nabla\ell_i^{(t,k)} \mid \mathcal{F}^t\right]\right\rangle \\ & \leq \frac{\eta}{2m}\sum_{i=1}^m \left(\eta s^2 \left\|\nabla F(\bar{\boldsymbol{x}}^t)\right\|_2^2 + \frac{1}{\eta s^2} \left\|\mathbb{E}\left[s\nabla\ell_i^{(t,0)} - \sum_{k=0}^{s-1}\nabla\ell_i^{(t,k)} \mid \mathcal{F}^t\right]\right\|_2^2\right) \\ & \stackrel{\text{(a)}}{\leq} \frac{\eta^2 s^2}{2} \left\|\nabla F(\bar{\boldsymbol{x}}^t)\right\|_2^2 + \underbrace{\frac{1}{2ms^2}\sum_{i=1}^m \mathbb{E}\left[\left\|s\nabla\ell_i^{(t,0)} - \sum_{k=0}^{s-1}\nabla\ell_i^{(t,k)}\right\|_2^2 \mid \mathcal{F}^t\right]}_{\text{(B.1)}}. \end{split}$$

From Lemma 3, we bound term (B.1) as follows

$$\begin{split} &\frac{1}{2ms^2} \sum_{i=1}^m \mathbb{E} \left[ \left\| s \nabla \ell_i^{(t,0)} - \sum_{k=0}^{s-1} \nabla \ell_i^{(t,k)} \right\|_2^2 \mid \mathcal{F}^t \right] \\ &\leq \frac{1}{2ms^2} \sum_{i=1}^m \mathbb{E} \left[ \kappa^2 \eta^2 \binom{s}{2}^2 L^2 \left\| \nabla \ell_i^{(t,0)} \right\|_2^2 \mid \mathcal{F}^t \right] \\ &= \frac{\kappa^2 \eta^2 \binom{s}{2}^2 L^2}{2ms^2} \sum_{i=1}^m \mathbb{E} \left[ \left\| \nabla \ell_i^{(t,0)} \right\|_2^2 \mid \mathcal{F}^t \right] \\ &= \frac{\kappa^2 \eta^2 \binom{s}{2}^2 L^2}{2ms^2} \sum_{i=1}^m \mathbb{E} \left[ \left\| \nabla \ell_i^{(t,0)} - \nabla F_i(\boldsymbol{x}_i^t) + \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 \mid \mathcal{F}^t \right] \\ &\leq \frac{\kappa^2 \eta^2 \binom{s}{2}^2 L^2}{ms^2} \sum_{i=1}^m \mathbb{E} \left[ \left\| \nabla \ell_i^{(t,0)} - \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 \mid \mathcal{F}^t \right] + \left\| \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 \\ &\leq \kappa^2 \eta^2 s^2 L^2 \sigma^2 + \frac{\kappa^2 \eta^2 s^2 L^2}{m} \sum_{i=1}^m \left\| \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 \\ &\leq \kappa^2 \eta^2 s^2 L^2 \frac{3L^2}{m} \sum_{i=1}^m \left\| \boldsymbol{x}_i^t - \bar{\boldsymbol{x}}^t \right\|_2^2 + 3\kappa^2 \eta^2 s^2 L^2 \left( \beta^2 + 1 \right) \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 + \kappa^2 \eta^2 s^2 L^2 (3\xi^2 + \sigma^2), \end{split}$$

where inequality (a) follows from Assumption 2, and the last inequality follows from Proposition 2. Thus, term (B) can be further bounded as

$$\begin{split} & \mathbb{E}\left[\frac{\eta}{m}\left\langle\nabla F(\bar{\boldsymbol{x}}^{t}), \sum_{i=1}^{m} s \nabla \ell_{i}^{(t,0)} - \sum_{k=0}^{s-1} \nabla \ell_{i}^{(t,k)}\right\rangle \mid \mathcal{F}^{t}\right] \\ & \leq \frac{\eta^{2} s^{2}}{2} \left\|\nabla F(\bar{\boldsymbol{x}}^{t})\right\|_{2}^{2} + \frac{3L^{4} \eta^{2} \kappa^{2} s^{2}}{m} \sum_{i=1}^{m} \left\|\boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t}\right\|_{2}^{2} + 3\kappa^{2} \eta^{2} s^{2} L^{2} \left(\beta^{2} + 1\right) \left\|\nabla F(\bar{\boldsymbol{x}}^{t})\right\|_{2}^{2} + \kappa^{2} \eta^{2} s^{2} L^{2} (3\xi^{2} + \sigma^{2}). \end{split}$$

