

Asynchronous Decentralized Parallel Stochastic Gradient Descent

Xiangru Lian^{*1}, Wei Zhang^{*2}, Ce Zhang³, and Ji Liu^{1, 4}

xiangru@yandex.com, weiz@us.ibm.com, ce.zhang@inf.ethz.ch,
ji.liu.uwisc@gmail.com

¹Department of Computer Science, University of Rochester

²IBM T. J. Watson Research Center

³Department of Computer Science, ETH Zurich

⁴Tencent AI Lab

February 10, 2018

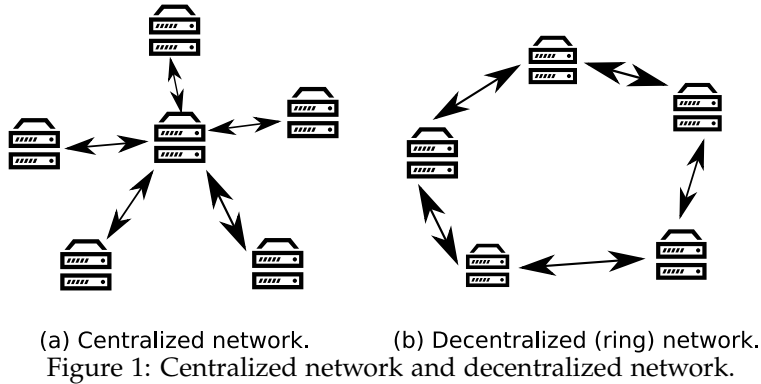
Abstract

Most commonly used distributed machine learning systems are either synchronous or centralized asynchronous. Synchronous algorithms like AllReduce-SGD perform poorly in a heterogeneous environment, while asynchronous algorithms using a parameter server suffer from 1) communication bottleneck at parameter servers when workers are many, and 2) significantly worse convergence when the traffic to parameter server is congested. Can we design an algorithm that is robust in a heterogeneous environment, while being communication efficient and maintaining the best-possible convergence rate? In this paper, we propose an asynchronous decentralized stochastic gradient descent algorithm (AD-PSGD) satisfying all above expectations. Our theoretical analysis shows AD-PSGD converges at the optimal $O(1/\sqrt{K})$ rate as SGD and has linear speedup w.r.t. number of workers. Empirically, AD-PSGD outperforms the best of decentralized parallel SGD (D-PSGD), asynchronous parallel SGD (A-PSGD), and standard data parallel SGD (AllReduce-SGD), often by orders of magnitude in a heterogeneous environment. When training ResNet-50 on ImageNet with up to 128 GPUs, AD-PSGD converges (w.r.t epochs) similarly to the AllReduce-SGD, but each epoch can be up to $4-8\times$ faster than its synchronous counterparts in a network-sharing HPC environment. To the best of our knowledge, AD-PSGD is the first asynchronous algorithm that achieves a similar epoch-wise convergence rate as AllReduce-SGD, at an over 100-GPU scale.

1 Introduction

It often takes hours to train large deep learning tasks such as ImageNet, even with hundreds of GPUs Goyal et al. [2017]. At this scale, how workers communicate becomes a crucial design choice. Most existing systems such as TensorFlow [Abadi et al., 2016], MXNet [Chen et al., 2015], and CNTK [Seide and Agarwal, 2016] support two communication modes: (1) synchronous communication via parameter servers or AllReduce, or (2) asynchronous communication via parameter servers. When there are stragglers (i.e., slower workers) in the system, which is common especially at the scale of hundreds devices, asynchronous approaches are more robust. However, most asynchronous implementations have a *centralized* design, as illustrated in Figure 1(a) — a central server holds the shared model for all other workers. Each worker

^{*}Contributed equally.



	Communication complexity (n.t./n.h.) ^a	Idle time
S-PSGD [Ghadimi et al., 2016]	Long ($O(n)/O(n)$)	Long
A-PSGD [Lian et al., 2015]	Long ($O(n)/O(n)$)	Short
AllReduce-SGD [Luehr, 2016]	Medium ($O(1)/O(n)$)	Long
D-PSGD [Lian et al., 2017]	Short ($O(\deg(G))/O(\deg(G))$)	Long
AD-PSGD (this paper)	Short ($O(\deg(G))/O(\deg(G))$)	Short

Table 1: Comparison of different distributed machine learning algorithms on a network graph G . *Long idle time* means in each iteration the whole system needs to wait for the slowest worker. *Short idle time* means the corresponding algorithm breaks this synchronization per iteration. Note that if G is a ring network as required in AllReduce-SGD, $O(\deg(G)) = O(1)$.

^an.t. means number of gradients/models transferred at the busiest worker per n (minibatches of) stochastic gradients updated. n.h. means number of handshakes at the busiest worker per n (minibatches of) stochastic gradients updated.

calculates its own gradients and updates the shared model asynchronously. The parameter server may become a communication bottleneck and slow down the convergence. We focus on the question: *Can we remove the central server bottleneck in asynchronous distributed learning systems while maintaining the best possible convergence rate?*

Recent work [Lian et al., 2017] shows that *synchronous decentralized* parallel stochastic gradient descent (D-PSGD) can achieve comparable convergence rate as its centralized counterparts without any central bottleneck. Figure 1-(b) illustrates one communication topology of D-PSGD in which each worker only talks to its neighbors. However, the synchronous nature of D-PSGD makes it vulnerable to stragglers because of the synchronization barrier at each iteration among *all* workers. *Is it possible to get the best of both worlds of asynchronous SGD and decentralized SGD?*

In this paper, we propose the *asynchronous decentralized* parallel stochastic gradient decent algorithm (AD-PSGD) that is theoretically justified to keep the advantages of both asynchronous SGD and decentralized SGD. In AD-PSGD, workers do not wait for all others and only communicate in a decentralized fashion. AD-PSGD can achieve linear speedup with respect to the number of workers and admit a convergence rate of $O(1/\sqrt{K})$, where K is the number of updates. This rate is consistent with D-PSGD and centralized parallel SGD. By design, AD-PSGD enables wait-free computation and communication, which ensures *AD-PSGD always converges better (w.r.t epochs or wall time) than D-PSGD as the former allows much more frequent information exchanging.*

In practice, we found that AD-PSGD is particularly useful in heterogeneous computing environments such as cloud-computing, where computing/communication devices' speed often varies. We implement AD-PSGD in Torch and MPI and evaluate it on an IBM S822LC cluster of up to 128 P100 GPUs. We show that, on real-world datasets such as ImageNet, AD-PSGD has the same empirical convergence rate as its

centralized and/or synchronous counterpart. In heterogeneous environments, AD-PSGD can be faster than its fastest synchronous counterparts by orders of magnitude. On an HPC cluster with homogeneous computing devices but shared network, AD-PSGD can still outperform its synchronous counterparts by 4X-8X.

Both the theoretical analysis and system implementations of AD-PSGD are non-trivial, and they form the two technical contributions of this work.

2 Related work

We review related work in this section. In the following, K and n refer to the number of iterations and the number of workers, respectively. A comparison of the algorithms can be found in Table 1.

The *Stochastic Gradient Descent (SGD)* Ghadimi and Lan [2013], Moulines and Bach [2011], Nemirovski et al. [2009] is a powerful approach to solve large scale machine learning problems, with the optimal convergence rate $O(1/\sqrt{K})$ on nonconvex problems.

For *Synchronous Parallel Stochastic Gradient Descent (S-PSGD)*, every worker fetches the model saved in a parameter server and computes a minibatch of stochastic gradients. Then they push the stochastic gradients to the parameter server. The parameter server synchronizes all the stochastic gradients and update their average into the model saved in the parameter server, which completes one iteration. The convergence rate is proved to be $O(1/\sqrt{nK})$ on nonconvex problems [Ghadimi et al., 2016]. Results on convex objectives can be found in Dekel et al. [2012]. Due to the synchronization step, all other workers have to stay idle to wait for the slowest one. In each iteration the parameter server needs to synchronize $O(n)$ workers, which causes high communication cost at the parameter server especially when n is large.

The *Asynchronous Parallel Stochastic Gradient Descent (A-PSGD)* [Agarwal and Duchi, 2011, Feysmahdavian et al., 2016, Paine et al., 2013, Recht et al., 2011] breaks the synchronization in S-PSGD by allowing workers to use stale weights to compute gradients. Asynchronous algorithms significantly reduce the communication overhead by avoiding idling any worker and can still work well when part of the computing workers are down. On nonconvex problems, when the staleness of the weights used is upper bounded, A-PSGD is proved to admit the same convergence rate as S-PSGD [Lian et al., 2015, 2016].

In *AllReduce Stochastic Gradient Descent implementation (AllReduce-SGD)* [Luehr, 2016, MPI contributors, 2015, Patarasuk and Yuan, 2009], the update rule per iteration is exactly the same as in S-PSGD, so they share the same convergence rate. However, there is no parameter server in AllReduce-SGD. The workers are connected with a ring network and each worker keeps the same local copy of the model. In each iteration, each worker calculates a minibatch of stochastic gradients. Then all the workers use AllReduce to synchronize the stochastic gradients, after which each worker will get the average of all stochastic gradients. In this procedure, only $O(1)$ amount of gradient is sent/received per worker, but $O(n)$ handshakes are needed on each worker. This makes AllReduce slow on high latency network. At the end of the iteration the averaged gradient is updated into the local model of each worker. Since we still have synchronization in each iteration, the idle time is still high as in S-PSGD.

In *Decentralized Parallel Stochastic Gradient Descent (D-PSGD)* [Lian et al., 2017], all workers are connected with a network that forms a connected graph G . Every worker has its local copy of the model. In each iteration, all workers compute stochastic gradients locally and at the same time average its local model with its neighbors. Finally the locally computed stochastic gradients are updated into the local models. In this procedure, the busiest worker only sends/receives $O(\deg(G))$ models and has $O(\deg(G))$ handshakes per iteration. Note that in D-PSGD the computation and communication can be done in parallel, which means, when communication time is smaller than the computation time, the communication can be completely hidden. The idle time is still high in D-PSGD because all workers need to finish updating before stepping into the next iteration. Before Lian et al. [2017] there are also previous studies on decentralized stochastic algorithms (both synchronous and asynchronous versions) though *none of*

them is proved to have speedup when the number of workers increases. For example, Lan et al. [2017] proposed a decentralized stochastic primal-dual type algorithm with a computational complexity of $O(n/\epsilon^2)$ for general convex objectives and $O(n/\epsilon)$ for strongly convex objectives. Sirb and Ye [2016] proposed an asynchronous decentralized stochastic algorithm with a $O(n/\epsilon^2)$ complexity for convex objectives. These bounds do not imply any speedup for decentralized algorithms. Bianchi et al. [2013] proposed a similar decentralized stochastic algorithm. The authors provided a convergence rate for the consensus of the local models when the local models are bounded. The convergence to a solution was provided by using central limit theorem. However, they did not provide the convergence rate to the solution. Ram et al. [2010] proposed an asynchronous subgradient variations of the decentralized stochastic optimization algorithm for convex problems. The asynchrony was modeled by viewing the update event as a Poisson process and the convergence to the solution was shown. Srivastava and Nedic [2011], Sundhar Ram et al. [2010] are similar. The main differences from this work are 1) we take the situation where a worker calculates gradients based on old model into consideration, which is the case in the asynchronous setting; 2) we prove that our algorithm can achieve linear speedup when we increase the number of workers, which is important if we want to use the algorithm to accelerate training; 3) Our implementation guarantees deadlock-free, wait-free computation and communication.

We next briefly review *decentralized algorithms*. Decentralized algorithms were initially studied by the control community for solving the consensus problem where the goal is to compute the mean of all the data distributed on multiple nodes [Aysal et al., 2009, Boyd et al., 2005, Carli et al., 2010, Fagnani and Zampieri, 2008, Olfati-Saber et al., 2007, Schenato and Gamba, 2007]. For decentralized algorithms used for optimization problems, Lu et al. [2010] proposed two non-gradient-based algorithms for solving one-dimensional unconstrained convex optimization problems where the objective on each node is strictly convex, by calculating the inverse function of the derivative of the local objectives and transmitting the gradients or local objectives to neighbors, and the algorithms can be used over networks with time-varying topologies. A convergence rate was not shown but the authors did prove the algorithms will converge to the solution eventually. Mokhtari and Ribeiro [2016] proposed a fast decentralized variance reduced algorithm for strongly convex optimization problems. The algorithm is proved to have linear convergence rate and a nice stochastic saddle point method interpretation is given. However, the speedup property is unclear and a table of stochastic gradients need to be stored. Yuan et al. [2016] studied decentralized gradient descent on convex and strongly convex objectives. The algorithm in each iteration averages the models of the nodes with their neighbors' and then updates the full gradient of the local objective function on each node. The subgradient version was considered in Nedic and Ozdaglar [2009], Ram et al. [2009]. The algorithm is intuitive and easy to understand. However, the limitation of the algorithm is that it does not converge to the exact solution because the exact solution is not a fixed point of the algorithm's update rule. This issue was fixed later by Shi et al. [2015a], Wu et al. [2016] by using the gradients of last two instead of one iterates in each iteration, which was later improved in Li et al. [2017], Shi et al. [2015b] by considering proximal gradients. Decentralized ADMM algorithms were analyzed in Aybat et al. [2015], Shi et al., Zhang and Kwok [2014]. Wang et al. [2016] develops a decentralized algorithm for recursive least-squares problems.

3 Algorithm

We introduce the AD-PSGD algorithm in this section.

Definitions and notations Throughout this paper, we use the following notation and definitions:

- $\|\cdot\|$ denotes the vector ℓ_2 norm or the matrix spectral norm depending on the argument.
- $\|\cdot\|_F$ denotes the matrix Frobenius norm.

- $\nabla f(\cdot)$ denotes the gradient of a function f .
- $\mathbf{1}_n$ denotes the column vector in \mathbb{R}^n with 1 for all elements.
- f^* denotes the optimal solution to (1).
- $\lambda_i(\cdot)$ denotes the i -th largest eigenvalue of a matrix.
- e_i denotes the i th element of the standard basis of \mathbb{R}^n .

3.1 Problem definition

The decentralized communication topology is represented as an undirected graph: (V, E) , where $V := \{1, 2, \dots, n\}$ denotes the set of n workers and $E \subseteq V \times V$ is the set of the edges in the graph. Each worker represents a machine/gpu owning its local data (or a sensor collecting local data online) such that each worker is associated with a local loss function

$$f_i(x) := \mathbb{E}_{\xi \sim \mathcal{D}_i} F_i(x; \xi),$$

where \mathcal{D}_i is a distribution associated with the local data at worker i and ξ is a data point sampled via \mathcal{D}_i . The edge means that the connected two workers can exchange information. For the AD-PSGD algorithm, the overall optimization problem it solves is

$$\min_{x \in \mathbb{R}^N} f(x) := \mathbb{E}_{i \sim \mathcal{I}} f_i(x) = \sum_{i=1}^n p_i f_i(x), \quad (1)$$

where p_i 's define a distribution, that is, $p_i \geq 0$ and $\sum_i p_i = 1$, and p_i indicates the updating frequency of worker i or the percentage of the updates performed by worker i . The faster a worker, the higher the corresponding p_i . The intuition is that if a worker is faster than another worker, then the faster worker will run more epochs given the same amount of time, and consequently the corresponding worker has a larger impact.

Remark 1. To solve the common form of objectives in machine learning using AD-PSGD

$$\min_{x \in \mathbb{R}^N} \mathbb{E}_{\xi \sim \mathcal{D}} F(x; \xi),$$

we can appropriately distribute data such that Eq. (1) solves the target objective above:

Strategy-1 Let $\mathcal{D}_i = \mathcal{D}$ and \mathcal{D} , that is, all worker can access all data, and consequently $F_i(\cdot; \cdot) = F(\cdot; \cdot)$, that is, all $f_i(\cdot)$'s are the same;

Strategy-2 Split the data into all workers appropriately such that the portion of data is p_i on worker i and define \mathcal{D}_i to be the uniform distribution over the assigned data samples.

3.2 AD-PSGD algorithm

The AD-PSGD algorithm can be described in the following: each worker maintains a local model x in its local memory and (using worker i as an example) repeats the following steps:

- **Sample data:** Sample a mini-batch of training data denoted by $\{\xi_m^i\}_{m=1}^M$, where M is the batch size.
- **Compute gradients:** Use the sampled data to compute the stochastic gradient $\sum_{m=1}^M \nabla F(\hat{x}^i; \xi_m^i)$, where \hat{x}^i is read from the model in the local memory.
- **Gradient update:** Update the model in the local memory by $x^i \leftarrow x^i - \gamma \sum_{m=1}^M \nabla F(\hat{x}^i; \xi_m^i)$. Note that \hat{x}^i may not be the same as x^i as it may be modified by other workers in the **averaging** step.

- **Averaging:** Randomly select a neighbor (e.g. worker i') and average the local model with the worker i' 's model $x^{i'}$ (both models on both workers are updated to the averaged model). More specifically, $x^i, x^{i'} \leftarrow \frac{x^i}{2} + \frac{x^{i'}}{2}$.

Note that each worker runs the procedure above on its own without any global synchronization. This reduces the idle time of each worker and the training process will still be fast even if part of the network or workers slow down.

The **averaging** step can be generalized into the following update for all workers:

$$[x^1, x^2, \dots, x^n] \leftarrow [x^1, x^2, \dots, x^n]W$$

where W can be an arbitrary doubly stochastic matrix. This generalization gives plenty flexibility to us in implementation without hurting our analysis.

All workers run the procedure above simultaneously, as shown in Algorithm 1. We use a virtual counter k to denote the iteration counter – every single **gradient update** happens no matter on which worker will increase k by 1. i_k denotes the worker performing the k th update.

3.3 Implementation details

We briefly describe two interesting aspects of system designs and leave more discussions to Appendix A.

3.3.1 Deadlock avoidance

A naive implementation of the above algorithm may cause deadlock — the averaging step needs to be atomic and involves updating two workers (the selected worker and one of its neighbors). As an example, given three fully connected workers A , B , and C , A sends its local model x_A to B and waits for x_B from B ; B has already sent out x_B to C and waits for C 's response; and C has sent out x_C to A and waits for x_A from A .

We prevent the deadlock in the following way: The communication network is designed to be a bipartite graph, that is, the worker set V can be split into two disjoint sets A (active set) and P (passive set) such that any edge in the graph connects one worker in A and one worker in P . Due to the property of the bipartite graph, the neighbors of any active worker can only be passive workers and the neighbors of any passive worker can only be active workers. This implementation avoids deadlock but still fits in the general algorithm Algorithm 1 we are analyzing. We leave more discussions and a detailed implementation for wait-free training to Appendix A.

3.3.2 Communication topology

The simplest realization of AD-PSGD algorithm is a ring-based topology. To accelerate information exchanging, we also implement a communication topology in which each sender communicates with a receiver that is $2^i + 1$ hops away in the ring, where i is an integer from 0 to $\log(n - 1)$ (n is the number of learners). It is easy to see it takes at most $O(\log(n))$ steps for any pair of workers to exchange information instead of $O(n)$ in the simple ring-based topology. In this way, ρ (as defined in Section 4) becomes smaller and the scalability of AD-PSGD improves. This implementation also enables robustness against slow or failed network links because there are multiple routes for a worker to disseminate its information.

4 Theoretical analysis

In this section we provide theoretical analysis for the AD-PSGD algorithm. We will show that the convergence rate of AD-PSGD is consistent with SGD and D-PSGD.

Algorithm 1 AD-PSGD (logical view)

b

Require: Initialize local models $\{x_0^i\}_{i=1}^n$ with the same initialization, learning rate γ , batch size M , and total number of iterations K .

- 1: **for** $k = 0, 1, \dots, K - 1$ **do**
- 2: Randomly sample a worker i_k of the graph G and randomly sample an averaging matrix W_k which can be dependent on i_k .
- 3: Randomly sample a batch $\tilde{\zeta}_{k,i_k} := (\tilde{\zeta}_{k,1}^{i_k}, \tilde{\zeta}_{k,2}^{i_k}, \dots, \tilde{\zeta}_{k,M}^{i_k})$ from local data of the i_k -th worker.
- 4: Compute the stochastic gradient locally $g_k(\hat{x}_k^{i_k}; \tilde{\zeta}_k^{i_k}) := \sum_{j=1}^M \nabla F(\hat{x}_k^{i_k}; \tilde{\zeta}_{k,j}^{i_k})$.
- 5: Average local modes by ^a $[x_{k+1/2}^1, x_{k+1/2}^2, \dots, x_{k+1/2}^n] \leftarrow [x_k^1, x_k^2, \dots, x_k^n] W_k$
- 6: Update the local model $x_{k+1}^{i_k} \leftarrow x_{k+1/2}^{i_k} - \gamma g_k(\hat{x}_k^{i_k}; \tilde{\zeta}_k^{i_k})$ and $x_{k+1}^j \leftarrow x_{k+1/2}^j, \forall j \neq i_k$.
- 7: **end for**
- 8: Output the average of the models on all workers.

^aNote that Line 4 and Line 5 can run in parallel.

Note that by counting each update of stochastic gradients as one iteration, the update of each iteration in Algorithm 1 can be viewed as

$$X_{k+1} = X_k W_k - \gamma \partial g(\hat{X}_k; \tilde{\zeta}_k^{i_k}, i_k),$$

where k is the iteration number, x_k^i is the local model of the i th worker at the k th iteration, and

$$\begin{aligned} X_k &= \begin{bmatrix} x_k^1 & \cdots & x_k^n \end{bmatrix} \in \mathbb{R}^{N \times n}, \\ \hat{X}_k &= \begin{bmatrix} \hat{x}_k^1 & \cdots & \hat{x}_k^n \end{bmatrix} \in \mathbb{R}^{N \times n}, \\ \partial g(\hat{X}_k; \tilde{\zeta}_k^{i_k}, i_k) &= \begin{bmatrix} 0 & \cdots & 0 & \sum_{j=1}^M \nabla F(\hat{x}_k^{i_k}; \tilde{\zeta}_{k,j}^{i_k}) & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{N \times n}, \end{aligned}$$

and $\hat{X}_k = X_{k-\tau_k}$ for some nonnegative integer τ_k .

Assumption 1. Throughout this paper, we make the following commonly used assumptions:

1. **Lipschitzian gradient:** All functions $f_i(\cdot)$'s are with L -Lipschitzian gradients.
2. **Doubly stochastic averaging:** W_k is doubly stochastic for all k .
3. **Spectral gap:** There exists a $\rho \in [0, 1)$ such that

$$\max\{|\lambda_2(\mathbb{E}[W_k^\top W_k])|, |\lambda_n(\mathbb{E}[W_k^\top W_k])|\} \leq \rho, \forall k. \quad (2)$$

4. **Unbiased estimation:**¹

$$\mathbb{E}_{\tilde{\zeta} \sim \mathcal{D}_i} \nabla F(x; \tilde{\zeta}) = \nabla f_i(x), \quad (3)$$

$$\mathbb{E}_{i \sim \mathcal{I}} \mathbb{E}_{\tilde{\zeta} \sim \mathcal{D}_i} \nabla F(x; \tilde{\zeta}) = \nabla f(x). \quad (4)$$

5. **Bounded variance:** Assume the variance of the stochastic gradient

$$\mathbb{E}_{i \sim \mathcal{I}} \mathbb{E}_{\tilde{\zeta} \sim \mathcal{D}_i} \|\nabla F(x; \tilde{\zeta}) - \nabla f(x)\|^2$$

¹Note that this is easily satisfied when all workers can access all data so that $\mathbb{E}_{\tilde{\zeta} \sim \mathcal{D}_i} \nabla F(x; \tilde{\zeta}) = \nabla f(x)$. When each worker can only access part of the data, we can also meet these assumptions by appropriately distributing data.

is bounded for any x with i sampled from the distribution \mathcal{I} and ξ from the distribution \mathcal{D}_i . This implies there exist constants σ and ζ such that

$$\mathbb{E}_{\xi \sim \mathcal{D}_i} \|\nabla F(x, \xi) - \nabla f_i(x)\|^2 \leq \sigma^2, \forall i, \forall x. \quad (5)$$

$$\mathbb{E}_{i \sim \mathcal{I}} \|\nabla f_i(x) - \nabla f(x)\|^2 \leq \zeta^2, \forall x. \quad (6)$$

Note that if all workers can access all data, then $\zeta = 0$.

6. **Dependence of random variables:** $\xi_k, i_k, k \in \{0, 1, 2, \dots\}$ are independent random variables. W_k is a random variable dependent on i_k .

7. **Bounded staleness:** $\hat{X}_k = X_{k-\tau_k}$ and there exists a constant T such that $\max_k \tau_k \leq T$.

Throughout this paper, we define the following notations for simpler notation

$$\bar{\rho} := \frac{n-1}{n} \left(\frac{1}{1-\rho} + \frac{2\sqrt{\rho}}{(1-\sqrt{\rho})^2} \right),$$

$$C_1 := 1 - 24M^2L^2\gamma^2 \left(T \frac{n-1}{n} + \bar{\rho} \right),$$

$$C_2 := \frac{\gamma M}{2n} - \frac{\gamma^2 LM^2}{n^2} - \frac{2M^3L^2T^2\gamma^3}{n^3} - \left(\frac{6\gamma^2L^3M^2}{n^2} + \frac{\gamma M}{n}L^2 + \frac{12M^3L^4T^2\gamma^3}{n^3} \right) \frac{4M^2\gamma^2(T \frac{n-1}{n} + \bar{\rho})}{C_1},$$

$$C_3 := \frac{1}{2} + 2C_1^{-1} \left(6\gamma^2L^2M^2 + \gamma nML + \frac{12M^3L^3T^2\gamma^3}{n} \right) \bar{\rho} + \frac{LT^2\gamma M}{n}.$$

Under Assumption 1 we have the following results:

Theorem 1 (Main theorem). While $C_3 \leq 1$ and $C_2 \geq 0$ and $C_1 > 0$ are satisfied we have

$$\frac{\sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2}{K} \leq \frac{2(\mathbb{E}f(x_0) - \mathbb{E}f^*)n}{\gamma KM} + \frac{2\gamma L}{n}(\sigma^2 + 6M\zeta^2).$$

Noting that $\frac{X_n \mathbf{1}_n}{n} = \frac{1}{n} \sum_{i=1}^n x_k^i$, this theorem characterizes the convergence of the average of all local models. By appropriately choosing the learning rate, we obtain the following corollary

Corollary 2. Let $\gamma = \frac{n}{10ML + \sqrt{\sigma^2 + 6M\zeta^2} \sqrt{KM}}$. We have the following convergence rate

$$\frac{\sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2}{K} \leq \frac{20(f(x_0) - f^*)L}{K} + \frac{2(f(x_0) - f^* + L)\sqrt{\sigma^2/M + 6\zeta^2}}{\sqrt{K}} \quad (7)$$

if the total number of iterations is sufficiently large, in particular,

$$K \geq \frac{ML^2n^2}{\sigma^2 + 6M\zeta^2} \max \left\{ 192 \left(T \frac{n-1}{n} + \bar{\rho} \right), \frac{64T^4}{n^2}, 1024n^2\bar{\rho}^2, \frac{(8\sqrt{6}T^{2/3} + 8)^2 (T + \bar{\rho} \frac{n-1}{n})^{2/3} (n-1)^{1/2}}{n^{1/6}} \right\}. \quad (8)$$

This corollary indicates that if the iteration number is big enough, AP-PSGD's convergence rate is $O(1/\sqrt{K})$. We compare the convergence rate of AD-PSGD with existing results for SGD and D-PSGD to show the tightness of the proved convergence rate. We will also show the efficiency and the linear speedup property for AD-PSGD w.r.t. batch size, number of workers, and staleness respectively. Further discussions on communication topology and intuition will be provided at the end of this section.