Combing the bounds of terms (A) and (B), we get

$$\mathbb{E}\left[\left\langle \nabla F(\bar{\boldsymbol{x}}^{t}), -\frac{\eta}{m}\boldsymbol{G}^{(t)}\boldsymbol{1}\right\rangle \mid \mathcal{F}^{t}\right] \leq -\left[\frac{s\eta}{2} - \frac{\eta^{2}s^{2}}{2} - 3\kappa^{2}\eta^{2}s^{2}L^{2}\left(\beta^{2} + 1\right)\right] \left\|\nabla F(\bar{\boldsymbol{x}}^{t})\right\|_{2}^{2} \\
-\frac{s\eta}{2} \left\|\frac{1}{m}\nabla \boldsymbol{F}^{(t)}\boldsymbol{1}\right\|_{2}^{2} + \kappa^{2}\eta^{2}s^{2}L^{2}(3\xi^{2} + \sigma^{2}) \\
+\left(\frac{s\eta L^{2}}{2m} + \kappa^{2}\eta^{2}s^{2}L^{2}\frac{3L^{2}}{m}\right) \sum_{i=1}^{m} \left\|\bar{\boldsymbol{x}}^{t} - \boldsymbol{x}_{i}^{t}\right\|_{2}^{2}.$$
(8)

b) Bounding  $\mathbb{E}\left[\left\|\frac{1}{m}\pmb{G}^{(t)}\pmb{1}\right\|_2^2\mid\mathcal{F}^t\right]$ . So, we have

$$\begin{split} \left\| \frac{1}{m} \boldsymbol{G}^{(t)} \mathbf{1} \right\|_{2}^{2} &= \left\| \frac{1}{m} \sum_{i=1}^{m} \sum_{k=0}^{s-1} \nabla \ell_{i}^{(t,k)} \right\|_{2}^{2} \\ &= \left\| \frac{1}{m} \sum_{i=1}^{m} \sum_{k=0}^{s-1} \left( \nabla \ell_{i}^{(t,k)} - \nabla \ell_{i}^{(t,0)} + \nabla \ell_{i}^{(t,0)} \right) \right\|_{2}^{2} \\ &\leq 2 \underbrace{\left\| \frac{1}{m} \sum_{i=1}^{m} \sum_{k=0}^{s-1} \left( \nabla \ell_{i}^{(t,k)} - \nabla \ell_{i}^{(t,0)} \right) \right\|_{2}^{2}}_{(\mathbf{C})} + 2 \underbrace{\left\| \frac{s}{m} \sum_{i=1}^{m} \nabla \ell_{i}^{(t,0)} \right\|_{2}^{2}}_{(\mathbf{D})}. \end{split}$$

For term (C), by Lemma 3, we have

$$\begin{split} \left\| \frac{1}{m} \sum_{i=1}^{m} \sum_{k=0}^{s-1} \left( \nabla \ell_i^{(t,k)} - \nabla \ell_i^{(t,0)} \right) \right\|_2 &\leq \frac{1}{m} \sum_{i=1}^{m} \left\| \sum_{k=0}^{s-1} \left( \nabla \ell_i^{(t,k)} - \nabla \ell_i^{(t,0)} \right) \right\|_2 \\ &\leq \frac{\kappa \eta s^2 L}{2m} \sum_{i=1}^{m} \left\| \nabla \ell_i^{(t,0)} \right\|_2. \end{split}$$

Thus, we get

$$\left\| \frac{1}{m} \sum_{i=1}^{m} \sum_{k=0}^{s-1} \left( \nabla \ell_i^{(t,k)} - \nabla \ell_i^{(t,0)} \right) \right\|_2^2 \leq \frac{\kappa^2 \eta^2 s^4 L^2}{4m} \sum_{i=1}^{m} \left\| \nabla \ell_i^{(t,0)} \right\|_2^2 \\
\leq \frac{\kappa^2 \eta^2 s^4 L^2}{2m} \left( \sum_{i=1}^{m} \left\| \nabla \ell_i^{(t,0)} - \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 + \sum_{i=1}^{m} \left\| \nabla F_i(\boldsymbol{x}_i^t) \right\|_2^2 \right).$$

By Assumption 2, we obtain

$$\mathbb{E}\left[\left\|\frac{1}{m}\sum_{i=1}^{m}\sum_{k=0}^{s-1}\left(\nabla \ell_{i}^{(t,k)} - \nabla \ell_{i}^{(t,0)}\right)\right\|_{2}^{2} \mid \mathcal{F}^{t}\right] \leq \frac{\kappa^{2}\eta^{2}s^{4}L^{2}\sigma^{2}}{2} + \frac{\kappa^{2}\eta^{2}s^{4}L^{2}}{2m}\sum_{i=1}^{m}\left\|\nabla F_{i}(\boldsymbol{x}_{i}^{t})\right\|_{2}^{2}.$$

For term (D), by Assumption 2, we have

$$\mathbb{E}\left[\frac{s^2}{m^2}\left\|\sum_{i=1}^{m}\nabla \ell_i^{(t,0)}\right\|_2^2\mid \mathcal{F}^t\right] \leq \mathbb{E}\left[\frac{2s^2}{m^2}\left\|\sum_{i=1}^{m}\nabla \ell_i^{(t,0)} - \nabla F_i(\boldsymbol{x}_i^t)\right\|_2^2\mid \mathcal{F}^t\right] + \frac{2s^2}{m^2}\left\|\sum_{i=1}^{m}\nabla F_i(\boldsymbol{x}_i^t)\right\|_2^2$$

$$\leq \frac{2s^2\sigma^2}{m} + \frac{2s^2}{m^2}\left\|\sum_{i=1}^{m}\nabla F_i(\boldsymbol{x}_i^t)\right\|_2^2.$$

Combing the above upper bounds of terms (C) and (D), we get

$$\mathbb{E}\left[\left\|\frac{1}{m}\boldsymbol{G}^{(t)}\mathbf{1}\right\|_{2}^{2} \mid \mathcal{F}^{t}\right] \leq 2\left[\frac{2s^{2}}{m}\sum_{i=1}^{m}\left\|\nabla F_{i}(\boldsymbol{x}_{i}^{t})\right\|_{2}^{2} + \frac{\kappa^{2}\eta^{2}s^{4}L^{2}}{2m}\sum_{i=1}^{m}\left\|\nabla F_{i}(\boldsymbol{x}_{i}^{t})\right\|_{2}^{2} + s^{2}\sigma^{2}\left(\frac{2}{m} + \frac{\kappa^{2}\eta^{2}s^{2}L^{2}}{2}\right)\right]$$

$$= s^{2}\left(4 + \kappa^{2}\eta^{2}s^{2}L^{2}\right)\frac{1}{m}\sum_{i=1}^{m}\left\|\nabla F_{i}(\boldsymbol{x}_{i}^{t})\right\|_{2}^{2} + s^{2}\sigma^{2}\left(\frac{4}{m} + \kappa^{2}\eta^{2}s^{2}L^{2}\right).$$

Applying Proposition 2, we get

$$\mathbb{E}\left[\left\|\frac{1}{m}\boldsymbol{G}^{(t)}\mathbf{1}\right\|_{2}^{2} \mid \mathcal{F}^{t}\right] \leq 6s^{2}L^{2}\left(2 + \frac{\kappa^{2}\eta^{2}s^{2}L^{2}}{2}\right) \frac{1}{m} \sum_{i=1}^{m} \left\|\boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t}\right\|_{2}^{2} + 6s^{2}\left(\beta^{2} + 1\right)\left(2 + \frac{\kappa^{2}\eta^{2}s^{2}L^{2}}{2}\right) \left\|\nabla F(\bar{\boldsymbol{x}}^{t})\right\|_{2}^{2} + 6s^{2}\xi^{2}\left(2 + \frac{\kappa^{2}\eta^{2}s^{2}L^{2}}{2}\right) + 2s^{2}\sigma^{2}\left(\frac{2}{m} + \frac{\kappa^{2}\eta^{2}s^{2}L^{2}}{2}\right). \tag{9}$$

c) Putting them together.: With Eq.(8) and (9), we have

$$\mathbb{E}\left[F(\bar{x}^{t+1}) - F(\bar{x}^{t}) \mid \mathcal{F}^{t}\right] \leq \mathbb{E}\left[\left\langle\nabla F(\bar{x}^{t}), -\frac{\eta}{m}G^{(t)}\mathbf{1}\right\rangle \mid \mathcal{F}^{t}\right] + \frac{L\eta^{2}}{2}\mathbb{E}\left[\left\|\frac{1}{m}G^{(t)}\mathbf{1}\right\|_{2}^{2} \mid \mathcal{F}^{t}\right] \\
\leq -\left[\frac{\eta s}{2} - \frac{\eta^{2}s^{2}}{2} - 3\kappa^{2}\eta^{2}s^{2}L^{2}\left(\beta^{2} + 1\right)\right] \left\|\nabla F(\bar{x}^{t})\right\|_{2}^{2} \\
- \frac{\eta s}{2}\left\|\frac{1}{m}\nabla F^{(t)}\mathbf{1}\right\|_{2}^{2} + \kappa^{2}\eta^{2}s^{2}L^{2}(3\xi^{2} + \sigma^{2}) \\
+ \left(\frac{s\eta L^{2}}{2m} + \kappa^{2}\eta^{2}s^{2}\frac{3L^{4}}{m}\right)\sum_{i=1}^{m}\left\|\bar{x}^{t} - x_{i}^{t}\right\|_{2}^{2} \\
+ \frac{L\eta^{2}}{2}6s^{2}L^{2}\left(2 + \frac{\kappa^{2}L^{2}}{2}\right)\frac{1}{m}\sum_{i=1}^{m}\left\|x_{i}^{t} - \bar{x}^{t}\right\|_{2}^{2} \\
+ \frac{L\eta^{2}}{2}6s^{2}\left(\beta^{2} + 1\right)\left(2 + \frac{\kappa^{2}L^{2}}{2}\right)\left\|\nabla F(\bar{x}^{t})\right\|_{2}^{2} \\
+ \frac{L\eta^{2}}{2}6s^{2}\xi^{2}\left(2 + \frac{\kappa^{2}L^{2}}{2}\right) + \frac{L\eta^{2}}{2}2s^{2}\sigma^{2}\left(\frac{2}{m} + \frac{\kappa^{2}L^{2}}{2}\right).$$