Remark 2 (Consistency with SGD). *Note that if $T = 0$ and $n = 1$ the proposed AD-PSGD reduces to the vanilla SGD algorithm Ghadimi and Lan [2013], Moulines and Bach [2011], Nemirovski et al. [2009]. Since $n = 1$, we do not have the variance among workers, that is, $\zeta = 0$, the convergence rate becomes $O(1/K + \sigma/\sqrt{KM})$ which is consistent with the convergence rate with SGD.*

Remark 3 (Linear speedup w.r.t. batch size). *When K is large enough the second term on the RHS of (7) dominates the first term. Note that the second term converges at a rate $O(1/\sqrt{MK})$ if $\zeta = 0$, which means the convergence efficiency gets boosted with a linear rate if increase the mini-batch size. This observation indicates the linear speedup w.r.t. the batch size and matches the results of mini-batch SGD.²*

Remark 4 (Linear speedup w.r.t. number of workers). *Note that every single stochastic gradient update counts one iteration in our analysis and our convergence rate in Corollary 2 is consistent with SGD / mini-batch SGD. It means that the number of required stochastic gradient updates to achieve a certain precision is consistent with SGD / mini-batch SGD, as long as the total number of iterations is large enough. It further indicates the linear speedup with respect to the number of workers n (n workers will make the iteration number advance n times faster in the sense of wall-clock time, which means we will converge n times faster). To the best of our knowledge, the linear speedup property w.r.t. to the number of workers for decentralized algorithms has not been recognized until the recent analysis for D-PSGD by Lian et al. [2017]. Our analysis reveals that by breaking the synchronization AD-PSGD can maintain linear speedup, reduce the idle time, and improve the robustness in heterogeneous computing environments.*

Remark 5 (Linear speedup w.r.t. the staleness). *From (8) we can also see that as long as the staleness T is bounded by $O(K^{1/4})$ (if other parameters are considered to be constants), linear speedup is achievable.*

5 Experiments

We describe our experimental methodologies in Section 5.1 and we evaluate the AD-PSGD algorithm in the following sections:

- Section 5.2: Compare AD-PSGD’s convergence rate (w.r.t epochs) with other algorithms.
- Section 5.3: Compare AD-PSGD’s convergence rate (w.r.t runtime) and its speedup with other algorithms.
- Section 5.4: Compare AD-PSGD’s robustness to other algorithms in heterogeneous computing and heterogeneous communication environments.
- Section 5.5: Evaluate AD-PSGD on ImageNet and ResNet-50 model.
- Appendix B: Evaluate AD-PSGD on IBM proprietary natural language processing dataset and model.

5.1 Experiments methodology

5.1.1 Dataset, model, and software

We use CIFAR10 and ImageNet-1K as the evaluation dataset and we use Torch-7 as our deep learning framework. We use MPI to implement the communication scheme.

For CIFAR10, we evaluate both VGG [Simonyan and Zisserman, 2015] and ResNet-20 [He et al., 2016] models. VGG, whose size is about 60MB, represents a communication intensive workload and ResNet-20, whose size is about 1MB, represents a computation intensive workload.

² Note that when $\zeta^2 \neq 0$, AD-PSGD does not admit this linear speedup w.r.t. batch size. It is unavoidable because increasing the minibatch size only decreases the variance of the stochastic gradients within each worker, while ζ^2 characterizes the variance of stochastic gradient among different workers, independent of the batch size.

Table 2: Testing accuracy comparison for VGG and ResNet-20 model on CIFAR10. 16 workers in total.

	AllReduce	D-PSGD	EAMSGD	AD-PSGD
VGG	87.04%	86.48%	85.75%	88.58%
ResNet-20	90.72%	90.81%	89.82%	91.49%

For the ImageNet-1K dataset, we use the ResNet-50 model whose size is about 100MB.

Additionally, we experimented on an IBM proprietary natural language processing datasets and models Zhang et al. [2017] in Appendix B.

5.1.2 Hardware

We evaluate AD-PSGD in two different environments:

- IBM S822LC HPC cluster: Each node with 4 Nvidia P100 GPUs, 160 Power8 cores (8-way SMT) and 500GB memory on each node. 100Gbit/s Mellanox EDR infiniband network. We use 32 such nodes.
- x86-based cluster: This cluster is a cloud-like environment with 10Gbit/s ethernet connection. Each node has 4 Nvidia P100 GPUs, 56 Xeon E5-2680 cores (2-way SMT), and 1TB DRAM. We use 4 such nodes.

5.1.3 Compared algorithms

We compare the proposed AD-PSGD algorithm to AllReduce-SGD, D-PSGD Lian et al. [2017] and a state of the art asynchronous SGD implementation EAMSGD. Zhang et al. [2015]³ In EAMSGD, each worker can communicate with the parameter server less frequently by increasing the “communication period” parameter su .

5.2 Convergence w.r.t. epochs

Figure 2 plots training loss w.r.t. epochs for each algorithm, which is evaluated for VGG and ResNet-20 models on CIFAR10 dataset with 16 workers. Table 2 reports the test accuracy of all algorithms.

For EAMSGD, we did extensive hyper-parameter tuning to get the best possible model, where $su = 1$. We set momentum moving average to be $0.9/n$ (where n is the number of workers) as recommended in Zhang et al. [2015] for EAMSGD.

For other algorithms, we use the following hyper-parameter setup as prescribed in Zagoruyko [2015] and FAIR [2017]:

- Batch size: 128 for VGG, 32 for ResNet-20.
- Learning rate: For VGG start from 1 and reduce by half every 25 epochs. For ResNet-20 start from 0.1 and decay by a factor of 10 at the 81st epoch and the 122nd epoch.
- Momentum: 0.9.
- Weight decay: 10^{-4} .

Figure 2 show that w.r.t epochs, AllReduce-SGD, D-PSGD and AD-PSGD converge similar, while ASGD converges worse. Table 2 shows AD-PSGD does not sacrifice test accuracy.

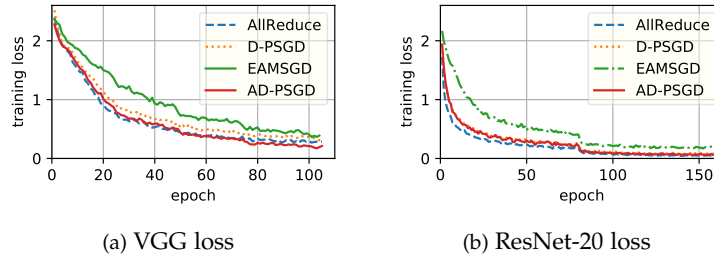


Figure 2: Training loss comparison for VGG and ResNet-20 model on CIFAR10. AllReduce-SGD, D-PSGD and AD-PSGD converge alike, EAMSGD converges the worst. 16 workers in total.

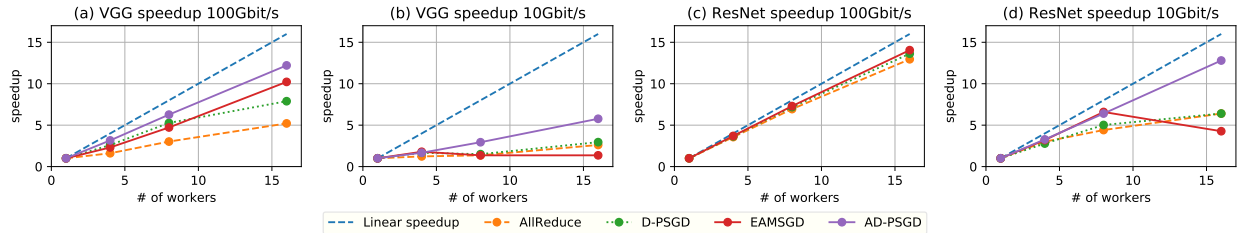


Figure 3: Runtime comparison for VGG (communication intensive) and ResNet-20 (computation intensive) models on CIFAR10. Experiments run on IBM HPC w/ 100Gbit/s network links and on x86 system w/ 10Gbit/s network links. AD-PSGD consistently converges the fastest. 16 workers in total.

5.3 Speedup and convergence w.r.t runtime

Figure 3 shows the runtime convergence results on both IBM HPC and x86 system. The EAMSGD implementation deploys parameter server sharding to mitigate the network bottleneck at the parameter servers. However, the central parameter server quickly becomes a bottleneck on a slow network with a large model as shown in Figure 3-(b).

Figure 5 shows the speedup for different algorithms w.r.t. number of workers. The speedup for ResNet-20 is better than VGG because ResNet-20 is a computation intensive workload.

Above results show that regardless of workload type (computation intensive or communication intensive) and communication networks (fast or slow), AD-PSGD consistently converges the fastest w.r.t. runtime and achieves the best speedup.

5.4 Robustness in a heterogeneous environment

In a heterogeneous environment, the speed of computation device and communication device may often vary, subject to architectural features (e.g., over/under-clocking, caching, paging), resource-sharing (e.g., cloud computing) and hardware malfunctions. Synchronous algorithms like AllReduce-SGD and D-PSGD perform poorly when workers' computation and/or communication speeds vary. Centralized asynchronous algorithms, such as A-PSGD, do poorly when the parameter server's network links slow down. In contrast, AD-PSGD localizes the impact of slower workers or network links. We evaluate AD-PSGD's robustness by randomly slowing down 1 of the 16 workers and its incoming/outgoing network links. Due to space limit, we show results for ResNet-20 model as the VGG results are similar.

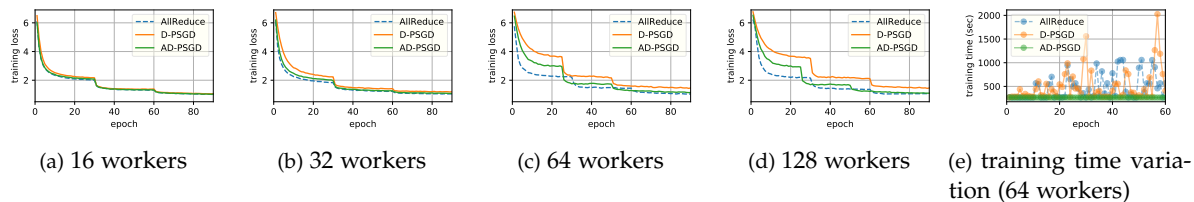


Figure 4: Training loss and training time per epoch comparison for ResNet-50 model on ImageNet dataset, evaluated up to 128 workers. AD-PSGD and AllReduce-SGD converge alike, better than D-PSGD. For 64 workers AD-PSGD finishes each epoch in 264 seconds, whereas AllReduce-SGD and D-PSGD can take over 1000 sec/epoch.

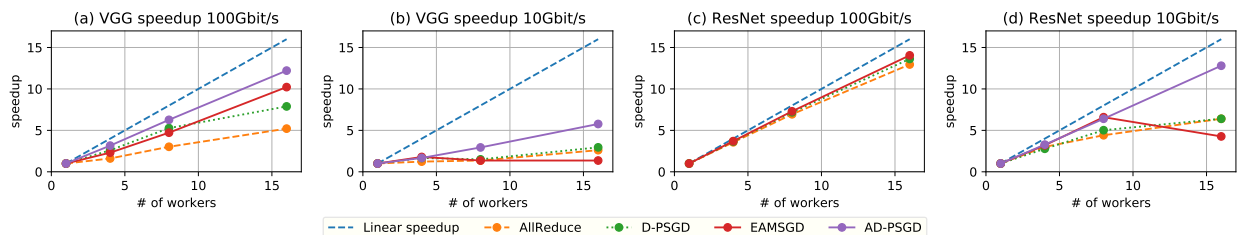


Figure 5: Speedup comparison for VGG (communication intensive) and ResNet-20 (computation intensive) models on CIFAR10. Experiments run on IBM HPC w/ 100Gbit/s network links and on x86 system w/ 10Gbit/s network links. AD-PSGD consistently achieves the best speedup.

Robustness against slow computation Figure 6 and Table 3 shows that AD-PSGD’s convergence is robust against slower workers. AD-PSGD can converge faster than AllReduce-SGD and D-PSGD by orders of magnitude when there is a very slow worker.

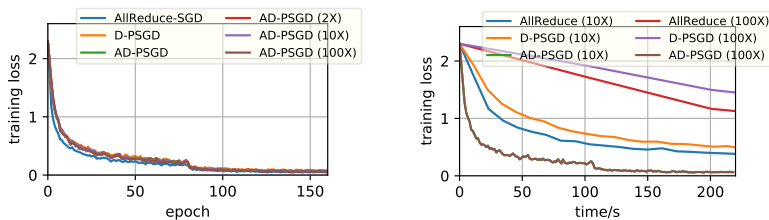
Robustness against slow communication Figure 7 shows that AD-PSGD is robust when one worker is connected to slower network links. In contrast, centralized asynchronous algorithm EAMSGD uses a larger communication period to overcome slower links, which significantly slows down the convergence.