We can choose  $\eta \leq \frac{1}{2s}$  so that

$$\mathbb{E}\left[F(\bar{\boldsymbol{x}}^{t+1}) - F(\bar{\boldsymbol{x}}^{t}) \mid \mathcal{F}^{t}\right] \leq \mathbb{E}\left[\left\langle\nabla F(\bar{\boldsymbol{x}}^{t}), -\frac{\eta}{m}\boldsymbol{G}^{(t)}\boldsymbol{1}\right\rangle \mid \mathcal{F}^{t}\right] + \frac{L\eta^{2}}{2}\mathbb{E}\left[\left\|\frac{1}{m}\boldsymbol{G}^{(t)}\boldsymbol{1}\right\|_{2}^{2} \mid \mathcal{F}^{t}\right]$$

$$\leq -\left\{\frac{s\eta}{4} - 3\eta^{2}s^{2}\left(\beta^{2} + 1\right)\left[\kappa^{2}L^{2} + 2L\left(1 + \frac{\kappa^{2}L^{2}}{4}\right)\right]\right\}\left\|\nabla F(\bar{\boldsymbol{x}}^{t})\right\|_{2}^{2}$$

$$+ 3\xi^{2}\eta^{2}s^{2}\left[\kappa^{2}L^{2} + 2L\left(1 + \frac{\kappa^{2}L^{2}}{4}\right)\right]$$

$$+ \sigma^{2}\eta^{2}s^{2}\left[\kappa^{2}L^{2} + 2L\left(\frac{1}{m} + \frac{\kappa^{2}L^{2}}{4}\right)\right]$$

$$+\left\{\eta sL^{2} + 3\eta^{2}s^{2}L^{2}\left[\kappa^{2}L^{2} + 2L\left(1 + \frac{\kappa^{2}L^{2}}{4}\right)\right]\right\}\frac{1}{m}\sum_{i=1}^{m}\left\|\boldsymbol{x}_{i}^{t} - \bar{\boldsymbol{x}}^{t}\right\|_{2}^{2}.$$

**Proof of Lemma 5.** Our proof shares the same outline as that in [12] yet with non-trivial adaptation to account for multiple local updates and the fact the stochastic gradients at a client within each round are *not independent*. Particularly,  $T_1$  in Eq. (10) does not exist in [12].

We have the following relations:

$$\boldsymbol{X}^{(t)}\left(\mathbf{I} - \mathbf{J}\right) = \left(\boldsymbol{X}^{(t-1)} - \eta \boldsymbol{G}^{(t-1)}\right) W^{(t-1)} \left(\mathbf{I} - \mathbf{J}\right)$$
$$= -\eta \sum_{q=0}^{t-1} \boldsymbol{G}^{(q)} \left( \prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J} \right),$$

where the last follows from the fact that all clients are initiated at the same weights. It follows that

$$\|\mathbf{X}^{(t)}\left(\mathbf{I} - \mathbf{J}\right)\|_{F}^{2} \leq 3\eta^{2} \|\underbrace{\sum_{q=0}^{t-1} \left(\mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)}\right) \left(\prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J}\right)\|_{F}^{2}}_{\mathbf{T}_{1}} + 3\eta^{2} \|\underbrace{\sum_{q=0}^{t-1} \left(\mathbf{G}_{0}^{(q)} - s\nabla\mathbf{F}^{(q)}\right) \left(\prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J}\right)\|_{F}^{2}}_{\mathbf{T}_{2}} + 3\eta^{2}s^{2} \|\underbrace{\sum_{q=0}^{t-1} \nabla\mathbf{F}^{(q)} \left(\prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J}\right)\|_{F}^{2}}_{\mathbf{T}_{3}}.$$
(10)

d) Bounding  $\mathbb{E}[T_1]$  .:

$$\mathbb{E}\left[\mathbf{T}_{1}\right] = \sum_{q=0}^{t-1} \mathbb{E}\left[ \left\| \left( \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)} \right) \left( \prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J} \right) \right\|_{F}^{2} \right] \\
+ \sum_{q=0}^{t-1} \sum_{p=0, p \neq q}^{t-1} \mathbb{E}\left[ \left\langle \left( \mathbf{G}^{(p)} - \mathbf{G}_{0}^{(p)} \right) \left( \prod_{\ell=p}^{t-1} W^{(\ell)} - \mathbf{J} \right), \left( \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)} \right) \left( \prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J} \right) \right\rangle \right] \\
\stackrel{(a)}{\leq} \sum_{q=0}^{t-1} \rho^{t-q} \mathbb{E}\left[ \left\| \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)} \right\|_{F}^{2} \right] \\
+ \sum_{q=0}^{t-1} \sum_{p=0, p \neq q}^{t-1} \mathbb{E}\left[ \left\| \left( \mathbf{G}^{(p)} - \mathbf{G}_{0}^{(p)} \right) \left( \prod_{\ell=p}^{t-1} W^{(\ell)} - \mathbf{J} \right) \right\|_{F} \right\| \left( \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)} \right) \left( \prod_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J} \right) \right\|_{F} \right] \\
\leq \sum_{q=0}^{t-1} \rho^{t-q} \mathbb{E}\left[ \left\| \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)} \right\|_{F}^{2} \right] \\
+ \sum_{q=0}^{t-1} \sum_{p=0, p \neq q}^{t-1} \mathbb{E}\left[ \left\| \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(p)} \right\|_{F}^{2} \right] \\
+ \sum_{q=0}^{t-1} \sum_{p=0, p \neq q}^{t-1} \mathbb{E}\left[ \left\| \mathbf{G}^{(p)} - \mathbf{G}_{0}^{(p)} \right\|_{F}^{2} \right] + \frac{\epsilon \rho^{t-q}}{2} \left\| \left( \mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)} \right) \right\|_{F}^{2} \right], \tag{11}$$

where inequality (a) follows from Lemma 2, and Cauchy-Schwarz inequality. Next, we bound the second term, choose  $\epsilon = \rho^{\frac{q-p}{2}}$ ,

$$\begin{split} &\sum_{q=0}^{t-1} \sum_{p=0, p \neq q}^{t-1} \frac{\sqrt{\rho}^{2t-p-q}}{2} \mathbb{E} \left[ \| \left( \boldsymbol{G}^{(p)} - \boldsymbol{G}_{0}^{(p)} \right) \|_{\mathrm{F}}^{2} + \| \left( \boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)} \right) \|_{\mathrm{F}}^{2} \right] \\ &\leq \sum_{q=0}^{t-1} \sum_{p=0}^{t-1} \frac{\sqrt{\rho}^{2t-p-q}}{2} \mathbb{E} \left[ \| \left( \boldsymbol{G}^{(p)} - \boldsymbol{G}_{0}^{(p)} \right) \|_{\mathrm{F}}^{2} + \| \left( \boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)} \right) \|_{\mathrm{F}}^{2} \right] \\ &= \sum_{p=0}^{t-1} \frac{\sqrt{\rho}^{t-p}}{2} \mathbb{E} \left[ \| \left( \boldsymbol{G}^{(p)} - \boldsymbol{G}_{0}^{(p)} \right) \|_{\mathrm{F}}^{2} \right] \sum_{q=0}^{t-1} \sqrt{\rho}^{t-q} + \sum_{q=0}^{t-1} \frac{\sqrt{\rho}^{t-p}}{2} \mathbb{E} \left[ \| \left( \boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)} \right) \|_{\mathrm{F}}^{2} \right] \sum_{p=0}^{t-1} \sqrt{\rho}^{t-p} \\ &= \frac{\sqrt{\rho} - \sqrt{\rho}^{t+1}}{1 - \sqrt{\rho}} \sum_{q=0}^{t-1} \sqrt{\rho}^{t-q} \mathbb{E} \left[ \| \left( \boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)} \right) \|_{\mathrm{F}}^{2} \right]. \end{split}$$

Plugging the above bound back in Eq.(11), we get

$$\begin{split} \mathbb{E}\left[\mathbf{T}_{1}\right] &\leq \sum_{q=0}^{t-1} \left[\sqrt{\rho}^{t-q} + \frac{\sqrt{\rho} - \sqrt{\rho}^{t+1}}{1 - \sqrt{\rho}}\right] \sqrt{\rho}^{t-q} \mathbb{E}\left[\|\boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)}\|_{\mathrm{F}}^{2}\right] \\ &= \sum_{q=0}^{t-1} \left[\frac{\sqrt{\rho} + \sqrt{\rho}^{t+1} \left(\frac{1 - \sqrt{\rho}}{\sqrt{\rho}^{q+1}} - 1\right)}{1 - \sqrt{\rho}}\right] \sqrt{\rho}^{t-q} \mathbb{E}\left[\|\boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)}\|_{\mathrm{F}}^{2}\right] \\ &\leq \sum_{q=0}^{t-1} \left[\frac{\sqrt{\rho} + \sqrt{\rho} \left(1 - \sqrt{\rho} - \sqrt{\rho}^{t}\right)}{1 - \sqrt{\rho}}\right] \sqrt{\rho}^{t-q} \mathbb{E}\left[\|\boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)}\|_{\mathrm{F}}^{2}\right] \\ &\leq \frac{2\sqrt{\rho}}{1 - \sqrt{\rho}} \sum_{q=0}^{t-1} \sqrt{\rho}^{t-q} \mathbb{E}\left[\|\boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)}\|_{\mathrm{F}}^{2}\right]. \end{split}$$

It remains to bound  $\mathbb{E}\left[\|m{G}^{(q)}-m{G}_0^{(q)}\|_{\mathrm{F}}^2
ight]$ 

$$\mathbb{E}\left[\|\boldsymbol{G}^{(q)} - \boldsymbol{G}_{0}^{(q)}\|_{\mathrm{F}}^{2}\right] \stackrel{(a)}{\leq} \kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} \mathbb{E}\left[\|\boldsymbol{G}_{0}^{(q)} - s \nabla \boldsymbol{F}^{(q)} + s \nabla \boldsymbol{F}^{(q)}\|_{\mathrm{F}}^{2}\right] \\
\leq 2\kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} \mathbb{E}\left[\|\boldsymbol{G}_{0}^{(q)} - s \nabla \boldsymbol{F}^{(q)}\|_{\mathrm{F}}^{2}\right] + 2\kappa^{2} s^{2} \eta^{2} \binom{s}{2}^{2} L^{2} \mathbb{E}\left[\|\nabla \boldsymbol{F}^{(q)}\|_{\mathrm{F}}^{2}\right] \\
\leq 2\kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} m \sigma^{2} + 2\kappa^{2} s^{2} \eta^{2} \binom{s}{2}^{2} L^{2} \mathbb{E}\left[\|\nabla \boldsymbol{F}^{(q)}\|_{\mathrm{F}}^{2}\right],$$