These results show only AD-PSGD is robust against both heterogeneous computation and heterogeneous communication.

³In this paper, we use ASGD and EAMSGD interchangeably.

Table 3: Runtime efficiency comparison for ResNet-20 model on CIFAR-10 dataset when a computation device slows down by 2X-100X. AD-PSGD converges faster than AllReduce-SGD and D-PSGD, by orders of magnitude. 16 workers in total.

Slowdown of one node	AD-PSGD		AllReduce/D-PSGD	
	Time/epoch(sec)	Speedup	Time/epoch (sec)	Speedup
no slowdown	1.22	14.78	1.47/1.45	12.27/12.44
2X	1.28	14.09	2.6/2.36	6.93/7.64
10X	1.33	13.56	11.51/11.24	1.56/1.60
100X	1.33	13.56	100.4/100.4	0.18/0.18



(a) AD-PSGD converges steadily w.r.t. epochs despite a slower computing device. (b) AD-PSGD converges much faster than AllReduce-SGD and D-PSGD in the presence of a slower computing device.

Figure 6: Training loss for ResNet-20 model on CIFAR10 w.r.t (a) epochs and (b) runtime, when a computation device slows down by 2X-100X. 16 workers in total.

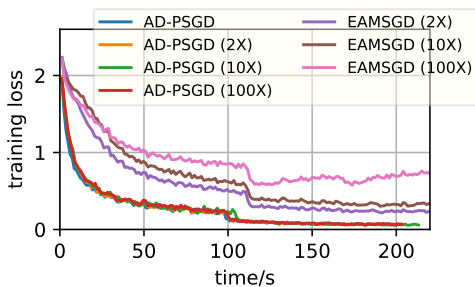


Figure 7: Training loss for ResNet-20 on CIFAR10 w.r.t. runtime, when a network link slows down by 2X-100X.

5.5 ImageNet experiments

We further evaluate the AD-PSGD’s convergence rate w.r.t. epochs using ImageNet-1K and ResNet-50 model. We compare AD-PSGD with AllReduce-SGD and D-PSGD as they tend to converge better than A-PSGD.

Figure 4 and Table 4 demonstrate that w.r.t. epochs AD-PSGD converges similarly to AllReduce and converges better than D-PSGD when running with 16,32,64,128 workers. How to maintain convergence while increasing $M \times n^4$ is an active ongoing research area Goyal et al. [2017], Zhang et al. [2016] and it is orthogonal to the topic of this paper. For 64 and 128 workers, we adopted similar learning rate tuning scheme as proposed in Goyal et al. [2017] (i.e., learning rate warm-up and linear scaling)⁵ *It worths noting*

⁴ M is mini-batch size per worker and n is the number of workers

⁵In AD-PSGD, we decay the learning rate every 25 epochs instead of 30 epochs as in AllReduce.

Table 4: Testing accuracy comparison for ResNet-50 model on ImageNet dataset for AllReduce, D-PSGD, and AD-PSGD. The ResNet-50 model is trained for 90 epochs. AD-PSGD and AllReduce-SGD achieve similar model accuracy.

	AllReduce	D-PSGD	AD-PSGD
16 Workers	74.86%	74.74%	75.28%
32 Workers	74.78%	73.66%	74.66%
64 Workers	74.90%	71.18%	74.20%
128 Workers	74.78%	70.90%	74.23%

that we could further increase the scalability of AD-PSGD by combining learners on the same computing node as a super-learner (via Nvidia NCCL AllReduce collectives). In this way, a 128-worker system can easily scale up to 512 GPUs or more, depending on the GPU count on a node.

It is common for a long running job to have exclusive access to computing devices but share the network links with other jobs in a cluster. Figure 4e shows the epoch-wise training time of the AD-PSGD, D-PSGD and AllReduce run over 64 GPUs (16 nodes) over a reserved window of 10 hours when the job shares network links with other jobs on IBM HPC. AD-PSGD finishes each epoch in 264 seconds, whereas AllReduce-SGD and D-PSGD can take over 1000 sec/epoch.

Above results show AD-PSGD converges similarly to AllReduce-SGD w.r.t epochs and better than D-PSGD. Techniques used for tuning learning rate for AllReduce-SGD can be applied to AD-PSGD when batch size is large. AD-PSGD is robust in a resource-sharing system.

6 Conclusion

This paper proposes an asynchronous decentralized stochastic gradient descent algorithm (AD-PSGD). The algorithm is not only robust in heterogeneous environments by combining both decentralization and asynchronization, but it is also theoretically justified to have the same convergence rate as its synchronous and/or centralized counterparts and can achieve linear speedup w.r.t. number of workers. Extensive experiments validate the proposed algorithm.

References

- M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- A. Agarwal and J. C. Duchi. Distributed delayed stochastic optimization. In *NIPS*, 2011.
- N. S. Aybat, Z. Wang, T. Lin, and S. Ma. Distributed linearized alternating direction method of multipliers for composite convex consensus optimization. *arXiv preprint arXiv:1512.08122*, 2015.
- T. C. Aysal, M. E. Yildiz, A. D. Sarwate, and A. Scaglione. Broadcast gossip algorithms for consensus. *IEEE Transactions on Signal processing*, 2009.
- P. Bianchi, G. Fort, and W. Hachem. Performance of a distributed stochastic approximation algorithm. *IEEE Transactions on Information Theory*, 2013.
- S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah. Gossip algorithms: Design, analysis and applications. In *INFOCOM*, 2005.
- R. Carli, F. Fagnani, P. Frasca, and S. Zampieri. Gossip consensus algorithms via quantized communication. *Automatica*, 2010.
- T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao. Optimal distributed online prediction using mini-batches. *Journal of Machine Learning Research*, 2012.
- F. Fagnani and S. Zampieri. Randomized consensus algorithms over large scale networks. *IEEE Journal on Selected Areas in Communications*, 2008.
- FAIR. ResNet in Torch. <https://github.com/facebook/fb.resnet.torch>, 2017.
- H. R. Feyzmahdavian, A. Aytekin, and M. Johansson. An asynchronous mini-batch algorithm for regularized stochastic optimization. *IEEE Transactions on Automatic Control*, 2016.
- S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 2013.
- S. Ghadimi, G. Lan, and H. Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 2016.
- P. Goyal, P. Dollár, R. B. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017. URL <http://arxiv.org/abs/1706.02677>.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- G. Lan, S. Lee, and Y. Zhou. Communication-efficient algorithms for decentralized and stochastic optimization. *arXiv preprint arXiv:1701.03961*, 2017.
- Z. Li, W. Shi, and M. Yan. A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates. *arXiv preprint arXiv:1704.07807*, 2017.

- X. Lian, Y. Huang, Y. Li, and J. Liu. Asynchronous parallel stochastic gradient for nonconvex optimization. In *NIPS*, 2015.
- X. Lian, H. Zhang, C.-J. Hsieh, Y. Huang, and J. Liu. A comprehensive linear speedup analysis for asynchronous stochastic parallel optimization from zeroth-order to first-order. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *NIPS*. Curran Associates, Inc., 2016.
- X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu. Can decentralized algorithms outperform centralized algorithms? A case study for decentralized parallel stochastic gradient descent, 2017.
- J. Lu, C. Y. Tang, P. R. Regier, and T. D. Bow. A gossip algorithm for convex consensus optimization over networks. In *ACC*. IEEE, 2010.
- N. Luehr. Fast multi-gpu collectives with nccl, 2016. URL <https://devblogs.nvidia.com/parallelforall/fast-multi-gpu-collectives-nccl/>.
- A. Mokhtari and A. Ribeiro. DSA: decentralized double stochastic averaging gradient algorithm. *Journal of Machine Learning Research*, 2016.
- E. Moulines and F. R. Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *NIPS*, 2011.
- MPI contributors. MPI AllReduce, 2015. URL <http://mpi-forum.org/docs/>.
- A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 2009.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 2009.
- R. Olfati-Saber, J. A. Fax, and R. M. Murray. Consensus and cooperation in networked multi-agent systems. *Proceedings of the IEEE*, 2007.
- T. Paine, H. Jin, J. Yang, Z. Lin, and T. Huang. Gpu asynchronous stochastic gradient descent to speed up neural network training. *arXiv preprint arXiv:1312.6186*, 2013.
- P. Patarasuk and X. Yuan. Bandwidth optimal all-reduce algorithms for clusters of workstations. *Journal of Parallel and Distributed Computing*, 2009.
- S. S. Ram, A. Nedic, and V. V. Veeravalli. Distributed subgradient projection algorithm for convex optimization. In *ICASSP*. IEEE, 2009.
- S. S. Ram, A. Nedić, and V. V. Veeravalli. Asynchronous gossip algorithm for stochastic optimization: Constant stepsize analysis. In *Recent Advances in Optimization and its Applications in Engineering*. Springer, 2010.
- B. Recht, C. Re, S. Wright, and F. Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *Advances in neural information processing systems*, 2011.
- L. Schenato and G. Gamba. A distributed consensus protocol for clock synchronization in wireless sensor network. In *CDC*. IEEE, 2007.
- F. Seide and A. Agarwal. CNTK: Microsoft’s open-source deep-learning toolkit. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*. ACM, 2016.

- W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin. On the linear convergence of the admm in decentralized consensus optimization.
- W. Shi, Q. Ling, G. Wu, and W. Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 2015a.
- W. Shi, Q. Ling, G. Wu, and W. Yin. A proximal gradient algorithm for decentralized composite optimization. *IEEE Transactions on Signal Processing*, 63(22):6013–6023, 2015b.
- K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015.
- B. Sirb and X. Ye. Consensus optimization with delayed and stochastic gradients on decentralized networks. In *Big Data*, 2016.
- K. Srivastava and A. Nedic. Distributed asynchronous constrained stochastic optimization. *IEEE Journal of Selected Topics in Signal Processing*, 2011.
- S. Sundhar Ram, A. Nedić, and V. Veeravalli. Distributed stochastic subgradient projection algorithms for convex optimization. *Journal of optimization theory and applications*, 2010.
- Z. Wang, Z. Yu, Q. Ling, D. Berberidis, and G. B. Giannakis. Decentralized rls with data-adaptive censoring for regressions over large-scale networks. *arXiv preprint arXiv:1612.08263*, 2016.
- T. Wu, K. Yuan, Q. Ling, W. Yin, and A. H. Sayed. Decentralized consensus optimization with asynchrony and delays. *arXiv preprint arXiv:1612.00150*, 2016.
- K. Yuan, Q. Ling, and W. Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*, 2016.
- S. Zagoruyko. CIFAR VGG in Torch. <https://github.com/szagoruyko/cifar.torch>, 2015.
- R. Zhang and J. Kwok. Asynchronous distributed admm for consensus optimization. In *ICML*, 2014.
- S. Zhang, A. E. Choromanska, and Y. LeCun. Deep learning with elastic averaging SGD. In *Advances in Neural Information Processing Systems*, pages 685–693, 2015.
- W. Zhang, S. Gupta, and F. Wang. Model accuracy and runtime tradeoff in distributed deep learning: A systematic study. In *IEEE International Conference on Data Mining*, 2016.
- W. Zhang, M. Feng, Y. Zheng, Y. Ren, Y. Wang, J. Liu, P. Liu, B. Xiang, L. Zhang, B. Zhou, and F. Wang. Gadei: On scale-up training as a service for deep learning. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*. The IEEE International Conference on Data Mining series(ICDM'2017), 2017.

A Wait-free (continuous) training and communication

The theoretical guarantee of AD-PSGD relies on the doubly stochastic property of matrix W . The implication is the averaging of the weights between two workers should be atomic. This brings a special challenge for current distributed deep learning frameworks where the computation (gradients calculation and weights update) runs on GPU devices and the communication runs on CPU (or its peripherals such as infiniband or RDMA), because when there is averaging happening on a worker, the GPU is not allowed to update gradients into the weights. This can be solve by using CPU to update weights while GPUs only calculate gradients. Every worker (including active and passive workers) runs two threads in parallel with a shared buffer g , one thread for computation and the other for communication. Algorithm 2, Algorithm 3, and Algorithm 4 illustrate the task on each thread. The communication thread is run by CPUs, while the computation thread is run by GPUs. In this way GPUs can continuously calculate new gradients by putting the results in CPUs' buffer regardless of whether there is averaging happening. Recall in D-PSGD, communication only occurs once in each iteration. In contrast, AD-PSGD can exchange weights at any time by using this implementation.

Algorithm 2 Computation thread on active or passive worker (worker index is i)

Require: Batch size M

- 1: **while** not terminated **do**
- 2: Pull model x^i from the communication thread.
- 3: Update locally in the thread $x^i \leftarrow x^i - \gamma g$.^a
- 4: Randomly sample a batch $\zeta^i := (\zeta_1^i, \zeta_2^i, \dots, \zeta_M^i)$ from local data of the i -th worker and compute the stochastic gradient $g^i(x^i; \zeta^i) := \sum_{m=1}^M \nabla F(x^i; \zeta_m^i)$ locally.
- 5: **wait until** $g = 0$ **then**
- 6: Local buffer $g \leftarrow g^i(x^i; \zeta^i)$.^b
- 7: **end wait until**
- 8: **end while**

^aAt this time the communication thread may have not update g into x^i so the computation thread pulls an old model. We compensate this by doing local update in computation thread. We observe this helps the scaling.