where inequality (a) follows from Lemma 3. Thus,

$$\mathbb{E}\left[\mathbf{T}_{1}\right] \leq \frac{2\sqrt{\rho}}{1-\sqrt{\rho}} \sum_{q=0}^{t-1} \sqrt{\rho^{t-q}} \mathbb{E}\left[\|\mathbf{G}^{(q)} - \mathbf{G}_{0}^{(q)}\|_{\mathrm{F}}^{2}\right] \\
\leq \frac{2\sqrt{\rho}}{1-\sqrt{\rho}} \sum_{q=0}^{t-1} \sqrt{\rho^{t-q}} \left[2\kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} m \sigma^{2} + 2\kappa^{2} s^{2} \eta^{2} \binom{s}{2}^{2} L^{2} \mathbb{E}\left[\|\nabla \mathbf{F}^{(q)}\|_{\mathrm{F}}^{2}\right]\right] \\
\leq \frac{4\kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} m \sigma^{2} \rho}{\left(1-\sqrt{\rho}\right)^{2}} + \frac{4\kappa^{2} s^{2} \eta^{2} \binom{s}{2}^{2} L^{2} \sqrt{\rho}}{1-\sqrt{\rho}} \sum_{q=0}^{t-1} \sqrt{\rho^{t-q}} \mathbb{E}\left[\|\nabla \mathbf{F}^{(q)}\|_{\mathrm{F}}^{2}\right]$$

e) Bounding  $\mathbb{E}[T_2]$  .:

$$\mathbb{E}\left[\mathbf{T}_{2}\right] = \mathbb{E}\left[\left\|\sum_{q=0}^{t-1}\left(\boldsymbol{G}_{0}^{(q)} - s\nabla\boldsymbol{F}^{(q)}\right)\left(\boldsymbol{\Pi}_{\ell=q}^{t-1}W^{(\ell)} - \mathbf{J}\right)\right\|_{\mathrm{F}}^{2}\right]$$

$$\leq \sum_{q=0}^{t-1}\rho^{t-q}\mathbb{E}\left[\left\|\left(\boldsymbol{G}_{0}^{(q)} - s\nabla\boldsymbol{F}^{(q)}\right)\right\|_{\mathrm{F}}^{2}\right]$$

$$\leq \frac{\rho ms^{2}\sigma^{2}}{1-\rho}.$$

f) Bounding  $\mathbb{E}[T_3]$ : Use a similar trick as in bounding  $\mathbb{E}[T_1]$ , and we get

$$\mathbb{E}\left[\mathbf{T}_{3}\right] = \mathbb{E}\left[\left\|\sum_{q=0}^{t-1} \nabla \boldsymbol{F}^{(q)} \left(\boldsymbol{\Pi}_{\ell=q}^{t-1} W^{(\ell)} - \mathbf{J}\right)\right\|_{\mathrm{F}}^{2}\right]$$

$$\leq \frac{2\sqrt{\rho}}{1 - \sqrt{\rho}} \sum_{q=0}^{t-1} \sqrt{\rho^{t-q}} \mathbb{E}\left[\left\|\nabla \boldsymbol{F}^{(q)}\right\|_{\mathrm{F}}^{2}\right].$$

For the last term, we have

$$\begin{split} &\frac{1}{mT} \sum_{t=0}^{T-1} \sum_{q=0}^{t-1} \sqrt{\rho}^{t-q} \mathbb{E} \left[ \| \nabla \boldsymbol{F}^{(q)} \|_{\mathrm{F}}^{2} \right] = \frac{1}{mT} \sum_{k=0}^{T-1} \mathbb{E} \left[ \| \nabla \boldsymbol{F}^{(t)} \|_{\mathrm{F}}^{2} \right] \sum_{q=1}^{T-1-t} \sqrt{\rho}^{q} \\ &\leq \frac{\sqrt{\rho}}{mT \left( 1 - \sqrt{\rho} \right)} \sum_{t=0}^{T-1} \mathbb{E} \left[ \| \nabla \boldsymbol{F}^{(t)} \|_{\mathrm{F}}^{2} \right]. \end{split}$$

g) Putting them together.:

$$\begin{split} &\frac{1}{mT} \sum_{t=0}^{T-1} \mathbb{E} \left[ \| \boldsymbol{X}^{(t)} \left( \mathbf{I} - \mathbf{J} \right) \|_{\mathrm{F}}^{2} \right] \\ &\leq \sigma^{2} \rho \left[ \frac{12 \kappa^{2} \eta^{4} \binom{s}{2}^{2} L^{2}}{\left( 1 - \sqrt{\rho} \right)^{2}} + \frac{3 \eta^{2} s^{2}}{1 - \rho} \right] + \left[ 2 \kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} + 1 \right] \frac{6 \eta^{2} s^{2} \sqrt{\rho}}{1 - \sqrt{\rho}} \frac{1}{mT} \sum_{t=0}^{T-1} \sum_{q=0}^{t-1} \sqrt{\rho^{t-q}} \mathbb{E} \left[ \| \nabla \boldsymbol{F}^{(q)} \|_{\mathrm{F}}^{2} \right] \\ &\leq \sigma^{2} \rho \left[ \frac{12 \kappa^{2} \eta^{4} \binom{s}{2}^{2} L^{2}}{\left( 1 - \sqrt{\rho} \right)^{2}} + \frac{3 \eta^{2} s^{2}}{1 - \rho} \right] + \left[ 2 \kappa^{2} \eta^{2} \binom{s}{2}^{2} L^{2} + 1 \right] \frac{6 \eta^{2} s^{2} \rho}{mT \left( 1 - \sqrt{\rho} \right)^{2}} \sum_{t=0}^{T-1} \mathbb{E} \left[ \| \nabla \boldsymbol{F}^{(t)} \|_{\mathrm{F}}^{2} \right]. \end{split}$$