^bWe can also make a queue of gradients here to avoid the waiting. Note that doing this will make the effective batch-size different from M .

Algorithm 3 Communication thread on active worker (worker index is i)

Require: Initialize local model x^i , learning rate γ .

- 1: **while** not terminated **do**
 - 2: **if** $g \neq 0$ **then**
 - 3: $x^i \leftarrow x^i - \gamma g, \quad g \leftarrow 0$.
 - 4: **end if**
 - 5: Randomly select a neighbor (namely worker j). Send x^i to worker j and fetch x^j from it.
 - 6: $x^i \leftarrow \frac{1}{2}(x^i + x^j)$.
 - 7: **end while**
-

B NLC experiments

In this section, we use IBM proprietary natural language processing dataset and model to evaluate AD-PSGD against other algorithms.

Algorithm 4 Communication thread on passive worker (worker index is j)

Require: Initialize local model x^j , learning rate γ .

- 1: **while** not terminated **do**
 - 2: **if** $g \neq 0$ **then**
 - 3: $x^j \leftarrow x^j - \gamma g, \quad g \leftarrow 0.$
 - 4: **end if**
 - 5: **if** receive the request of reading local model (say from worker i) **then**
 - 6: Send x^j to worker i .
 - 7: $x^j \leftarrow \frac{1}{2}(x^i + x^j).$
 - 8: **end if**
 - 9: **end while**
-

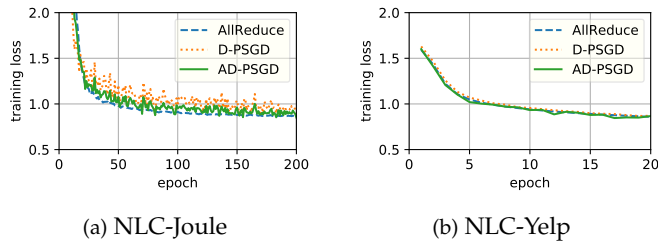


Figure 8: Training loss comparison for IBM NLC model on Joule and Yelp datasets. AllReduce-SGD, D-PSGD and AD-PSGD converge alike w.r.t. epochs.

The IBM NLC task is to classify input sentences into a target category in a predefined label set. The NLC model is a CNN model that has a word-embedding lookup table layer, a convolutional layer and a fully connected layer with a softmax output layer. We use two datasets in our evaluation. The first dataset Joule is an in-house customer dataset that has 2.5K training samples, 1K test samples, and 311 different classes. The second dataset Yelp, which is a public dataset, has 500K training samples, 2K test samples and 5 different classes. Figure 8 shows that AD-PSGD converges (w.r.t epochs) similarly to AllReduce-SGD and D-PSGD on NLC tasks.

Above results show AD-PSGD converges similarly (w.r.t) to AllReduce-SGD and D-PSGD for IBM NLC workload, which is an example of proprietary workloads.

C Appendix: proofs

In the following analysis we define

$$M_k := \sum_{i=1}^n p_i \left\| \frac{X_k \mathbf{1}_n}{n} - X_k e_i \right\|^2, \quad (9)$$

and

$$\hat{M}_k := M_{k-\tau_k}. \quad (10)$$

We also define

$$\begin{aligned} \partial f(X_k) &:= n \begin{bmatrix} p_1 \nabla f_1(x_k^1) & p_2 \nabla f_2(x_k^2) & \cdots & p_n \nabla f_n(x_k^n) \end{bmatrix} \in \mathbb{R}^{N \times n}, \\ \partial f(X_k, i) &:= \begin{bmatrix} 0 & \cdots & \nabla f_i(x_k^i) & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{N \times n}, \\ \partial g(\hat{X}_k, \xi_k) &:= n \begin{bmatrix} p_1 \sum_{j=1}^M \nabla F(\hat{x}_k^1, \xi_{k,j}^1) & \cdots & p_n \sum_{j=1}^M \nabla F(\hat{x}_k^n, \xi_{k,j}^n) \end{bmatrix} \in \mathbb{R}^{N \times n}. \end{aligned}$$

$$\begin{aligned} \bar{\rho} &:= \frac{n-1}{n} \left(\frac{1}{1-\rho} + \frac{2\sqrt{\rho}}{(1-\sqrt{\rho})^2} \right), \\ C_1 &:= 1 - 24M^2 L^2 \gamma^2 \left(T \frac{n-1}{n} + \bar{\rho} \right), \\ C_2 &:= \frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} - \frac{2M^3 L^2 T^2 \gamma^3}{n^3} - \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 + \frac{12M^3 L^4 T^2 \gamma^3}{n^3} \right) \frac{4M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho})}{C_1}, \\ C_3 &:= \frac{1}{2} + \frac{2 \left(6\gamma^2 L^2 M^2 + \gamma n M L + \frac{12M^3 L^3 T^2 \gamma^3}{n} \right) \bar{\rho}}{C_1} + \frac{L T^2 \gamma M}{n}. \end{aligned}$$

Proof to Theorem 1. We start from

$$\begin{aligned} & \mathbb{E} f \left(\frac{X_{k+1} \mathbf{1}_n}{n} \right) \\ &= \mathbb{E} f \left(\frac{X_k W_k \mathbf{1}_n}{n} - \gamma \frac{\partial g(\hat{X}_k; \xi_k^{i_k}, i_k) \mathbf{1}_n}{n} \right) = \mathbb{E} f \left(\frac{X_k \mathbf{1}_n}{n} - \gamma \frac{\partial g(\hat{X}_k; \xi_k^{i_k}, i_k) \mathbf{1}_n}{n} \right) \\ &\leq \mathbb{E} f \left(\frac{X_k \mathbf{1}_n}{n} \right) - \gamma \mathbb{E} \left\langle \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right), \frac{\partial g(\hat{X}_k; \xi_k^{i_k}, i_k) \mathbf{1}_n}{n} \right\rangle + \frac{\gamma^2 L}{2} \mathbb{E} \left\| \frac{\partial g(\hat{X}_k; \xi_k^{i_k}, i_k) \mathbf{1}_n}{n} \right\|^2 \\ &\stackrel{(3), \text{Lemma 4}}{\leq} \mathbb{E} f \left(\frac{X_k \mathbf{1}_n}{n} \right) - \frac{\gamma M}{n} \mathbb{E} \left\langle \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right), \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\rangle + \frac{\gamma^2 L \sigma^2 M}{2n^2} + \frac{\gamma^2 L M^2}{2n^2} \sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2 \\ &= \mathbb{E} f \left(\frac{X_k \mathbf{1}_n}{n} \right) + \frac{\gamma M}{2n} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) - \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2 - \frac{\gamma M}{2n} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\ &\quad + \frac{\gamma^2 L M^2}{2n^2} \sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2 + \frac{\gamma^2 L \sigma^2 M}{2n^2}. \end{aligned}$$

Using the upper bound of $\sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2$ in Lemma 5:

$$\mathbb{E} f \left(\frac{X_{k+1} \mathbf{1}_n}{n} \right)$$

$$\begin{aligned}
&\leq \mathbb{E}f\left(\frac{X_k \mathbf{1}_n}{n}\right) + \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 - \frac{\gamma M}{2n} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
&\quad + \frac{\gamma^2 L M^2}{2n^2} \left(12L^2 \hat{M}_k + 6\zeta^2 + 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2 \right) + \frac{\gamma^2 L \sigma^2 M}{2n^2} \\
&= \mathbb{E}f\left(\frac{X_k \mathbf{1}_n}{n}\right) + \frac{\gamma M}{2n} \underbrace{\mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2}_{T_1} - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 \\
&\quad - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} + \frac{6\gamma^2 L^3 M^2}{n^2} \hat{M}_k. \tag{11}
\end{aligned}$$

For T_1 we have

$$\begin{aligned}
T_1 &= \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
&\leq 2\mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \nabla f\left(\frac{\hat{X}_k \mathbf{1}_n}{n}\right) \right\|^2 + 2\mathbb{E} \left\| \nabla f\left(\frac{\hat{X}_k \mathbf{1}_n}{n}\right) - \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
&= 2\mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \nabla f\left(\frac{\hat{X}_k \mathbf{1}_n}{n}\right) \right\|^2 + 2\mathbb{E} \left\| \sum_i p_i \left(\nabla f_i\left(\frac{\hat{X}_k \mathbf{1}_n}{n}\right) - \nabla f_i(\hat{x}_k^i) \right) \right\|^2 \\
&\leq 2\mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \nabla f\left(\frac{\hat{X}_k \mathbf{1}_n}{n}\right) \right\|^2 + 2\mathbb{E} \sum_i p_i \left\| \nabla f_i\left(\frac{\hat{X}_k \mathbf{1}_n}{n}\right) - \nabla f_i(\hat{x}_k^i) \right\|^2 \\
&\stackrel{\text{Assumption 1:1}}{\leq} 2L^2 \mathbb{E} \left\| \frac{(X_k - \hat{X}_k) \mathbf{1}_n}{n} \right\|^2 + 2L^2 \mathbb{E} \hat{M}_k. \tag{12}
\end{aligned}$$

From (11) and (12) we obtain

$$\begin{aligned}
\mathbb{E}f\left(\frac{X_{k+1} \mathbf{1}_n}{n}\right) &\leq \mathbb{E}f\left(\frac{X_k \mathbf{1}_n}{n}\right) + \frac{\gamma M}{2n} \mathbb{E} \left(2L^2 \left\| \frac{(X_k - \hat{X}_k) \mathbf{1}_n}{n} \right\|^2 + 2L^2 \hat{M}_k \right) \\
&\quad - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
&\quad + \frac{6\gamma^2 L^3 M^2}{n^2} \mathbb{E} \hat{M}_k + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} \\
&= \mathbb{E}f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
&\quad + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \mathbb{E} \hat{M}_k + \frac{\gamma M}{n} L^2 \mathbb{E} \left\| \frac{(X_k - \hat{X}_k) \mathbf{1}_n}{n} \right\|^2 + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} \\
&\stackrel{\text{Lemma 8}}{\leq} \mathbb{E}f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
&\quad + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \mathbb{E} \hat{M}_k + \frac{\gamma M}{n} L^2 \left(\frac{\tau_k^2 \gamma^2 \sigma^2 M}{n^2} + \tau_k \gamma^2 \sum_{t=1}^{\tau_k} \left(\frac{M^2}{n^2} \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\hat{x}_{k-t}^i) \right\|^2 \right) \right) \\
&\quad + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} \\
&\leq \mathbb{E}f\left(\frac{X_k \mathbf{1}_n}{n}\right) - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2
\end{aligned}$$

$$\begin{aligned}
& + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \mathbb{E} \hat{M}_k + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} + \frac{L^2 T^2 \gamma^3 \sigma^2 M^2}{n^3} \\
& + \frac{M^3 L^2 \tau_k \gamma^3}{n^3} \sum_{t=1}^{\tau_k} \left(\sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_{k-t}^i)\|^2 \right) \\
\stackrel{\text{Lemma 5}}{\leq} & \mathbb{E} f \left(\frac{X_k \mathbf{1}_n}{n} \right) - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \mathbb{E} \hat{M}_k + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} + \frac{L^2 T^2 \gamma^3 \sigma^2 M^2}{n^3} \\
& + \frac{M^3 L^2 \tau_k \gamma^3}{n^3} \sum_{t=1}^{\tau_k} \left(12L^2 \hat{M}_{k-t} + 6\zeta^2 + 2\mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_{k-t}^j) \right\|^2 \right) \\
= & \mathbb{E} f \left(\frac{X_k \mathbf{1}_n}{n} \right) - \frac{\gamma M}{2n} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \mathbb{E} \hat{M}_k + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2)}{2n^2} + \frac{L^2 T^2 \gamma^3 M (\sigma^2 M + 6\zeta^2 M^2)}{n^3} \\
& + \frac{2M^3 L^2 T \gamma^3}{n^3} \sum_{t=1}^{\tau_k} \left(6L^2 \mathbb{E} \hat{M}_{k-t} + \mathbb{E} \left\| \frac{\partial f(\hat{X}_{k-t}) \mathbf{1}_n}{n} \right\|^2 \right).
\end{aligned}$$