We know that

$$2\kappa^2\eta^2\binom{s}{2}^2L^2+1 \leq \frac{L^2\kappa^2s^2\eta^2}{2}s^2+1 \leq s^2+1 \leq 2s^2.$$

Put all the parts together, we get

$$\frac{1}{mT}\sum_{k=0}^{T-1}\mathbb{E}\left[\left\|\boldsymbol{X}^{(t)}\left(\mathbf{I}-\mathbf{J}\right)\right\|_{\mathrm{F}}^{2}\right]\leq 6s^{2}\eta^{2}\sigma^{2}\rho\left[\frac{2}{\left(1-\sqrt{\rho}\right)^{2}}+\frac{1}{1-\rho}\right]+\frac{72\xi^{2}\eta^{2}s^{4}\rho}{\left(1-\sqrt{\rho}\right)^{2}}+\frac{72\left(\beta^{2}+1\right)\eta^{2}s^{4}\rho}{\left(1-\sqrt{\rho}\right)^{2}}\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\left[\left\|\nabla F(\bar{\boldsymbol{x}}^{t})\right\|_{2}^{2}\right],$$

which follows from the step size  $\eta$ 

$$\frac{36L^2\eta^2 s^4 \rho}{\left(1 - \sqrt{\rho}\right)^2} \le \frac{1}{2}.$$

**Proof of Theorem 1.** By taking an extra expectation over the remaining randomness and telescoping sum, we get

$$\begin{split} &\frac{F^* - F(\bar{x}^0)}{T} \\ &\leq -\left[\frac{s\eta}{4} - 3\eta^2 s^2 \left(\beta^2 + 1\right) \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right]\right] \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\left\|\nabla F(\bar{x}^t)\right\|_2^2\right] \\ &+ 3\xi^2 \eta^2 s^2 \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right] \\ &+ \left\{s\eta L^2 + 3\eta^2 s^2 L^2 \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right]\right\} \left\{6s^2 \eta^2 \sigma^2 \left[\frac{2\rho}{\left(1 - \sqrt{\rho}\right)^2} + \frac{\rho}{1 - \rho}\right] + \frac{72\xi^2 \eta^2 s^4 \rho}{\left(1 - \sqrt{\rho}\right)^2}\right\} \\ &+ \left\{s\eta L^2 + 3\eta^2 s^2 L^2 \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right]\right\} \left[\frac{72 \left(\beta^2 + 1\right) \eta^2 s^4 \rho}{\left(1 - \sqrt{\rho}\right)^2} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\left\|\nabla F(\bar{x}^t)\right\|_2^2\right] \\ &= -s\eta \left\{\frac{1}{4} - 3\eta s \left(\beta^2 + 1\right) \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right] \left[1 + \frac{72\eta^2 L^2 s^4 \rho}{\left(1 - \sqrt{\rho}\right)^2}\right] - \frac{72 \left(\beta^2 + 1\right) L^2 \eta^2 s^4 \rho}{\left(1 - \sqrt{\rho}\right)^2}\right\} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\left\|\nabla F(\bar{x}^t)\right\|_2^2\right] \\ &+ \eta^2 s^2 \left\{\left[\kappa^2 L^2 + 2L \left(\frac{1}{m} + \frac{\kappa^2 L^2}{4}\right)\right] + 6\left\{s\eta L^2 + 3\eta^2 s^2 L^2 \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right]\right\} \left[\frac{2\rho}{\left(1 - \sqrt{\rho}\right)^2} + \frac{\rho}{1 - \rho}\right]\right\} \sigma^2 \\ &+ 3\eta^2 s^2 \left\{\left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right] + \left\{s\eta L^2 + 3\eta^2 s^2 L^2 \left[\kappa^2 L^2 + 2L \left(1 + \frac{\kappa^2 L^2}{4}\right)\right]\right\} \frac{24s^2 \rho}{\left(1 - \sqrt{\rho}\right)^2}\right\} \xi^2 \end{split}$$

What follows refines the choice of the step-size:

$$\frac{1}{4} - 3\eta s \left(\beta^{2} + 1\right) \left[\kappa^{2} L^{2} + 2L\left(1 + \frac{\kappa^{2} L^{2}}{4}\right)\right] \left[1 + \frac{72\eta^{2} L^{2} s^{4} \rho}{\left(1 - \sqrt{\rho}\right)^{2}}\right] - \frac{72\left(\beta^{2} + 1\right) L^{2} \eta^{2} s^{4} \rho}{\left(1 - \sqrt{\rho}\right)^{2}} \\
\stackrel{\text{(a)}}{\geq} \frac{1}{4} - 3\eta s \left(\beta^{2} + 1\right) \left[\kappa^{2} L^{2} + 2L\left(1 + \frac{\kappa^{2} L^{2}}{4}\right)\right] \left[1 + \frac{18s^{2} L^{2} \rho}{\left(1 - \sqrt{\rho}\right)^{2}}\right] - \frac{18s^{2} L^{2} \left(\beta^{2} + 1\right) \rho}{\left(1 - \sqrt{\rho}\right)^{2}} \stackrel{\text{(b)}}{\geq} \frac{1}{8},$$

where (a) follows because  $\eta s \leq \frac{1}{2}$  , while (b) because

$$\eta \leq \frac{1}{24\left(\beta^2 + 1\right)\left[\kappa^2 L^2 + 2L\left(1 + \frac{\kappa^2 L^2}{4}\right)\right]\left[1 + \frac{18s^2 L^2 \rho}{\left(1 - \sqrt{\rho}\right)^2}\right] + \frac{144(\beta^2 + 1)s^2 L^2 \rho}{\left(1 - \sqrt{\rho}\right)^2}}.$$

$$\begin{split} & \left[\kappa^2 L^2 + 2L\left(\frac{1}{m} + \frac{\kappa^2 L^2}{4}\right)\right] + 6\left[s\eta L^2 + 3\eta^2 s^2 L^2\left(\kappa^2 L^2 + 2L\left(1 + \frac{\kappa^2 L^2}{4}\right)\right)\right] \left[\frac{2\rho}{\left(1 - \sqrt{\rho}\right)^2} + \frac{\rho}{1 - \rho}\right] \\ & \leq \left[\kappa^2 L^2 + 2L\left(\frac{1}{m} + \frac{\kappa^2 L^2}{4}\right)\right] + 6\left[s\eta L^2 + 3L^2\left(\kappa^2 L^2 + 2L\left(1 + \frac{\kappa^2 L^2}{4}\right)\right)\right] \left[\frac{2\rho}{\left(1 - \sqrt{\rho}\right)^2} + \frac{\rho}{1 - \rho}\right]. \end{split}$$

In addition, we need to ensure that  $\eta \rho s^3 \leq 1$ , with such an additional choice, we get

$$\left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] + \left\{ s \eta L^2 + 3 \eta^2 s^2 L^2 \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] \right\} \frac{24 s^2 \rho}{\left( 1 - \sqrt{\rho} \right)^2}$$

$$\leq \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] + \left\{ 1 + 3 \eta s \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] \right\} \frac{24 L^2}{\left( 1 - \sqrt{\rho} \right)^2}.$$

A little rearrangement, and applying the fact that

$$1 - \rho = (1 - \sqrt{\rho})(1 + \sqrt{\rho}) \ge (1 - \sqrt{\rho})^2$$

we arrive at

$$\begin{split} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 \right] &\leq \frac{8F(\bar{\boldsymbol{x}}^0) - 8F^\star}{s\eta T} \\ &+ 8\eta s \kappa^2 L^2 \left( 1 + \frac{L}{2} \right) \sigma^2 \\ &+ \frac{144\rho \left( \eta s \right)^2}{\left( 1 - \sqrt{\rho} \right)^2} \left\{ L^2 + 3L^2 \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] \right\} \sigma^2 \\ &+ \frac{16\eta s L \sigma^2}{m} \\ &+ 24\eta s \left\{ \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] + \frac{24L^2}{\left( 1 - \sqrt{\rho} \right)^2} \right\} \xi^2 \\ &+ \frac{1728L^2 \left( \eta s \right)^2}{\left( 1 - \sqrt{\rho} \right)^2} \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] \xi^2. \end{split}$$

Choose the step size to be  $\eta=c_1\sqrt{\frac{m}{sT}}.$  When T is sufficiently large such that

$$\eta \leq \min \left\{ \frac{1}{24 \left(\beta^2+1\right) \left[\kappa^2 L^2+2 L \left(1+\frac{\kappa^2 L^2}{4}\right)\right] \left[1+\frac{18 s^2 L^2 \rho}{\left(1-\sqrt{\rho}\right)^2}\right]+\frac{144 \left(\beta^2+1\right) s^2 L^2 \rho}{\left(1-\sqrt{\rho}\right)^2}}, \frac{1}{2 s}, \frac{\sqrt{2}}{\kappa s L}, \frac{1}{\rho s^3}, \frac{1-\sqrt{\rho}}{6 \sqrt{2 \rho} L s^2}}\right\},$$

we have

$$\begin{split} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla F(\bar{\boldsymbol{x}}^t) \right\|_2^2 \right] &\leq O\left\{ \frac{8F(\bar{\boldsymbol{x}}^0) - 8F^\star}{\sqrt{msT}} \right. \\ &+ 8\kappa^2 L^2 \left( 1 + \frac{L}{2} \right) \sigma^2 \sqrt{\frac{ms}{T}} \\ &+ \frac{144\rho L^2}{\left( 1 - \sqrt{\rho} \right)^2} \left\{ 1 + 3 \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] \right\} \sigma^2 \frac{ms}{T} \\ &+ 16L\sigma^2 \sqrt{\frac{s}{mT}} \\ &+ 24 \left\{ \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] + \frac{24L^2}{\left( 1 - \sqrt{\rho} \right)^2} \right\} \xi^2 \sqrt{\frac{ms}{T}} \\ &+ \frac{1728L^2}{\left( 1 - \sqrt{\rho} \right)^2} \left[ \kappa^2 L^2 + 2L \left( 1 + \frac{\kappa^2 L^2}{4} \right) \right] \xi^2 \frac{ms}{T} \right\}. \end{split}$$