Summing from $k = 0$ to $k = K - 1$ we obtain

$$\begin{aligned}
\mathbb{E} f \left(\frac{X_K \mathbf{1}_n}{n} \right) & \leq \mathbb{E} f \left(\frac{X_0 \mathbf{1}_n}{n} \right) - \frac{\gamma M}{2n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \sum_{k=0}^{K-1} \mathbb{E} \hat{M}_k + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2) K}{2n^2} + \frac{L^2 T^2 \gamma^3 M (\sigma^2 M + 6\zeta^2 M^2) K}{n^3} \\
& + \frac{2M^3 L^2 T \gamma^3}{n^3} \sum_{k=0}^{K-1} \sum_{t=1}^{\tau_k} \left(6L^2 \mathbb{E} \hat{M}_{k-t} + \mathbb{E} \left\| \frac{\partial f(\hat{X}_{k-t}) \mathbf{1}_n}{n} \right\|^2 \right) \\
\leq & \mathbb{E} f \left(\frac{X_0 \mathbf{1}_n}{n} \right) - \frac{\gamma M}{2n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2 - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} \right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 \right) \sum_{k=0}^{K-1} \mathbb{E} \hat{M}_k + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2) K}{2n^2} + \frac{L^2 T^2 \gamma^3 M (\sigma^2 M + 6\zeta^2 M^2) K}{n^3} \\
& + \frac{2M^3 L^2 T^2 \gamma^3}{n^3} \sum_{k=0}^{K-1} \left(6L^2 \mathbb{E} \hat{M}_k + \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \right) \\
= & \mathbb{E} f \left(\frac{X_0 \mathbf{1}_n}{n} \right) - \frac{\gamma M}{2n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2 \\
& - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} - \frac{2M^3 L^2 T^2 \gamma^3}{n^3} \right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 + \frac{12M^3 L^4 T^2 \gamma^3}{n^3} \right) \sum_{k=0}^{K-1} \mathbb{E} \hat{M}_k \\
& + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2) K}{2n^2} + \frac{L^2 T^2 \gamma^3 M (\sigma^2 M + 6\zeta^2 M^2) K}{n^3}.
\end{aligned}$$

$$\begin{aligned}
& \leq \mathbb{E} f\left(\frac{X_0 \mathbf{1}_n}{n}\right) - \frac{\gamma M}{2n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 \\
& \quad - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} - \frac{2M^3 L^2 T^2 \gamma^3}{n^3} \right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& \quad + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 + \frac{12M^3 L^4 T^2 \gamma^3}{n^3} \right) \sum_{k=0}^{K-1} \mathbb{E} \hat{M}_k \\
& \quad + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2) K}{2n^2} + \frac{L^2 T^2 \gamma^3 M (\sigma^2 M + 6\zeta^2 M^2) K}{n^3} \\
& \stackrel{C_1 > 0, \text{Lemma 7}}{\leq} \mathbb{E} f\left(\frac{X_0 \mathbf{1}_n}{n}\right) - \frac{\gamma M}{2n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 \\
& \quad - \left(\frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} - \frac{2M^3 L^2 T^2 \gamma^3}{n^3} \right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 \\
& \quad + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 + \frac{12M^3 L^4 T^2 \gamma^3}{n^3} \right) K \frac{2\gamma^2 (M\sigma^2 + 6M^2 \zeta^2) \bar{\rho}}{C_1} \\
& \quad + \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 + \frac{12M^3 L^4 T^2 \gamma^3}{n^3} \right) \frac{4M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho}) \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_k^i) \right\|^2}{C_1} \\
& \quad + \frac{\gamma^2 L (\sigma^2 M + 6\zeta^2 M^2) K}{2n^2} + \frac{L^2 T^2 \gamma^3 M (\sigma^2 M + 6\zeta^2 M^2) K}{n^3} \\
& = \mathbb{E} f\left(\frac{X_0 \mathbf{1}_n}{n}\right) - \frac{\gamma M}{2n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2 \\
& \quad - C_2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{\partial f(\hat{X}_k) \mathbf{1}_n}{n} \right\|^2 + C_3 \frac{\gamma^2 L K}{n^2} (M\sigma^2 + 6M^2 \zeta^2).
\end{aligned}$$

Thus while $C_3 \leq 1$ and $C_2 \geq 0$ we have

$$\begin{aligned}
\frac{\sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f\left(\frac{X_k \mathbf{1}_n}{n}\right) \right\|^2}{K} & \leq \frac{2 \left(\mathbb{E} f\left(\frac{X_0 \mathbf{1}_n}{n}\right) - \mathbb{E} f\left(\frac{X_K \mathbf{1}_n}{n}\right) \right)}{\gamma K M / n} + \frac{2\gamma L}{Mn} (M\sigma^2 + 6M^2 \zeta^2) \\
& \leq \frac{2(\mathbb{E} f(x_0) - \mathbb{E} f^*)}{\gamma K M / n} + \frac{2\gamma L}{n} (\sigma^2 + 6M\zeta^2).
\end{aligned}$$

It completes the proof. \square

Lemma 3. Define $\prod_{k=1}^0 W_k = I$, where I is the identity matrix. Then

$$\mathbb{E} \left\| \frac{\mathbf{1}_n}{n} - \prod_{k=1}^K W_k e_i \right\|^2 \leq \frac{n-1}{n} \rho^K, \quad \forall K \geq 0.$$

Proof. Let $y_K = \frac{\mathbf{1}_n}{n} - \prod_{k=1}^K W_k e_i$. Then noting that $y_{K+1} = W_{K+1} y_K$ we have

$$\begin{aligned}
& \mathbb{E} \|y_{K+1}\|^2 \\
& = \mathbb{E} \|W_{K+1} y_K\|^2 \\
& = \mathbb{E} \langle W_{K+1} y_K, W_{K+1} y_K \rangle \\
& = \mathbb{E} \langle y_K, W_{K+1}^\top W_{K+1} y_K \rangle
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \langle y_K, \mathbb{E} \mathbb{E}_{i_{K+1}} (W_{K+1}^\top W_{K+1}) y_K \rangle \\
&= \mathbb{E} \langle y_K, \mathbb{E} (W_{K+1}^\top W_{K+1}) y_K \rangle.
\end{aligned}$$

Note that $\mathbb{E}(W_{K+1}^\top W_{K+1})$ is symmetric and doubly stochastic and $\mathbf{1}_n$ is an eigenvector of $\mathbb{E}(W_{K+1}^\top W_{K+1})$ with eigenvalue 1. Starting from $\mathbf{1}_n$ we construct a basis of \mathbb{R}^n composed by the eigenvectors of $\mathbb{E}(W_{K+1}^\top W_{K+1})$, which is guaranteed to exist by the spectral theorem of Hermitian matrices. From (2) the magnitude of all other eigenvectors' associated eigenvalues should be smaller or equal to ρ . Noting y_K is orthogonal to $\mathbf{1}_n$, we decompose y_K using this constructed basis and it follows that

$$\mathbb{E} \|y_{K+1}\|^2 \leq \rho \mathbb{E} \|y_K\|^2.$$

Noting that $\|y_0\|^2 = \|\mathbf{1}_n/n - e_i\|^2 = \frac{(n-1)^2}{n^2} + \sum_{i=1}^{n-1} \frac{1}{n^2} = \frac{n^2 - 2n + 1 + n - 1}{n^2} = \frac{n-1}{n}$, by induction, we complete the proof. \square

Lemma 4.

$$\mathbb{E} \left\| \frac{\partial g(\hat{X}_k; \zeta_k^{i_k}, i_k) \mathbf{1}_n}{n} \right\|^2 \leq \frac{\sigma^2 M}{n^2} + \frac{M^2}{n^2} \sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2, \quad \forall k \geq 0.$$

Proof. The LHS can be bounded by

$$\begin{aligned}
\mathbb{E} \left\| \frac{\partial g(\hat{X}_k; \zeta_k^{i_k}, i_k) \mathbf{1}_n}{n} \right\|^2 &\stackrel{(1)}{=} \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\sum_{j=1}^M \nabla F(\hat{x}_k^i, \zeta_{k,j}^i)}{n} \right\|^2 \\
&= \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\sum_{j=1}^M (\nabla F(\hat{x}_k^i, \zeta_{k,j}^i) - \nabla f_i(\hat{x}_k^i))}{n} \right\|^2 + \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{M \nabla f_i(\hat{x}_k^i)}{n} \right\|^2 \\
&\stackrel{(5)}{\leq} \frac{\sigma^2 M}{n^2} + \frac{M^2}{n^2} \sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2.
\end{aligned}$$

\square

Lemma 5.

$$\sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2 \leq 12L^2 \mathbb{E} \hat{M}_k + 6\zeta^2 + 2\mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2, \quad \forall k \geq 0.$$

Proof. The LHS can be bounded by

$$\begin{aligned}
\sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i)\|^2 &\leq \sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_k^i) - \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) + \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j)\|^2 \\
&\leq 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\hat{x}_k^i) - \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2 + 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2 \\
&= 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\hat{x}_k^i) - \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2 + 2\mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2. \tag{13}
\end{aligned}$$

For the first term on the RHS we have

$$\sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\hat{x}_k^i) - \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) \right\|^2$$

$$\begin{aligned}
&\leq 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\hat{x}_k^i) - \nabla f_i \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) - \sum_{j=1}^n p_j \nabla f_j \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 \\
&\quad + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) - \sum_{j=1}^n p_j \nabla f_j \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 \\
&\leq 3L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \hat{x}_k^i - \frac{\hat{X}_k \mathbf{1}_n}{n} \right\|^2 + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) - \sum_{j=1}^n p_j \nabla f_j \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 \\
&\quad + 3 \mathbb{E} \left\| \sum_{j=1}^n p_j \nabla f_j(\hat{x}_k^j) - \sum_{j=1}^n p_j \nabla f_j \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 \\
&\leq 3L^2 \mathbb{E} \hat{M}_k + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) - \nabla f \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 + 3 \sum_{j=1}^n p_j \mathbb{E} \left\| \nabla f_j(\hat{x}_k^j) - \nabla f_j \left(\frac{\hat{X}_k \mathbf{1}_n}{n} \right) \right\|^2 \\
&\leq 6L^2 \mathbb{E} \hat{M}_k + 3\zeta^2.
\end{aligned}$$

Plugging this upper bound into (13) we complete the proof. \square

Lemma 6. For any $k \geq -1$ we have

$$\begin{aligned}
&\mathbb{E} \left\| \frac{X_{k+1} \mathbf{1}_n}{n} - X_{k+1} e_i \right\|^2 \\
&\leq 2\gamma^2 (M\sigma^2 + 6M^2\zeta^2) \bar{\rho} \\
&\quad + 2 \frac{n-1}{n} M^2 \gamma^2 \mathbb{E} \sum_{j=0}^k \left(12L^2 \hat{M}_j + 2 \mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \left(\rho^{k-j} + 2(k-j) \rho^{\frac{k-j}{2}} \right).
\end{aligned}$$

Proof. Note that for $k = -1$, we have

$$\mathbb{E} \left\| \frac{X_{k+1} \mathbf{1}_n}{n} - X_{k+1} e_i \right\|^2 = 0.$$

Also note that the columns of X_0 are the same (all workers start with the same model), we have $X_0 W_k = X_0$ for all k and $X_0 \mathbf{1}_n / n - X_0 e_i = 0, \forall i$. It follows that

$$\begin{aligned}
&\mathbb{E} \left\| \frac{X_{k+1} \mathbf{1}_n}{n} - X_{k+1} e_i \right\|^2 \\
&= \mathbb{E} \left\| \frac{X_k \mathbf{1}_n - \gamma \partial g(\hat{X}_k; \zeta_k^{i_k}, i_k) \mathbf{1}_n}{n} - (X_k W_k e_i - \gamma \partial g(\hat{X}_k, \zeta_k^{i_k}, i_k) e_i) \right\|^2 \\
&= \mathbb{E} \left\| \frac{X_0 \mathbf{1}_n - \sum_{j=0}^k \gamma \partial g(\hat{X}_j; \zeta_j^{i_j}, i_j) \mathbf{1}_n}{n} - \left(X_0 \prod_{j=0}^k W_j e_i - \sum_{j=0}^k \gamma \partial g(\hat{X}_j; \zeta_j^{i_j}, i_j) \prod_{q=j+1}^k W_q e_i \right) \right\|^2 \\
&= \mathbb{E} \left\| - \sum_{j=0}^k \gamma \partial g(\hat{X}_j; \zeta_j^{i_j}, i_j) \frac{\mathbf{1}_n}{n} + \sum_{j=0}^k \gamma \partial g(\hat{X}_j; \zeta_j^{i_j}, i_j) \prod_{q=j+1}^k W_q e_i \right\|^2 \\
&= \gamma^2 \mathbb{E} \left\| \sum_{j=0}^k \partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right) \right\|^2
\end{aligned}$$

$$\begin{aligned}
& \leq 2\gamma^2 \mathbb{E} \left\| \underbrace{\sum_{j=0}^k (\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right)}_{A_1} \right\|^2 \\
& + 2M^2\gamma^2 \mathbb{E} \left\| \underbrace{\sum_{j=0}^k \partial f(\hat{X}_j, i_j) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right)}_{A_2} \right\|^2.
\end{aligned} \tag{14}$$

For A_1 ,

$$\begin{aligned}
A_1 &= \mathbb{E} \left\| \sum_{j=0}^k (\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right) \right\|^2 \\
&= \sum_{j=0}^k \mathbb{E} \left\| \underbrace{(\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right)}_{A_3} \right\|^2 \\
&+ 2 \mathbb{E} \underbrace{\sum_{k \geq j > j' \geq 0} \left\langle (\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right), \right.}_{A_4} \\
&\quad \left. (\partial g(\hat{X}_{j'}, \zeta_{j'}^{i_{j'}}, i_{j'}) - M\partial f(\hat{X}_{j'}, i_{j'})) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right) \right\rangle}.
\end{aligned}$$

A_3 can be bounded by a constant:

$$\begin{aligned}
A_3 &= \sum_{j=0}^k \mathbb{E} \left\| (\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right) \right\|^2 \\
&\leq \sum_{j=0}^k \mathbb{E} \|\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)\|^2 \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\|^2 \\
&\stackrel{\text{Lemma 3}}{\leq} \frac{n-1}{n} \sum_{j=0}^k \mathbb{E} \|\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)\|^2 \rho^{k-j} \\
&\stackrel{\text{Assumption 1:5}}{\leq} \frac{n-1}{n} M\sigma^2 \sum_{j=0}^k \rho^{k-j} \leq \frac{n-1}{n} \frac{M\sigma^2}{1-\rho}.
\end{aligned}$$

A_4 can be bounded by another constant:

$$\begin{aligned}
A_4 &= \sum_{k \geq j > j' \geq 0} \mathbb{E} \left\langle (\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right), \right. \\
&\quad \left. (\partial g(\hat{X}_{j'}, \zeta_{j'}^{i_{j'}}, i_{j'}) - M\partial f(\hat{X}_{j'}, i_{j'})) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right) \right\rangle \\
&\leq \sum_{k \geq j > j' \geq 0} \mathbb{E} \|\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\| \times \\
&\quad \|\partial g(\hat{X}_{j'}, \zeta_{j'}^{i_{j'}}, i_{j'}) - M\partial f(\hat{X}_{j'}, i_{j'})\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\| \\
&\leq \mathbb{E} \sum_{k \geq j > j' \geq 0} \left(\frac{\left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\|^2 \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\|^2}{2\alpha_{j,j'}} \right. \\
&\quad \left. + \frac{\|\partial g(\hat{X}_j, \zeta_j^{i_j}, i_j) - M\partial f(\hat{X}_j, i_j)\|^2 \|\partial g(\hat{X}_{j'}, \zeta_{j'}^{i_{j'}}, i_{j'}) - M\partial f(\hat{X}_{j'}, i_{j'})\|^2}{2/\alpha_{j,j'}} \right), \forall \alpha_{j,j'} > 0
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(5)}{\leq} \mathbb{E} \sum_{k \geq j > j' \geq 0} \left(\frac{\left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\|^2 \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\|^2}{2\alpha_{j,j'}} + \frac{M^2 \sigma^4}{2/\alpha_{j,j'}} \right), \forall \alpha_{j,j'} > 0 \\
&\leq \mathbb{E} \sum_{k \geq j > j' \geq 0} \left(\frac{n-1}{n} \frac{\left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\|^2}{2\alpha_{j,j'}} + \frac{M^2 \sigma^4}{2/\alpha_{j,j'}} \right), \forall \alpha_{j,j'} > 0 \\
&\stackrel{\text{Lemma 3}}{\leq} \mathbb{E} \sum_{k \geq j > j' \geq 0} \left(\left(\frac{n-1}{n} \right)^2 \frac{\rho^{k-j'}}{2\alpha_{j,j'}} + \frac{M^2 \sigma^4}{2/\alpha_{j,j'}} \right), \forall \alpha_{j,j'} > 0.
\end{aligned}$$

We can choose $\alpha_{j,j'} > 0$ to make the term in the last step become $\frac{n-1}{n} \sum_{k \geq j > j' \geq 0} \rho^{\frac{k-j'}{2}} M \sigma^2$ (by applying inequality of arithmetic and geometric means). Thus

$$\begin{aligned}
A_4 &\leq \frac{n-1}{n} \sum_{k \geq j > j' \geq 0} \rho^{\frac{k-j'}{2}} M \sigma^2 \leq \frac{n-1}{n} M \sigma^2 \sum_{j'=0}^k \sum_{j=j'+1}^k \rho^{\frac{k-j'}{2}} \\
&= \frac{n-1}{n} M \sigma^2 \sum_{j'=0}^k (k-j') \rho^{\frac{k-j'}{2}} \leq \frac{n-1}{n} M \sigma^2 \frac{\sqrt{\rho}}{(1-\sqrt{\rho})^2}.
\end{aligned}$$

Putting A_3 and A_4 back into A_1 we obtain:

$$A_1 \leq \frac{n-1}{n} M \sigma^2 \left(\frac{1}{1-\rho} + \frac{2\sqrt{\rho}}{(1-\sqrt{\rho})^2} \right) = M \sigma^2 \bar{\rho}. \quad (15)$$

We then start bounding A_2 :

$$\begin{aligned}
A_2 &= \mathbb{E} \left\| \sum_{j=0}^k \partial f(\hat{X}_j, i_j) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right) \right\|^2 \\
&= \mathbb{E} \sum_{j=0}^k \left\| \partial f(\hat{X}_j, i_j) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right) \right\|^2 \\
&\quad + 2 \mathbb{E} \sum_{j=0}^k \sum_{j'=j+1}^k \left\langle \partial f(\hat{X}_j, i_j) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right), \partial f(\hat{X}_{j'}, i_{j'}) \left(\frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right) \right\rangle \\
&\stackrel{\text{Lemma 3(1)}}{\leq} \frac{n-1}{n} \mathbb{E} \sum_{j=0}^k \left(\sum_{i=1}^n p_i \|\nabla f_i(\hat{x}_j^i)\|^2 \right) \rho^{k-j} \\
&\quad + 2 \mathbb{E} \sum_{j=0}^k \sum_{j'=j+1}^k \|\partial f(\hat{X}_j, i_j)\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\| \|\partial f(\hat{X}_{j'}, i_{j'})\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\|. \quad (16)
\end{aligned}$$

For the second term:

$$\begin{aligned}
&\mathbb{E} \sum_{j=0}^k \sum_{j'=j+1}^k \|\partial f(\hat{X}_j, i_j)\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\| \|\partial f(\hat{X}_{j'}, i_{j'})\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\| \\
&\leq \mathbb{E} \sum_{j=0}^k \sum_{j'=j+1}^k \left(\frac{\|\partial f(\hat{X}_j, i_j)\|^2 \|\partial f(\hat{X}_{j'}, i_{j'})\|^2}{2\alpha_{j,j'}} + \frac{\left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\|^2 \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\|^2}{2/\alpha_{j,j'}} \right), \forall \alpha_{j,j'} > 0
\end{aligned}$$

$$\stackrel{\text{Lemma 3}}{\leq} \frac{1}{2} \mathbb{E} \sum_{j \neq j'}^k \left(\frac{\|\partial f(\hat{X}_j, i_j)\|^2 \|\partial f(\hat{X}_{j'}, i_{j'})\|^2}{2\alpha_{j,j'}} + \frac{\rho^{k-\min\{j,j'\}}}{2/\alpha_{j,j'}} \left(\frac{n-1}{n} \right)^2 \right), \quad \forall \alpha_{j,j'} > 0, \alpha_{j,j'} = \alpha_{j',j}.$$

By applying inequality of arithmetic and geometric means to the term in the last step we can choose $\alpha_{j,j'} > 0$ such that

$$\begin{aligned} & \mathbb{E} \sum_{j=0}^k \sum_{j'=j+1}^k \|\partial f(\hat{X}_j, i_j)\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j+1}^k W_q e_i \right\| \|\partial f(\hat{X}_{j'}, i_{j'})\| \left\| \frac{\mathbf{1}_n}{n} - \prod_{q=j'+1}^k W_q e_i \right\| \\ & \leq \frac{n-1}{2n} \mathbb{E} \sum_{j \neq j'}^k \left(\|\partial f(\hat{X}_j, i_j)\| \|\partial f(\hat{X}_{j'}, i_{j'})\| \rho^{\frac{k-\min\{j,j'\}}{2}} \right) \\ & \leq \frac{n-1}{2n} \mathbb{E} \sum_{j \neq j'}^k \left(\frac{\|\partial f(\hat{X}_j, i_j)\|^2 + \|\partial f(\hat{X}_{j'}, i_{j'})\|^2}{2} \rho^{\frac{k-\min\{j,j'\}}{2}} \right) \\ & = \frac{n-1}{2n} \mathbb{E} \sum_{j \neq j'}^k \left(\|\partial f(\hat{X}_j, i_j)\|^2 \rho^{\frac{k-\min\{j,j'\}}{2}} \right) = \frac{n-1}{n} \sum_{j=0}^k \sum_{j'=j+1}^k \left(\mathbb{E} \|\partial f(\hat{X}_j, i_j)\|^2 \rho^{\frac{k-j}{2}} \right) \\ & = \frac{n-1}{n} \sum_{j=0}^k \left(\sum_{i=1}^n p_i \mathbb{E} \|\nabla f_i(\hat{x}_j^i)\|^2 \right) (k-j) \rho^{\frac{k-j}{2}}. \end{aligned} \tag{17}$$

It follows from (17) and (16) that

$$\begin{aligned} A_2 & \leq \frac{n-1}{n} \mathbb{E} \sum_{j=0}^k \left(\sum_{i=1}^n p_i \|\nabla f_i(\hat{x}_j^i)\|^2 \right) \left(\rho^{k-j} + 2(k-j) \rho^{\frac{k-j}{2}} \right) \\ & \stackrel{\text{Lemma 5}}{\leq} \frac{n-1}{n} \sum_{j=0}^k \left(12L^2 \mathbb{E} \hat{M}_j + 6\zeta^2 + 2 \mathbb{E} \left\| \sum_{j'=1}^n p_{j'} \nabla f_{j'}(\hat{x}_j^{j'}) \right\|^2 \right) \left(\rho^{k-j} + 2(k-j) \rho^{\frac{k-j}{2}} \right) \\ & \leq \frac{n-1}{n} \sum_{j=0}^k \left(12L^2 \mathbb{E} \hat{M}_j + 2 \mathbb{E} \left\| \sum_{j'=1}^n p_{j'} \nabla f_{j'}(\hat{x}_j^{j'}) \right\|^2 \right) \left(\rho^{k-j} + 2(k-j) \rho^{\frac{k-j}{2}} \right) \\ & \quad + 6\zeta^2 \underbrace{\frac{n-1}{n} \left(\frac{1}{1-\rho} + \frac{2\sqrt{\rho}}{(1-\sqrt{\rho})^2} \right)}_{=\bar{\rho}}. \end{aligned} \tag{18}$$

Finally from (15), (18) and (14) we obtain

$$\begin{aligned} & \mathbb{E} \left\| \frac{X_{k+1} \mathbf{1}_n}{n} - X_{k+1} e_i \right\|^2 \\ & \leq 2\gamma^2 A_1 + 2M^2 \gamma^2 A_2 \\ & \leq 2\gamma^2 M \sigma^2 \bar{\rho} \\ & \quad + 2\gamma^2 M^2 \mathbb{E} \sum_{j=0}^k \left(12L^2 \hat{M}_j + 2 \mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \left(\rho^{k-j} + 2(k-j) \rho^{\frac{k-j}{2}} \right) + 12\gamma^2 M^2 \zeta^2 \bar{\rho} \\ & = 2\gamma^2 (M \sigma^2 + 6M^2 \zeta^2) \bar{\rho} \\ & \quad + 2 \frac{n-1}{n} M^2 \gamma^2 \mathbb{E} \sum_{j=0}^k \left(12L^2 \hat{M}_j + 2 \mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \left(\rho^{k-j} + 2(k-j) \rho^{\frac{k-j}{2}} \right). \end{aligned}$$

This completes the proof. □

Lemma 7. While $C_1 > 0$, we have

$$\frac{\sum_{k=0}^{K-1} \mathbb{E} \hat{M}_k}{K} \leq \frac{2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} + \frac{4\gamma^2 M^2}{K} \left(T \frac{n-1}{n} + \bar{\rho}\right) \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_k^i) \right\|^2}{C_1}, \quad \forall K \geq 1.$$

Proof. From Lemma 6 and noting that $\hat{X}_k = X_{k-\tau_k}$, we have

$$\begin{aligned} & \mathbb{E} \left\| \frac{\hat{X}_k \mathbf{1}_n}{n} - \hat{X}_k e_i \right\|^2 \\ & \leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\ & \quad + 2\frac{n-1}{n} M^2 \gamma^2 \sum_{j=0}^{k-\tau_k-1} \left(12L^2 \mathbb{E} \hat{M}_j + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \left(\rho^{k-\tau_k-1-j} + 2(k-\tau_k-1-j)\rho^{\frac{k-\tau_k-1-j}{2}} \right). \end{aligned}$$

By averaging from $i = 1$ to n with distribution \mathcal{I} we obtain

$$\begin{aligned} \mathbb{E} \hat{M}_k &= \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\hat{X}_k \mathbf{1}_n}{n} - \hat{X}_k e_i \right\|^2 \\ & \leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\ & \quad + 2M^2 \gamma^2 \frac{n-1}{n} \sum_{j=0}^{k-\tau_k-1} \left(12L^2 \mathbb{E} \hat{M}_j + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \left(\rho^{k-\tau_k-1-j} + 2(k-\tau_k-1-j)\rho^{\frac{k-\tau_k-1-j}{2}} \right). \end{aligned}$$

It follows that

$$\begin{aligned} \frac{\sum_{k=0}^{K-1} \mathbb{E} \hat{M}_k}{K} & \leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\ & \quad + \frac{2\gamma^2}{K} \frac{n-1}{n} M^2 \sum_{k=0}^{K-1} \sum_{j=0}^{k-\tau_k-1} \left(12L^2 \mathbb{E} \hat{M}_j + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \times \\ & \quad \left(\rho^{k-\tau_k-1-j} + 2(k-\tau_k-1-j)\rho^{\frac{k-\tau_k-1-j}{2}} \right) \\ & = 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\ & \quad + \frac{2\gamma^2}{K} \frac{n-1}{n} M^2 \sum_{k=0}^{K-1} \sum_{j=0}^{k-\tau_k-1} \left(12L^2 \mathbb{E} \hat{M}_j + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \times \\ & \quad \left(\rho^{\max\{k-\tau_k-1-j, 0\}} + 2(\max\{k-\tau_k-1-j, 0\})\rho^{\frac{\max\{k-\tau_k-1-j, 0\}}{2}} \right) \\ & \leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\ & \quad + \frac{2\gamma^2}{K} \frac{n-1}{n} M^2 \sum_{j=0}^{K-1} \sum_{k=j+1}^{\infty} \left(12L^2 \mathbb{E} \hat{M}_j + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \times \\ & \quad \left(\rho^{\max\{k-\tau_k-1-j, 0\}} + 2(\max\{k-\tau_k-1-j, 0\})\rho^{\frac{\max\{k-\tau_k-1-j, 0\}}{2}} \right) \\ & \leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\ & \quad + \frac{2\gamma^2}{K} \frac{n-1}{n} M^2 \sum_{j=0}^{K-1} \left(12L^2 \mathbb{E} \hat{M}_j + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \nabla f_i(\hat{x}_j^i) \right\|^2 \right) \left(T + \sum_{h=0}^{\infty} (\rho^h + 2h\rho^{\frac{h}{2}}) \right) \end{aligned}$$

$$\begin{aligned}
&\leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\
&\quad + \frac{2\gamma^2}{K}M^2\left(T\frac{n-1}{n} + \bar{\rho}\right)\sum_{j=0}^{K-1}\left(12L^2\mathbb{E}\hat{M}_j + 2\mathbb{E}\left\|\sum_{i=1}^n p_i\nabla f_i(\hat{x}_{j,i})\right\|^2\right) \\
&\leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} \\
&\quad + \frac{4\gamma^2M^2}{K}\left(T\frac{n-1}{n} + \bar{\rho}\right)\sum_{k=0}^{K-1}\mathbb{E}\left\|\sum_{i=1}^n p_i\nabla f_i(\hat{x}_k^i)\right\|^2 \\
&\quad + \frac{24L^2\gamma^2M^2}{K}\left(T\frac{n-1}{n} + \bar{\rho}\right)\sum_{k=0}^{K-1}\mathbb{E}\hat{M}_k.
\end{aligned}$$

By rearranging the terms we obtain

$$\begin{aligned}
&\underbrace{\left(1 - 24L^2M^2\gamma^2\left(T\frac{n-1}{n} + \bar{\rho}\right)\right)}_{C_1}\frac{\sum_{k=0}^{K-1}\mathbb{E}\hat{M}_k}{K} \\
&\leq 2\gamma^2(M\sigma^2 + 6M^2\zeta^2)\bar{\rho} + \frac{4\gamma^2M^2}{K}\left(T\frac{n-1}{n} + \bar{\rho}\right)\sum_{k=0}^{K-1}\mathbb{E}\left\|\sum_{i=1}^n p_i\nabla f_i(\hat{x}_k^i)\right\|^2,
\end{aligned}$$

we complete the proof. \square

Lemma 8. For all $k \geq 0$ we have

$$\mathbb{E}\left\|\frac{X_k\mathbf{1}_n - \hat{X}_k\mathbf{1}_n}{n}\right\|^2 \leq \frac{\tau_k^2\gamma^2\sigma^2M}{n^2} + \tau_k\gamma^2\sum_{t=1}^{\tau_k}\left(\frac{M^2}{n^2}\sum_{i=1}^n p_i\mathbb{E}\|\nabla f_i(\hat{x}_{k-t}^i)\|^2\right).$$

Proof.

$$\begin{aligned}
\mathbb{E}\left\|\frac{X_k\mathbf{1}_n - \hat{X}_k\mathbf{1}_n}{n}\right\|^2 &\stackrel{\text{Assumption 1-7}}{=} \mathbb{E}\left\|\frac{\sum_{t=1}^{\tau_k}\gamma\partial g(\hat{X}_{k-t}; \zeta_{k-t}^{i_{k-t}}, i_{k-t})\mathbf{1}_n}{n}\right\|^2 \\
&\leq \tau_k\sum_{t=1}^{\tau_k}\gamma^2\mathbb{E}\left\|\frac{\partial g(\hat{X}_{k-t}; \zeta_{k-t}^{i_{k-t}}, i_{k-t})\mathbf{1}_n}{n}\right\|^2 \\
&\stackrel{\text{Lemma 4}}{\leq} \tau_k\sum_{t=1}^{\tau_k}\gamma^2\left(\frac{\sigma^2M}{n^2} + \frac{M^2}{n^2}\sum_{i=1}^n p_i\mathbb{E}\|\nabla f_i(\hat{x}_{k-t}^i)\|^2\right),
\end{aligned}$$

where the first step comes from any $n \times n$ doubly stochastic matrix multiplied by $\mathbf{1}_n$ equals $\mathbf{1}_n$ and Assumption 1-7. \square

Proof to Corollary 2. To prove this result, we will apply Theorem 1. We first verify that all conditions can be satisfied in Theorem 1.

First $C_1 > 0$ can be satisfied by a stronger condition $C_1 \geq 1/2$ which can be satisfied by $\gamma \leq \frac{1}{4\sqrt{6}ML}\left(T\frac{n-1}{n} + \bar{\rho}\right)^{-1/2}$. Second $C_3 \leq 1$ can be satisfied by :

$$\gamma \leq \min\left\{\frac{n}{8MT^2L}, \frac{1}{8\sqrt{3}LM}\bar{\rho}^{-1/2}, \frac{1}{32nML}\bar{\rho}^{-1}, \frac{n^{1/3}}{8\sqrt{6}MLT^{2/3}}\bar{\rho}^{-1/3}\right\}$$

and $C_1 \geq 1/2$, which can be seen from

$$C_3 = \frac{1}{2} + \frac{2 \left(6\gamma^2 L^2 M^2 + \gamma n M L + \frac{12M^3 L^3 T^2 \gamma^3}{n} \right) \bar{\rho}}{C_1} + \frac{L T^2 \gamma M}{n}$$

$$\stackrel{C_1 \geq \frac{1}{2}}{\leq} \frac{1}{2} + 24\gamma^2 L^2 M^2 + 4\gamma n M L + \frac{48M^3 L^3 T^2 \gamma^3}{n} \bar{\rho} + \frac{L T^2 \gamma M}{n}.$$

The requirements on γ are given by making each of the last four terms smaller than $1/8$:

$$\frac{L T^2 \gamma M}{n} \leq \frac{1}{8} \iff \gamma \leq \frac{n}{8 M T^2 L},$$

$$24\gamma^2 L^2 M^2 \bar{\rho} \leq \frac{1}{8} \iff \gamma \leq \frac{1}{8\sqrt{3} L M} \bar{\rho}^{-1/2},$$

$$4\gamma n M L \bar{\rho} \leq \frac{1}{8} \iff \gamma \leq \frac{1}{32 n M L} \bar{\rho}^{-1},$$

and

$$\frac{48M^3 L^3 T^2 \gamma^3}{n} \bar{\rho} \leq \frac{1}{8} \iff \gamma \leq \frac{n^{1/3}}{8\sqrt{6} M L T^{2/3}} \bar{\rho}^{-1/3}.$$

Third $C_2 \geq 0$ can be satisfied by

$$\gamma \leq \min \left\{ \frac{n}{10 L M}, \frac{n}{2\sqrt{5} M L T}, \frac{n^{1/3}}{8 M L} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/3}, \frac{1}{4\sqrt{5} L M} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/2}, \frac{n^{1/2}}{6 M L T^{1/2}} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/4} \right\}$$

and $C_1 \geq 1/2$, which can be seen from

$$C_2 := \frac{\gamma M}{2n} - \frac{\gamma^2 L M^2}{n^2} - \frac{2M^3 L^2 T^2 \gamma^3}{n^3} - \left(\frac{6\gamma^2 L^3 M^2}{n^2} + \frac{\gamma M}{n} L^2 + \frac{12M^3 L^4 T^2 \gamma^3}{n^3} \right) \frac{4M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho})}{C_1} \geq 0$$

$$\stackrel{C_1 \geq \frac{1}{2}}{\iff} 1 \geq \frac{2\gamma L M}{n} + \frac{4M^2 L^2 T^2 \gamma^2}{n^2} + \frac{96\gamma L^3 M}{n} + 16L^2 + \frac{192M^2 L^4 T^2 \gamma^2}{n^2} M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho}).$$

The last inequality is satisfied given the requirements on γ because each term on the RHS is bounded by $1/5$:

$$\frac{2\gamma L M}{n} \leq \frac{1}{5} \iff \gamma \leq \frac{n}{10 L M},$$

$$\frac{4M^2 L^2 T^2 \gamma^2}{n^2} \leq \frac{1}{5} \iff \gamma \leq \frac{n}{2\sqrt{5} M L T},$$

$$\frac{96\gamma L^3 M}{n} M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho}) \leq \frac{1}{5} \iff \gamma \leq \frac{n^{1/3}}{8 M L} (T \frac{n-1}{n} + \bar{\rho})^{-1/3},$$

$$16L^2 M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho}) \leq \frac{1}{5} \iff \gamma \leq \frac{1}{4\sqrt{5} L M} (T \frac{n-1}{n} + \bar{\rho})^{-1/2},$$

$$\frac{192M^2 L^4 T^2 \gamma^2}{n^2} M^2 \gamma^2 (T \frac{n-1}{n} + \bar{\rho}) \leq \frac{1}{5} \iff \gamma \leq \frac{n^{1/2}}{6 M L T^{1/2}} (T \frac{n-1}{n} + \bar{\rho})^{-1/4}.$$

Combining all above the requirements on γ to satisfy $C_1 \geq 1/2, C_2 \geq 0$ and $C_3 \leq 1$ are

$$\gamma \leq \frac{1}{M L} \min \left\{ \begin{array}{l} \frac{1}{4\sqrt{6}} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/2}, \frac{n}{8 T^2}, \frac{1}{8\sqrt{3}} \bar{\rho}^{-1/2}, \\ \frac{1}{32 n} \bar{\rho}^{-1}, \frac{n^{1/3}}{8\sqrt{6} T^{2/3}} \bar{\rho}^{-1/3}, \\ \frac{n}{10}, \frac{n}{2\sqrt{5} T}, \frac{n^{1/3}}{8} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/3}, \\ \frac{1}{4\sqrt{5}} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/2}, \frac{n^{1/2}}{6 T^{1/2}} \left(T \frac{n-1}{n} + \bar{\rho} \right)^{-1/4} \end{array} \right\}.$$

Note that the RHS is larger than

$$U := \frac{1}{ML} \min \left\{ \frac{1}{8\sqrt{3}\sqrt{T\frac{n-1}{n}+\bar{\rho}}}, \frac{n}{8T^2}, \frac{1}{32n\bar{\rho}}, \frac{n}{10}, \frac{n^{1/12}(n-1)^{-1/4}}{(8\sqrt{6}T^{2/3}+8)(T+\bar{\rho}\frac{n}{n-1})^{1/3}} \right\}.$$

Let $\gamma = \frac{n}{10ML + \sqrt{\sigma^2 + 6M\zeta^2}\sqrt{KM}}$ then if $\gamma \leq U$ we will have $C_1 \geq 1/2, C_2 \geq 0$ and $C_3 \leq 1$. Further investigation gives us

$$\begin{aligned} \gamma = \frac{n}{10ML + \sqrt{\sigma^2 + 6M\zeta^2}\sqrt{KM}} &\leq \frac{1}{ML} \min \left\{ \frac{1}{8\sqrt{3}\sqrt{T\frac{n-1}{n}+\bar{\rho}}}, \frac{n}{8T^2}, \frac{1}{32n\bar{\rho}}, \frac{n}{10}, \frac{n^{1/12}(n-1)^{-1/4}}{(8\sqrt{6}T^{2/3}+8)(T+\bar{\rho}\frac{n}{n-1})^{1/3}} \right\} \\ \Leftrightarrow 10ML + \sqrt{\sigma^2 + 6M\zeta^2}\sqrt{KM} &\geq nML \max \left\{ \frac{8\sqrt{3}\sqrt{T\frac{n-1}{n} + \bar{\rho}}, \frac{8T^2}{n}, 32n\bar{\rho},}{(8\sqrt{6}T^{2/3}+8)(T+\bar{\rho}\frac{n}{n-1})^{1/3}}, \frac{n^{1/12}(n-1)^{-1/4}}{n^{1/12}(n-1)^{-1/4}} \right\} \\ \Leftrightarrow K &\geq \frac{ML^2n^2}{\sigma^2 + 6M\zeta^2} \max \left\{ \frac{192 \left(T\frac{n-1}{n} + \bar{\rho} \right), \frac{64T^4}{n^2}, 1024n^2\bar{\rho}^2,}{(8\sqrt{6}T^{2/3}+8)^2(T+\bar{\rho}\frac{n}{n-1})^{2/3}}, \frac{n^{1/6}(n-1)^{-1/2}}{n^{1/6}(n-1)^{-1/2}} \right\}. \end{aligned}$$

It follows from Theorem 1 that if the last inequality is satisfied and $\gamma = \frac{n}{10ML + \sqrt{\sigma^2 + 6M\zeta^2}\sqrt{KM}}$, we have

$$\begin{aligned} \frac{\sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f \left(\frac{X_k \mathbf{1}_n}{n} \right) \right\|^2}{K} &\leq \frac{2(\mathbb{E}f(x_0) - f^*)n}{\gamma KM} + \frac{2\gamma L}{Mn} (M\sigma^2 + 6M^2\zeta^2) \\ &= \frac{20(\mathbb{E}f(x_0) - f^*)L}{K} + \frac{2(\mathbb{E}f(x_0) - f^*)\sqrt{\sigma^2 + 6M\zeta^2}}{\sqrt{KM}} \\ &\quad + \frac{2L}{M \left(10ML + \sqrt{\sigma^2 + 6M\zeta^2}\sqrt{KM} \right)} (M\sigma^2 + 6M^2\zeta^2) \\ &\leq \frac{20(\mathbb{E}f(x_0) - f^*)L}{K} + \frac{2(\mathbb{E}f(x_0) - f^* + L)\sqrt{\sigma^2 + 6M\zeta^2}}{\sqrt{KM}}. \end{aligned}$$

This completes the proof. □